

# How Do I Revise My Agent's Action Theory?

Ivan José Varzinczak

Meraka Institute, CSIR  
Pretoria, South Africa  
ivan.varzinczak@meraka.org.za

## Abstract

Logical theories in reasoning about actions may also evolve, and knowledge engineers need revision tools to incorporate new incoming laws about the dynamic environment. We here fill this gap by providing an algorithmic approach for action theory revision. We give a well defined semantics that ensures minimal change, and show correctness of our algorithms w.r.t. the semantic constructions.

## Introduction

Like any logical theory, action theories in reasoning about actions may evolve, and thus need revision methods to adequately accommodate new information about the behavior of actions. In (Eiter et al. 2005; Herzig, Perrussel, and Varzinczak 2006; Varzinczak 2008) update and contraction-based methods for action theory repair are defined. Here we continue this important though quite new thread of investigation and develop a minimal change approach for *revising* a domain description.

The motivation is as follows. Consider an agent designed to interact with a coffee machine. Among her beliefs, the agent may know that a coffee is a hot drink, that after buying she gets a coffee, and that with a token it is possible to buy. We can see the agent's beliefs about the behavior of actions in this scenario as a transition system (Figure 1).

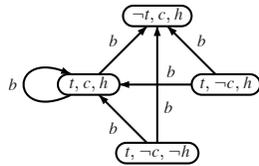


Figure 1: A transition system depicting the agent's knowledge about the dynamics of the coffee machine.  $b$ ,  $t$ ,  $c$ , and  $h$  stand for, respectively, *buy*, *token*, *coffee*, and *hot*.

Well, at some stage the agent may learn that coffee is the only hot drink available at the machine, or that even without a token she can still buy, or that all possible executions of *buy* should lead to states where  $\neg token$  is the case. These are examples of *revision* with new laws about the dynamics of the environment under consideration. And here we are interested in exactly these kinds of theory modification.

The contributions of the present work are as follows:

- What is the semantics of revising an action theory by a law? How to get minimal change, i.e., how to keep as much knowledge about other laws as possible?
- How to syntactically revise an action theory so that its result corresponds to the intended semantics?

Here we answer these questions.

## Logical Preliminaries

Our base formalism is multimodal logic  $K_n$  (Popkorn 1994).

### Action Theories in Multimodal K

Let  $\mathfrak{A} = \{a_1, a_2, \dots\}$  be the set of *atomic actions* of a domain. To each action  $a$  there is associated a modal operator  $[a]$ .  $\mathfrak{P} = \{p_1, p_2, \dots\}$  denotes the set of *propositions*, or *atoms*.  $\mathfrak{L} = \{p, \neg p : p \in \mathfrak{P}\}$  is the set of *literals*.  $\ell$  denotes a literal and  $|\ell|$  the atom in  $\ell$ .

We use  $\varphi, \psi, \dots$  to denote *Boolean formulas*.  $\mathfrak{F}$  is the set of all Boolean formulas. A propositional valuation  $v$  is a *maximally consistent* set of literals. We denote by  $v \Vdash \varphi$  the fact that  $v$  satisfies  $\varphi$ . By  $val(\varphi)$  we denote the set of all valuations satisfying  $\varphi$ .  $\models_{\text{CPL}}$  is the classical consequence relation.  $Cn(\varphi)$  denotes all logical consequences of  $\varphi$ .

With  $IP(\varphi)$  we denote the set of *prime implicants* (Quine 1952) of  $\varphi$ . By  $\pi$  we denote a prime implicant, and  $atm(\pi)$  is the set of atoms occurring in  $\pi$ . Given  $\ell$  and  $\pi$ ,  $\ell \in \pi$  abbreviates ' $\ell$  is a literal of  $\pi$ '.

We use  $\Phi, \Psi, \dots$  to denote complex formulas (possibly with modal operators).  $\langle a \rangle$  is the dual operator of  $[a]$  ( $\langle a \rangle \Phi =_{\text{def}} \neg[a]\neg\Phi$ ).

A  $K_n$ -*model* is a tuple  $\mathcal{M} = \langle W, R \rangle$  where  $W$  is a set of valuations, and  $R$  maps action constants  $a$  to accessibility relations  $R_a \subseteq W \times W$ . Given  $\mathcal{M}$ ,  $\models_w^{\mathcal{M}} p$  ( $p$  is true at world  $w$  of model  $\mathcal{M}$ ) if  $w \Vdash p$ ;  $\models_w^{\mathcal{M}} [a]\Phi$  if  $\models_{w'}^{\mathcal{M}} \Phi$  for every  $w'$  s.t.  $(w, w') \in R_a$ ; truth conditions for the other connectives are as usual. By  $\mathcal{M}$  we will denote a set of  $K_n$ -models.

$\mathcal{M}$  is a model of  $\Phi$  (noted  $\models^{\mathcal{M}} \Phi$ ) if and only if  $\models_w^{\mathcal{M}} \Phi$  for all  $w \in W$ .  $\mathcal{M}$  is a model of a set of formulas  $\Sigma$  (noted  $\models^{\mathcal{M}} \Sigma$ ) if and only if  $\models_w^{\mathcal{M}} \Phi$  for every  $\Phi \in \Sigma$ .  $\Phi$  is a *consequence* of

the global axioms  $\Sigma$  in all  $K_n$ -models (noted  $\Sigma \models_{K_n} \Phi$ ) if and only if for every  $\mathcal{M}$ , if  $\models^{\mathcal{M}} \Sigma$ , then  $\models^{\mathcal{M}} \Phi$ .

In  $K_n$  we can state laws describing the behavior of actions. Here we distinguish three types of them.

**Static Laws** A *static law* is a formula  $\varphi \in \mathfrak{F}$  that characterizes the possible states of the world. An example is  $\text{coffee} \rightarrow \text{hot}$ : if the agent holds a coffee, then she holds a hot drink. The set of static laws of a domain is denoted by  $\mathcal{S}$ .

**Effect Laws** An *effect law* for  $a$  has the form  $\varphi \rightarrow [a]\psi$ , with  $\varphi, \psi \in \mathfrak{F}$ . Effect laws relate an action to its effects, which can be conditional. The consequent  $\psi$  is the effect that always obtains when  $a$  is executed in a state where the antecedent  $\varphi$  holds. An example is  $\text{token} \rightarrow [\text{buy}]\text{hot}$ : whenever the agent has a token, after buying, she has a hot drink. If  $\psi$  is inconsistent we have a special kind of effect law that we call an *inexecutability law*. For example,  $\neg \text{token} \rightarrow [\text{buy}]\perp$  says that *buy* cannot be executed if the agent has no token. The set of effect laws is denoted by  $\mathcal{E}$ .

**Executability Laws** An *executability law* for  $a$  has the form  $\varphi \rightarrow \langle a \rangle \top$ , with  $\varphi \in \mathfrak{F}$ . It stipulates the context in which  $a$  is guaranteed to be executable. (In  $K_n$   $\langle a \rangle \top$  reads “ $a$ ’s execution is possible”.) For instance,  $\text{token} \rightarrow \langle \text{buy} \rangle \top$  says that buying can be executed whenever the agent has a token. The set of executability laws of a domain is denoted by  $\mathcal{X}$ .

Given  $a$ ,  $\mathcal{E}_a$  (resp.  $\mathcal{X}_a$ ) will denote the set of only those effect (resp. executability) laws about  $a$ .

**Action Theories**  $\mathcal{T} = \mathcal{S} \cup \mathcal{E} \cup \mathcal{X}$  is an *action theory*.

To make the presentation more clear to the reader, we here assume that the agent’s theory contains all frame axioms. However, all we shall say here can be defined within a formalism with a solution to the frame and ramification problems like (Herzig, Perrussel, and Varzinczak 2006) do. The action theory of our example will thus be:

$$\mathcal{T} = \left\{ \begin{array}{l} \text{coffee} \rightarrow \text{hot}, \text{token} \rightarrow \langle \text{buy} \rangle \top, \\ \neg \text{coffee} \rightarrow [\text{buy}]\text{coffee}, \neg \text{token} \rightarrow [\text{buy}]\perp, \\ \text{coffee} \rightarrow [\text{buy}]\text{coffee}, \text{hot} \rightarrow [\text{buy}]\text{hot} \end{array} \right\}$$

Figure 1 above shows a  $K_n$ -model for the theory  $\mathcal{T}$ .

Sometimes it will be useful to consider models whose possible worlds are *all* the possible states allowed by  $\mathcal{S}$ :

**Definition 1**  $\mathcal{M} = \langle W, R \rangle$  is a *big frame* of  $\mathcal{T}$  if and only if:

- $W = \text{val}(\mathcal{S})$ ; and
- $R_a = \{(w, w') : \forall \varphi \rightarrow [a]\psi \in \mathcal{E}_a, \text{ if } \models_w^{\mathcal{M}} \varphi \text{ then } \models_{w'}^{\mathcal{M}} \psi\}$

Big frames of  $\mathcal{T}$  are not always models of  $\mathcal{T}$ .

**Definition 2**  $\mathcal{M}$  is a *supra-model* of  $\mathcal{T}$  iff  $\models^{\mathcal{M}} \mathcal{T}$  and  $\mathcal{M}$  is a *big frame* of  $\mathcal{T}$ .

Figure 2 depicts a supra-model of our example  $\mathcal{T}$ .

### Prime Valuations

An atom  $p$  is *essential* to  $\varphi$  if and only if  $p \in \text{atm}(\varphi')$  for all  $\varphi'$  such that  $\models_{\text{CPL}} \varphi \leftrightarrow \varphi'$ . For instance,  $p_1$  is essential to  $\neg p_1 \wedge (\neg p_1 \vee p_2)$ .  $\text{atm}!(\varphi)$  will denote the essential atoms of  $\varphi$ . (If  $\varphi$  is a tautology or a contradiction, then  $\text{atm}!(\varphi) = \emptyset$ .)

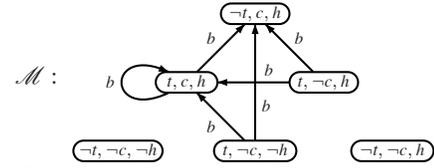


Figure 2: Supra-model for the coffee machine scenario.

For  $\varphi \in \mathfrak{F}$ ,  $\varphi^*$  is the set of all  $\varphi' \in \mathfrak{F}$  such that  $\varphi \models_{\text{CPL}} \varphi'$  and  $\text{atm}(\varphi') \subseteq \text{atm}!(\varphi)$ . For instance,  $p_1 \vee p_2 \notin p_1^*$ , as  $p_1 \models_{\text{CPL}} p_1 \vee p_2$  but  $\text{atm}(p_1 \vee p_2) \not\subseteq \text{atm}!(p_1)$ . Clearly,  $\text{atm}(\bigwedge \varphi^*) = \text{atm}!(\bigwedge \varphi^*)$ . Moreover, whenever  $\models_{\text{CPL}} \varphi \leftrightarrow \varphi'$ , then  $\text{atm}!(\varphi) = \text{atm}!(\varphi')$  and also  $\varphi^* = \varphi'^*$ .

**Theorem 1 (Parikh 1999)**  $\models_{\text{CPL}} \varphi \leftrightarrow \bigwedge \varphi^*$ , and  $\text{atm}(\varphi^*) \subseteq \text{atm}(\varphi')$  for every  $\varphi'$  s.t.  $\models_{\text{CPL}} \varphi \leftrightarrow \varphi'$ .

Thus for every  $\varphi \in \mathfrak{F}$  there is a unique least set of elementary atoms such that  $\varphi$  may equivalently be expressed using only atoms from that set. Hence,  $\text{Cn}(\varphi) = \text{Cn}(\varphi^*)$ .

Given a valuation  $v$ ,  $v' \subseteq v$  is a *subvaluation*. For  $W$  a set of valuations, a subvaluation  $v'$  *satisfies*  $\varphi \in \mathfrak{F}$  modulo  $W$  (noted  $v' \models_W \varphi$ ) if and only if  $v \models \varphi$  for all  $v \in W$  such that  $v' \subseteq v$ . A subvaluation  $v$  *essentially satisfies*  $\varphi$  modulo  $W$  ( $v \models_W^! \varphi$ ) if and only if  $v \models_W \varphi$  and  $\{\ell : \ell \in v\} \subseteq \text{atm}!(\varphi)$ .

**Definition 3** Let  $\varphi \in \mathfrak{F}$  and  $W$  be a set of valuations. A *subvaluation*  $v$  is a *prime subvaluation* of  $\varphi$  (modulo  $W$ ) if and only if  $v \models_W^! \varphi$  and there is no  $v' \subseteq v$  s.t.  $v' \models_W^! \varphi$ .

A prime subvaluation of a formula  $\varphi$  is one of the weakest states of truth in which  $\varphi$  is true. (Notice the similarity with the syntactical notion of prime implicant (Quine 1952).) We denote all prime subvaluations of  $\varphi$  modulo  $W$  by  $\text{base}(\varphi, W)$ .

**Theorem 2** Let  $\varphi \in \mathfrak{F}$  and  $W$  be a set of valuations. Then for all  $w \in W$ ,  $w \models \varphi$  if and only if  $w \models \bigvee_{v \in \text{base}(\varphi, W)} \bigwedge_{\ell \in v} \ell$ .

### Closeness Between Models

When revising a model, we perform a change in its structure. Because there can be several ways of modifying a model (not all minimal), we need a notion of distance between models to identify those closest to the original one.

As we are going to see in more depth in the sequel, changing a model amounts to modifying its possible worlds or its accessibility relation. Hence, the distance between two  $K_n$ -models will depend upon the distance between their sets of worlds and accessibility relations. These here will be based on the *symmetric difference* between sets, defined as  $X \dot{-} Y = (X \setminus Y) \cup (Y \setminus X)$ .

**Definition 4** Let  $\mathcal{M} = \langle W, R \rangle$ .  $\mathcal{M}' = \langle W', R' \rangle$  is at least as close to  $\mathcal{M}$  as  $\mathcal{M}'' = \langle W'', R'' \rangle$ , noted  $\mathcal{M}' \preceq_{\mathcal{M}} \mathcal{M}''$ , iff

- either  $W \dot{-} W' \subseteq W \dot{-} W''$
- or  $W \dot{-} W' = W \dot{-} W''$  and  $R \dot{-} R' \subseteq R \dot{-} R''$

This is an extension of Burger and Heidema’s relation (Burger and Heidema 2002) to our modal case. Note that other distance notions are also possible, like e.g. the *cardinality* of symmetric differences or Hamming distance.

## Semantics of Revision

Contrary to contraction, where we want the negation of a law to be *satisfiable*, in revision we want a new law to be *valid*. Thus we must eliminate all cases satisfying its negation.

The idea in our semantics is as follows: we initially have a set of models  $\mathcal{M}$  in which a given formula  $\Phi$  is (potentially) not valid, i.e.,  $\Phi$  is (possibly) not true in every model in  $\mathcal{M}$ . In the result we want to have only models of  $\Phi$ . Adding  $\Phi$ -models to  $\mathcal{M}$  is of no help. Moreover, adding models makes us lose laws: the resulting theory would be more liberal.

One solution amounts to deleting from  $\mathcal{M}$  those models that are not  $\Phi$ -models. Of course removing only some of them does not solve the problem, we must delete every such a model. By doing that, all resulting models will be models of  $\Phi$ . (This corresponds to *theory expansion*, when the resulting theory is satisfiable.) However, if  $\mathcal{M}$  contains no model of  $\Phi$ , we will end up with  $\emptyset$ . Consequence: the resulting theory is inconsistent. (This is the main revision problem.) In this case the solution is to *substitute* each model  $\mathcal{M}$  in  $\mathcal{M}$  by its *nearest modifications*  $\mathcal{M}_\Phi^*$  that makes  $\Phi$  true. This lets us to keep as close as possible to the original models that we had.

Before defining revision of sets of models, we present what modifications of (individual) models are.

### Revising a Model by a Static Law

Suppose that our coffee deliverer agent discovers that the only hot drink that is served on the machine is coffee. In this case, we might want to revise her beliefs with the new static law  $coffee \leftrightarrow hot$ .

Considering the model in Figure 2, we see that  $\neg coffee \wedge hot$  is satisfiable. As we do not want this, the first step is to *remove* all worlds in which  $\neg coffee \wedge hot$  is true. The second step is to guarantee all the remaining worlds satisfy the new law. This issue has been largely addressed in the literature on belief revision and update (Gärdenfors 1988; Winslett 1988; Katsuno and Mendelzon 1992; Herzig and Rifi 1999). Here we can achieve that with a semantics similar to that of classical revision operators: basically one can change the set of possible valuations, by removing or adding worlds.

In our example, removing the possible worlds  $\{t, \neg c, h\}$  and  $\{\neg t, \neg c, h\}$  would do the job (there is no need to add new valuations since the new static law is satisfied in at least one world of the original model).

The delicate point in removing worlds is that it may result in the loss of some executability laws: in the example, if there were only one arrow leaving some world  $w$  and pointing to  $\{\neg t, \neg c, h\}$ , then removing the latter from the model would make the action under concern no longer executable in  $w$ . Here we claim that this is intuitive: if the state of the world to which we could move is no longer possible, then we do not have a transition to that state anymore. Hence, if that transition was the only one we had, it is natural to lose it.

One could also ask what to do with the accessibility relation if new worlds must be added (revision case). We claim that it is reckless to blindly add new elements to  $R$ . Instead, we shall postpone correction of executability laws, if needed. This approach is debatable, but with the information we have at hand, it is the safest way of changing static laws.

**Definition 5** Let  $\mathcal{M} = \langle W, R \rangle$ .  $\mathcal{M}' = \langle W', R' \rangle \in \mathcal{M}_\varphi^*$  iff  $W' = (W \setminus val(\neg\varphi)) \cup val(\varphi)$  and  $R' \subseteq R$ .

Clearly  $\models^{\mathcal{M}'} \varphi$  for all  $\mathcal{M}' \in \mathcal{M}_\varphi^*$ . The minimal models of the revision of  $\mathcal{M}$  by  $\varphi$  are those closest to  $\mathcal{M}$  w.r.t.  $\preceq_{\mathcal{M}}$ :

**Definition 6**  $rev(\mathcal{M}, \varphi) = \bigcup \min\{\mathcal{M}_\varphi^*, \preceq_{\mathcal{M}}\}$ .

In the example of Figure 2,  $rev(\mathcal{M}, coffee \leftrightarrow hot)$  is the singleton  $\{\mathcal{M}'\}$ , with  $\mathcal{M}'$  as shown in Figure 3.

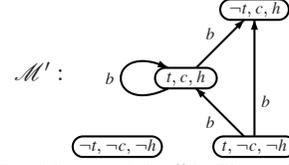


Figure 3: Revising model  $\mathcal{M}$  in Figure 2 with  $coffee \leftrightarrow hot$ .

### Revising a Model by an Effect Law

Let's suppose now that our agent eventually discovers that after buying coffee she does not keep her token. This means that her theory should now be revised by the new effect law  $token \rightarrow [buy]\neg token$ . Looking at model  $\mathcal{M}$  in Figure 2, this amounts to guaranteeing that  $token \wedge \langle buy \rangle token$  is satisfiable in none of its worlds. To do that, we have to look at all the worlds satisfying this formula (if any) and

- either make  $token$  false in each of these worlds,
- or make  $\langle buy \rangle token$  false in all of them.

If we chose the first option, we will essentially flip the truth value of literal  $token$  in the respective worlds, which changes the set of valuations of the model. If we chose the latter, we will basically remove  $buy$ -arrows leading to  $token$ -worlds, which amounts to changing the accessibility relation.

In our example, worlds  $w_1 = \{token, coffee, hot\}$ ,  $w_2 = \{token, \neg coffee, hot\}$  and  $w_3 = \{token, \neg coffee, \neg hot\}$  satisfy the formula  $token \wedge \langle buy \rangle token$ . Flipping  $token$  in all of them to  $\neg token$  would do the job, but this would also have as consequence the introduction of a new static law:  $\neg token$  would now be valid, i.e., the agent never has a token! Do we want this?

We claim that changing action laws should not have as side effect a change in the static laws. These have a special status (Shanahan 1997), and should change only if required. Hence each world satisfying  $token \wedge \langle buy \rangle token$  has to be changed so that  $\langle buy \rangle token$  becomes untrue in it. In the example, we thus should remove  $(w_1, w_1)$ ,  $(w_2, w_1)$  and  $(w_3, w_1)$  from  $R$ .

**Definition 7** Let  $\mathcal{M} = \langle W, R \rangle$ .  $\mathcal{M}' = \langle W', R' \rangle \in \mathcal{M}_{\varphi \rightarrow [a]\psi}^*$  iff:

- $W' = W, R' \subseteq R, \models^{\mathcal{M}'} \varphi \rightarrow [a]\psi$ , and
- If  $(w, w') \in R \setminus R'$ , then  $\not\models_w^{\mathcal{M}} \varphi$

The minimal models resulting from revision of a model  $\mathcal{M}$  by a new effect law are those closest to  $\mathcal{M}$  w.r.t.  $\preceq_{\mathcal{M}}$ :

**Definition 8**  $rev(\mathcal{M}, \varphi \rightarrow [a]\psi) = \bigcup \min\{\mathcal{M}_{\varphi \rightarrow [a]\psi}^*, \preceq_{\mathcal{M}}\}$ .

Taking  $\mathcal{M}$  as in Figure 2,  $rev(\mathcal{M}, token \rightarrow [buy]\neg token)$  will be the singleton  $\{\mathcal{M}'\}$  depicted in Figure 4.

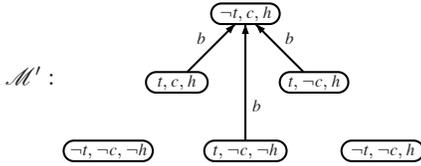


Figure 4: Revising  $\mathcal{M}$  in Figure 2 with  $token \rightarrow [buy]\neg token$ .

### Revising a Model by an Executability Law

Let us now suppose that at some stage it has been decided to grant free coffee to everybody. Faced with this information, we have to revise the agent's laws to reflect the fact that *buy* can also be executed in  $\neg token$ -contexts:  $\neg token \rightarrow \langle buy \rangle \top$  is a new executability law.

Considering model  $\mathcal{M}$  in Figure 2, we observe that  $\neg token \wedge [buy]\perp$  is satisfiable. Hence we must throw  $\neg token \wedge [buy]\perp$  away to ensure the new law becomes true.

To remove  $\neg token \wedge [buy]\perp$  we have to look at all worlds satisfying it and modify  $\mathcal{M}$  so that they no longer satisfy that formula. Given worlds  $w_4 = \{\neg token, \neg coffee, \neg hot\}$  and  $w_5 = \{\neg token, \neg coffee, hot\}$ , we have two options: change the interpretation of *token* in both or add new arrows leaving these worlds. A question that arises is 'what choice is more drastic: change a world or an arrow?'. Again, here we claim that changing the world's content (the valuation) is more drastic, as the existence of such a world is foreseen by some static law and is hence assumed to be as it is, unless we have enough information supporting the contrary, in which case we explicitly change the static laws (see above). Thus we shall add a new *buy*-arrow from each of  $w_4$  and  $w_5$ .

Having agreed on that, the issue now is: which worlds should the new arrows point to? In order to comply with minimal change, the new arrows shall point to worlds that are relevant targets of each of the  $\neg token$ -worlds in question.

**Definition 9** Let  $\mathcal{M} = \langle W, R \rangle$ ,  $w, w' \in W$ , and  $\mathcal{M}$  be a set of models s.t.  $\mathcal{M} \in \mathcal{M}$ . Then  $w'$  is a relevant target world of  $w$  w.r.t.  $\varphi \rightarrow \langle a \rangle \top$  for  $\mathcal{M}$  in  $\mathcal{M}$  iff  $\models_w^{\mathcal{M}} \varphi$  and

- If there is  $\mathcal{M}' = \langle W', R' \rangle \in \mathcal{M}$  such that  $R'_a(w) \neq \emptyset$ :
  - for all  $\ell \in w' \setminus w$ , there is  $\psi' \in \mathfrak{F}$  s.t. there is  $v' \in base(\psi', W)$  s.t.  $v' \subseteq w'$ ,  $\ell \in v'$ , and  $\models_w^{\mathcal{M}_i} [a]\psi'$  for every  $\mathcal{M}_i \in \mathcal{M}$
  - for all  $\ell \in w \cap w'$ , either there is  $\psi' \in \mathfrak{F}$  s.t. there is  $v' \in base(\psi', W)$  s.t.  $v' \subseteq w'$ ,  $\ell \in v'$ , and  $\models_w^{\mathcal{M}_i} [a]\psi'$  for all  $\mathcal{M}_i \in \mathcal{M}$ ; or there is  $\mathcal{M}_i \in \mathcal{M}$  s.t.  $\not\models_w^{\mathcal{M}_i} [a]\neg \ell$
- If  $R'_a(w) = \emptyset$  for every  $\mathcal{M}' = \langle W', R' \rangle \in \mathcal{M}$ :
  - for all  $\ell \in w' \setminus w$ , there is  $\mathcal{M}_i = \langle W_i, R_i \rangle \in \mathcal{M}$  s.t. there is  $u, v \in W_i$  s.t.  $(u, v) \in R_{i,a}$  and  $\ell \in v \setminus u$
  - for all  $\ell \in w \cap w'$ , there is  $\mathcal{M}_i = \langle W_i, R_i \rangle \in \mathcal{M}$  s.t. there is  $u, v \in W_i$  s.t.  $(u, v) \in R_{i,a}$  and  $\ell \in u \cap v$ , or for all  $\mathcal{M}_i = \langle W_i, R_i \rangle \in \mathcal{M}$ , if  $(u, v) \in R_{i,a}$ , then  $\neg \ell \notin v \setminus u$

By  $rt(w, \varphi \rightarrow \langle a \rangle \top, \mathcal{M}, \mathcal{M})$  we denote the set of all relevant target worlds of  $w$  w.r.t.  $\varphi \rightarrow \langle a \rangle \top$  for  $\mathcal{M}$  in  $\mathcal{M}$ .

In our example,  $w_6 = \{\neg token, coffee, hot\}$  is the only relevant target world here: the two other  $\neg token$ -worlds violate the effect *coffee* of *buy*, while the three *token*-worlds would make us violate the frame axiom  $\neg token \rightarrow [buy]\neg token$ .

**Definition 10** Let  $\mathcal{M} = \langle W, R \rangle$ .  $\mathcal{M}' = \langle W', R' \rangle \in \mathcal{M}_{\varphi \rightarrow \langle a \rangle \top}^*$  iff:

- $W' = W, R \subseteq R', \models^{\mathcal{M}'} \varphi \rightarrow \langle a \rangle \top$ , and
- If  $(w, w') \in R' \setminus R$ , then  $w' \in rt(w, \varphi \rightarrow \langle a \rangle \top, \mathcal{M}, \mathcal{M})$

The minimal models resulting from revising a model  $\mathcal{M}$  by a new executability law are those closest to  $\mathcal{M}$  w.r.t.  $\preceq_{\mathcal{M}}$ :

**Definition 11**  $rev(\mathcal{M}, \varphi \rightarrow \langle a \rangle \top) = \bigcup \min\{\mathcal{M}_{\varphi \rightarrow \langle a \rangle \top}^*, \preceq_{\mathcal{M}}\}$ .

In our running example,  $rev(\mathcal{M}, \neg token \rightarrow \langle buy \rangle \top)$  is the singleton  $\{\mathcal{M}'\}$ , where  $\mathcal{M}'$  is as shown in Figure 5.

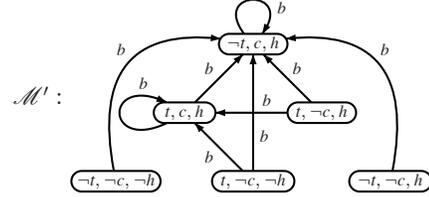


Figure 5: The result of revising model  $\mathcal{M}$  in Figure 2 by the new executability law  $\neg token \rightarrow \langle buy \rangle \top$ .

### Revising Sets of Models

Up until now we have seen what the revision of single models means. Now we are ready for a unified definition of revision of a set of models  $\mathcal{M}$  by a new law  $\Phi$ :

**Definition 12** Let  $\mathcal{M}$  be a set of models and  $\Phi$  a law. Then

$$\mathcal{M}_{\Phi}^* = (\mathcal{M} \setminus \{\mathcal{M} : \not\models^{\mathcal{M}} \Phi\}) \cup \bigcup_{\mathcal{M} \in \mathcal{M}} rev(\mathcal{M}, \Phi)$$

Definition 12 comprises both *expansion* and *revision*: in the former, addition of the new law gives a satisfiable theory; in the latter a deeper change is required to get rid of the inconsistency.

### Syntactic Operators for Revision

We now turn our attention to the syntactical counterpart of revision. Our endeavor here is to perform minimal change also at the syntactical level. By  $\mathcal{T}_{\Phi}^*$  we denote the result of revising an action theory  $\mathcal{T}$  with a new law  $\Phi$ .

### Revising a Theory by a Static Law

Looking at the semantics of revision by Boolean formulas, we see that revising an action theory by a new static law may conflict with the executability laws: some of them may be lost and thus have to be changed as well. The approach here is to preserve as many executability laws as we can in the old possible states. To do that, we look at each possible valuation that is common to the new  $\mathcal{S}$  and the old one. Every time an executability used to hold in that state and no inexecutability holds there now, we make the action executable in such a context. For those contexts not allowed by

the old  $\mathcal{S}$ , we make  $a$  inexecutable (cf. the semantics). Algorithm 1 deals with that ( $\mathcal{S} \star \varphi$  denotes the classical revision of  $\mathcal{S}$  by  $\varphi$  built upon some well established method from the literature (Winslett 1988; Katsuno and Mendelzon 1992; Herzig and Rifi 1999)).

---

**Algorithm 1** Revision by a Static Law
 

---

**input:**  $T, \varphi$   
**output:**  $T_\varphi^*$   
 $S' := \mathcal{S} \star \varphi, \mathcal{E}' := \mathcal{E}, \mathcal{X}' := \emptyset$   
**for all**  $\pi \in IP(S')$  **do**  
   **for all**  $A \subseteq atm(\pi)$  **do**  
      $\varphi_A := \bigwedge_{\substack{p_i \in atm(\pi) \\ p_i \in A}} p_i \wedge \bigwedge_{\substack{p_i \in atm(\pi) \\ p_i \notin A}} \neg p_i$   
     **if**  $S' \not\vdash_{\mathcal{CPL}} (\pi \wedge \varphi_A) \rightarrow \perp$  **then**  
       **if**  $S \not\vdash_{\mathcal{CPL}} (\pi \wedge \varphi_A) \rightarrow \perp$  **then**  
         **if**  $T \not\vdash_{\mathcal{R}_n} (\pi \wedge \varphi_A) \rightarrow \langle a \rangle \top$  **and**  $S', \mathcal{E}', \mathcal{X}' \not\vdash_{\mathcal{R}_n} \neg(\pi \wedge \varphi_A)$  **then**  
            $\mathcal{X}'_a := \{(\varphi_i \wedge \pi \wedge \varphi_A) \rightarrow \langle a \rangle \top : \varphi_i \rightarrow \langle a \rangle \top \in \mathcal{X}'_a\}$   
         **else**  
            $\mathcal{E}' := \mathcal{E}' \cup \{(\pi \wedge \varphi_A) \rightarrow [a] \perp\}$   
        $T_\varphi^* := S' \cup \mathcal{E}' \cup \mathcal{X}'$

---

**Revising a Theory by an Effect Law**

When revising a theory by a new effect law  $\varphi \rightarrow [a]\psi$ , we want to eliminate all possible executions of  $a$  leading to  $\neg\psi$ -states. To achieve that, we look at all  $\varphi$ -contexts and every time a transition to some  $\neg\psi$ -context is not always the case, i.e.,  $T \not\vdash_{\mathcal{R}_n} \varphi \rightarrow \langle a \rangle \neg\psi$ , we can safely force  $[a]\psi$  for that context. On the other hand, if in such a context there is always an execution of  $a$  to  $\neg\psi$ , then we should strengthen the executability laws to make room for the new effect in that context we want to add. Algorithm 2 below does the job.

---

**Algorithm 2** Revision by an Effect Law
 

---

**input:**  $T, \varphi \rightarrow [a]\psi$   
**output:**  $T_{\varphi \rightarrow [a]\psi}^*$   
 $T' := T$   
**for all**  $\pi \in IP(S \wedge \varphi)$  **do**  
   **for all**  $A \subseteq atm(\pi)$  **do**  
      $\varphi_A := \bigwedge_{\substack{p_i \in atm(\pi) \\ p_i \in A}} p_i \wedge \bigwedge_{\substack{p_i \in atm(\pi) \\ p_i \notin A}} \neg p_i$   
     **if**  $S \not\vdash_{\mathcal{CPL}} (\pi \wedge \varphi_A) \rightarrow \perp$  **then**  
       **for all**  $\pi' \in IP(S \wedge \neg\psi)$  **do**  
         **if**  $T' \not\vdash_{\mathcal{R}_n} (\pi \wedge \varphi_A) \rightarrow \langle a \rangle \pi'$  **then**  
            $T' := (T' \setminus \mathcal{X}'_a) \cup \{(\varphi_i \wedge \neg(\pi \wedge \varphi_A)) \rightarrow \langle a \rangle \top : \varphi_i \rightarrow \langle a \rangle \top \in \mathcal{X}'_a\}$   
        $T' := T' \cup \{(\pi \wedge \varphi_A) \rightarrow [a]\psi\}$   
       **if**  $T' \not\vdash_{\mathcal{R}_n} (\pi \wedge \varphi_A) \rightarrow [a] \perp$  **then**  
          $T' := T' \cup \{(\varphi_i \wedge \pi \wedge \varphi_A) \rightarrow \langle a \rangle \top : \varphi_i \rightarrow \langle a \rangle \top \in T\}$   
        $T_{\varphi \rightarrow [a]\psi}^* := T'$

---

**Revising a Theory by an Executability Law**

Revision of a theory by a new executability law has as consequence a change in the effect laws: all those laws preventing

the execution of  $a$  shall be weakened. Moreover, to comply with minimal change, we must ensure that in all models of the resulting theory there will be at most *one* transition by  $a$  from those worlds in which  $T$  precluded  $a$ 's execution.

Let  $(\mathcal{E}_a^{\varphi, \perp})_1, \dots, (\mathcal{E}_a^{\varphi, \perp})_n$  denote minimum subsets (w.r.t. set inclusion) of  $\mathcal{E}_a$  such that  $\mathcal{S}, (\mathcal{E}_a^{\varphi, \perp})_i \not\vdash_{\mathcal{R}_n} \varphi \rightarrow [a] \perp$ . (According to (Herzig and Varzinczak 2007), one can ensure at least one such a set always exists.) Let  $\mathcal{E}_a^- = \bigcup_{1 \leq i \leq n} (\mathcal{E}_a^{\varphi, \perp})_i$ . The effect laws in  $\mathcal{E}_a^-$  will serve as guidelines to get rid of  $[a] \perp$  in each  $\varphi$ -world allowed by  $T$ : they are the laws to be weakened to allow for  $\langle a \rangle \top$  in  $\varphi$ -contexts.

Our algorithm works as follows. To force  $\varphi \rightarrow \langle a \rangle \top$  to be true in all models of the resulting theory, we visit every possible  $\varphi$ -context allowed by it and make the following operations to ensure  $\langle a \rangle \top$  is the case for that context: Given a  $\varphi$ -context, if  $T$  does not always preclude  $a$  from being executed in it, we can safely force  $\langle a \rangle \top$  without modifying other laws. On the other hand, if  $a$  is always inexecutable in that context, then we should weaken the laws in  $\mathcal{E}_a^-$ . The first thing we must do is to preserve all old effects in all other  $\varphi$ -worlds. To achieve that we specialize the above laws to each possible valuation (maximal conjunction of literals) satisfying  $\varphi$  but the actual one. Then, in the current  $\varphi$ -valuation, we must ensure that action  $a$  may have any effect, i.e., from this  $\varphi$ -world we can reach any other possible world. We achieve that by weakening the *consequent* of the laws in  $\mathcal{E}_a^-$  to the exclusive disjunction of all possible contexts in  $T$ . Finally, to get minimal change, we must ensure that all literals in this  $\varphi$ -valuation that are not forced to change are preserved. We do this by stating a conditional frame axiom of the form  $(\varphi_k \wedge \ell) \rightarrow [a]\ell$ , where  $\varphi_k$  is the above-mentioned  $\varphi$ -valuation.

Algorithm 3 gives the pseudo-code for that.

**Correctness of the Algorithms**

Suppose we have two atoms  $p_1$  and  $p_2$ , and one action  $a$ . Let  $\mathcal{T}_1 = \{\neg p_2, p_1 \rightarrow [a]p_2, \langle a \rangle \top\}$ . The only model of  $\mathcal{T}_1$  is  $\mathcal{M}$  in Figure 6. Revising such a model by  $p_1 \vee p_2$  gives us the models  $\mathcal{M}'_i, 1 \leq i \leq 3$ , in Figure 6. Now, revising  $\mathcal{T}_1$  by  $p_1 \vee p_2$  will give us  $\mathcal{T}_{1, p_1 \vee p_2}^* = \{p_1 \wedge \neg p_2, p_1 \rightarrow [a]p_2\}$ . The only model of  $\mathcal{T}_{1, p_1 \vee p_2}^*$  is  $\mathcal{M}'_1$  in Figure 6. This means that the semantic revision may produce models (viz.  $\mathcal{M}'_2$  and  $\mathcal{M}'_3$  in Figure 6) that are not models of the revised theories.

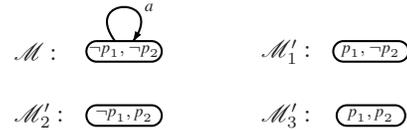


Figure 6: Model  $\mathcal{M}$  of  $\mathcal{T}_1$  and revision of  $\mathcal{M}$  by  $p_1 \vee p_2$ .

The other way round the algorithms may give theories whose models do not result from revision of models of the initial theory: let  $\mathcal{T}_2 = \{(p_1 \vee p_2) \rightarrow [a] \perp, \langle a \rangle \top\}$ . Its only model is  $\mathcal{M}$  (Figure 6). Revising  $\mathcal{M}$  by  $p_1 \vee p_2$  is as above. However  $\mathcal{T}_{2, p_1 \vee p_2}^* = \{p_1 \vee p_2, (p_1 \vee p_2) \rightarrow [a] \perp\}$  has a model  $\mathcal{M}'' = \{\{p_1, p_2\}, \{p_1, \neg p_2\}, \{\neg p_1, p_2\}\}, \emptyset\}$  that is not in  $\mathcal{M}_{p_1 \vee p_2}^*$ .

---

**Algorithm 3** Revision by an executability law

---

**input:**  $\mathcal{T}, \varphi \rightarrow \langle a \rangle \top$ **output:**  $\mathcal{T}_{\varphi \rightarrow \langle a \rangle \top}^*$  $\mathcal{T}' := \mathcal{T}$ **for all**  $\pi \in IP(\mathcal{S} \wedge \varphi)$  **do****for all**  $A \subseteq atm(\pi)$  **do** $\varphi_A := \bigwedge_{\substack{p_i \in atm(\pi) \\ p_i \in A}} p_i \wedge \bigwedge_{\substack{p_i \in atm(\pi) \\ p_i \notin A}} \neg p_i$ **if**  $\mathcal{S} \not\models_{\text{CPL}} (\pi \wedge \varphi_A) \rightarrow \perp$  **then****if**  $\mathcal{T}' \models_{\mathbb{K}_n} (\pi \wedge \varphi_A) \rightarrow [a] \perp$  **then** $(\mathcal{T}' \setminus \mathcal{E}'_a) \cup \{(\varphi_i \wedge \neg(\pi \wedge \varphi_A)) \rightarrow [a]\psi_i :$  $\varphi_i \rightarrow [a]\psi_i \in \mathcal{E}'_a\}$   $\cup$  $\mathcal{T}' := \{(\varphi_i \wedge \pi \wedge \varphi_A) \rightarrow [a] \bigoplus_{\substack{\pi' \in IP(\mathcal{S}) \\ A' \subseteq atm(\pi')}} (\pi' \wedge \varphi_{A'}) :$  $\varphi_i \rightarrow [a]\psi_i \in \mathcal{E}'_a\}$ **for all**  $L \subseteq \mathcal{L}$  **do****if**  $\mathcal{S} \models_{\text{CPL}} (\pi \wedge \varphi_A) \rightarrow \bigwedge_{\ell \in L} \ell$  **then****for all**  $\ell \in L$  **do****if**  $\mathcal{T}' \models_{\mathbb{K}_n} \ell \rightarrow [a] \perp$  **or**  $(\mathcal{T}' \not\models_{\mathbb{K}_n} \ell \rightarrow [a] \neg \ell$  **and** $\mathcal{T}' \models_{\mathbb{K}_n} \ell \rightarrow [a] \ell)$  **then** $\mathcal{T}' := \mathcal{T}' \cup \{(\pi \wedge \varphi_A \wedge \ell) \rightarrow [a] \ell\}$  $\mathcal{T}' := \mathcal{T}' \cup \{(\pi \wedge \varphi_A) \rightarrow \langle a \rangle \top\}$  $\mathcal{T}_{\varphi \rightarrow \langle a \rangle \top}^* := \mathcal{T}'$ 

---

All this happens because the possible states are not completely characterized by the static laws. Fortunately, concentrating on supra-models of  $\mathcal{T}$ , we get the right result.

**Theorem 3** If  $\mathcal{M} = \{\mathcal{M} : \mathcal{M} \text{ is a supra-model of } \mathcal{T}\}$  and there is  $\mathcal{M}' \in \mathcal{M}$  s.t.  $\models^{\mathcal{M}'} \Phi$ , then  $\bigcup_{\mathcal{M} \in \mathcal{M}} rev(\mathcal{M}, \Phi) \subseteq \mathcal{M}$ .

Then, revision of models of  $\mathcal{T}$  by a law  $\Phi$  in the semantics produces models of the output of the algorithms  $\mathcal{T}_{\Phi}^*$ :

**Theorem 4** If  $\mathcal{M} = \{\mathcal{M} : \mathcal{M} \text{ is a supra-model of } \mathcal{T}\} \neq \emptyset$ , then for every  $\mathcal{M}' \in \mathcal{M}_{\Phi}^*$ ,  $\models^{\mathcal{M}'} \mathcal{T}_{\Phi}^*$ .

Also, models of  $\mathcal{T}_{\Phi}^*$  result from revision of models of  $\mathcal{T}$  by  $\Phi$ :

**Theorem 5** If  $\mathcal{M} = \{\mathcal{M} : \mathcal{M} \text{ is a supra-model of } \mathcal{T}\} \neq \emptyset$ , then for every  $\mathcal{M}'$ , if  $\models^{\mathcal{M}'} \mathcal{T}_{\Phi}^*$ , then  $\mathcal{M}' \in \mathcal{M}_{\Phi}^*$ .

Sticking to supra-models of  $\mathcal{T}$  is not a big deal. We can use the algorithms in (Herzig and Varzinczak 2007) to ensure  $\mathcal{T}$  is characterized by its supra-models and that  $\mathcal{M} \neq \emptyset$ .

## Conclusion and Perspectives

The problem of action theory change has only recently received attention in the literature, both in action languages (Baral and Lobo 1997; Eiter et al. 2005) and in modal logic (Herzig, Perrussel, and Varzinczak 2006; Varzinczak 2008).

Here we have studied what revising action theories by a law means, both in the semantics and at the syntactical (algorithmic) level. We have defined a semantics based on distances between models that also captures minimal change w.r.t. the preservation of effects of actions. With our algorithms and the correctness results we have established the

link between the semantics and the syntax for theories with supra-models. (Due to page limits, proofs are omitted here.)

Our next step on the subject is analyze the behavior of our operators w.r.t. AGM-like postulates (Alchourrón, Gärdenfors, and Makinson 1985) for modal theories and the relationship between our revision method and contraction. What is known is that Levi identity (Levi 1977),  $\mathcal{T}_{\Phi}^* = \mathcal{T}_{\neg\Phi} \cup \{\Phi\}$ , in general does not hold for action laws  $\Phi$ . The reason is that up to now there is no contraction operator for  $\neg\Phi$  where  $\Phi$  is an action law. Indeed this is the general contraction problem for action theories: contraction of a theory  $\mathcal{T}$  by a general formula (like  $\neg\Phi$  above) is still an open problem in the area. The definition of a general method will certainly mostly benefit from the semantic modifications we studied here (addition/removal of arrows and worlds).

Given the relationship between modal logics and description logics, a revision method for DL TBoxes would also benefit from the constructions we defined here.

## References

- Alchourrón, C.; Gärdenfors, P.; and Makinson, D. 1985. On the logic of theory change: Partial meet contraction and revision functions. *J. of Symbolic Logic* 50:510–530.
- Baral, C., and Lobo, J. 1997. Defeasible specifications in action theories. In *Proc. IJCAI*, 1441–1446.
- Burger, I., and Heidema, J. 2002. Merging inference and conjecture by information. *Synthese* 131(2):223–258.
- Eiter, T.; Erdem, E.; Fink, M.; and Senko, J. 2005. Updating action domain descriptions. In *Proc. IJCAI*, 418–423.
- Gärdenfors, P. 1988. *Knowledge in Flux: Modeling the Dynamics of Epistemic States*. MIT Press.
- Herzig, A., and Rifi, O. 1999. Propositional belief base update and minimal change. *Artificial Intelligence* 115(1):107–138.
- Herzig, A., and Varzinczak, I. 2007. Metatheory of actions: beyond consistency. *Artificial Intelligence* 171:951–984.
- Herzig, A.; Perrussel, L.; and Varzinczak, I. 2006. Elaborating domain descriptions. In *Proc. ECAI*, 397–401.
- Katsuno, H., and Mendelzon, A. 1992. On the difference between updating a knowledge base and revising it. In *Belief revision*. Cambridge. 183–203.
- Levi, I. 1977. Subjunctives, dispositions and chances. *Synthese* 34:423–455.
- Parikh, R. 1999. Beliefs, belief revision, and splitting languages. In *Logic, Language and Computation*, 266–278.
- Popkorn, S. 1994. *First Steps in Modal Logic*. Cambridge University Press.
- Quine, W. V. O. 1952. The problem of simplifying truth functions. *American Mathematical Monthly* 59:521–531.
- Shanahan, M. 1997. *Solving the frame problem*. Cambridge, MA: MIT Press.
- Varzinczak, I. 2008. Action theory contraction and minimal change. In *Proc. KR*, 651–661.
- Winslett, M.-A. 1988. Reasoning about action using a possible models approach. In *Proc. AAAI*, 89–93.