# Defeasible Specifications in Action Theories

Chitta Baral*
Computer Science Department
University of Texas at El Paso
El Paso, Texas 79968
U.S.A
*chitta@cs. utep. edu*

Jorge Lobo[+]
Department of EECS
Univ. of Illinois at Chicago
851 S. Morgan St
Chicago, IL 60607, USA

## Abstract

Recent research in formalizing effects of actions on a world in the presence of constraints have mostly concentrated on non-defeasible specifications, where the effect of actions and constraints are strictly specified. In this paper we show how to incorporate defeasibility into the specifications. In our approach we consider extensions of the high level language A of Gelfond and Lifschitz and introduce defeasible constraints and effect propositions. While direct semantics of A does not need a logical language, our semantics is defined using extended logic programming. This is due to the defeasibility of the specification.

## 1   Introduction and Motivation

Recent research in formalizing effects of actions on a world in the presence of constraints have mostly concentrated on non-defeasible specifications, where the effect of actions and constraints are strictly specified[1]. The focus of these research was more on solving the frame problem and other problems such as the qualification and the ramification problem, and it made sense to start with non-defeasible specifications.

In this paper we show how to incorporate defeasibility in the specification without loosing any of the good properties of previous specification languages. It is important to consider defeasible specifications for the same reason

[1] Among the various approaches, the strictness of the high level specification language of A [5] and its successors [6], and the language used by Sandewal [10] is clear. Although, Reiter [9] and his group do not use a particular specification language per se, the syntactic restriction they put on how the effect axioms and constraints can be initially specified in first-order logic, and the dependence of the compilation process that follows on the input syntax amounts to a specification language.

jorge@eecs.uic.eduas the necessity of non-rnonotonic theories in knowled representation and common-sense reasoning. The specifications are not definite and un-modifiable and may evolve. Our language should allow easy evolution or in the words of McCarthy, they should be 'elaboration tolerant[1]. We take a first step in this direction by allowing defeasible specifications.

In our approach we consider[2] extensions of the high level language A of Gelfond and Lifschitz [5] and introduce defeasible constraints and effect propositions. Direct semantics of A and many of its early successors (particularly, the formalization of transition functions that map state-action pairs to states) do not use a logical language. But our formalization of the transition function is done using extended logic programming. This is due to the defeasibility of the specification. (Earlier, Baral [1] and McCain and Turner [8] use logic programming and inference rules respectively, to formalize transition functions in presence of causal constraints.) We then show how with minor modifications to the logic programming specification plus the addition of a few rules we can reason about the effect of a sequence of actions.

The necessity of defeasible specifications may be illustrated by the following example from [12]. Suppose we have a lamp with a switch on the side. *Turning the switch* causes an internal circuit in the lamp to flip between close *(not open)* and *open.* We can say then that the effect of executing the action *turning.switch* causes the circuit to close if it was *open;* or it will cause the switch to *open* if it was *not open.* It is also the case that when the circuit is close *(not open)* the *light* in the lamp is *on,* and when the circuit is *open* the light is *off.* Notice that although the state of the light is determined by the state of the circuit, the state of the light is indirectly affected by the action of turning the switch. This indirect effect is known in action theories as a *ramification* of the action. In an action description language such as AC [11] this action domain may be described by the

[2] We consider this language due to its simplicity. Our ideas are important for other specification languages and their formalizations. We would consider the impact of defeasible specifications on them in the full version of this paper.

following propositions:[3]

$$turning\_switch \text{ causes } open \text{ if } \neg open \quad (i)$$
$$turning\_switch \text{ causes } \neg open \text{ if } open \quad (2)$$
$$open \text{ suffices for } \neg light\_on \quad (3)$$
$$\neg open \text{ suffices for } light\_on \quad (4)$$

Ramification effects are usually derived from the state of fluents,[4] while direct effects come from the execution of an action. An equivalent formulation can be written using Lin's formalization [7] or the state specifications of Baral [1j.

Suppose now that we can also *unplug* the lamp and that a consequence of this action is that there will be no *power* going through the lamp. We can add to our description the following proposition to cover this action:

$$unplug\_lamp \text{ causes } \neg power \quad (5)$$

If there is no power going through the lamp the light in the lamp can never be on (this is regardless of the state of the circuit). We then must add the following ramification proposition:

$$\neg power \text{ suffices for } \neg light\_on \quad (6)$$

Notice too, this proposition together with the previous propositions may create contradictions or forbid certain sequences of actions. For example, executing the sequence *[unplug Jamp, turning switch]* creates problems because there will be a state where the circuit is not *open* that causes the light to be on, and at the same time the light can not be on since the lamp is unplugged and there is no power going through the lamp. The problem is that a closed circuit not always, but *normally,* causes the light in lamp to be on.[5] Thus, the causality in proposition (4) must be *defeasible.*

In the following section we present *ADC,* an action description language with defeasible causality propositions. To make it easier for the reader we have excluded many of the other features of action description languages such as constraints, that are described in the literature. Most of these can be easily added to *ADC. Following our example, we start by describing how we can specify defeasible ramification effects. Later in Section 5 we show how our framework can also be used to capture defeasibtlity of the direct effect of actions.*

---

[3] At this point we appeal to the intuition of the reader with the meaning of the language. In the following sections we will present a formal description of an action language similar to AC.

[4]Symbols that are used to represent the state of the world are called fluents. In this example the fluents are *open* and *light_on.*

[5]There are many other reasons why a closed circuit may not cause the light to be on. It could be that the bulb in the lamp is burned, or there was a storm that blacked out the entire neighborhood.

## 2  The language *ATX*

The language *ADC* has two disjoint non-empty sets of symbols called *fluents* and *actions,* and four kinds of *propositions: initial* propositions, *effect* propositions, *sufficiency* propositions and *defeasible sufficiency* propositions. An initial proposition is an expression of the form

$$\text{initially } l \quad (7)$$

where *l* is a fluent literal. A *fluent literal is* either a fluent or a fluent preceded by ¬. A fluent literal is *negative* if it is preceded by ¬, otherwise it is called positive. An effect proposition is an expression of the form

$$a \text{ causes } f \text{ if } p_1, \ldots, p_n \quad (8)$$

where *a* is an action, and *l* and and each of $p_1, \ldots, p_n$ $(0 \leq n)$ is a fluent literal. If n = 0 we will write:

$$a \text{ causes } l \quad (9)$$

Two effect propositions with preconditions $p_1, \ldots, p_n$ and $q_1, \ldots, q_m$ respectively are said to be *contradictory* if they describe the effect of the same action *a* on complementary fs, and $\{p_1, \ldots, p_n\} \cap \{\overline{q_1}, \ldots, \overline{q_m}\} = \emptyset$

Non-defeasible and defeasible sufficiency propositions are expressions of the form:

$$p_1, \ldots, p_n \text{ suffices for } l \quad (10)$$
$$p_1, \ldots, p_n \text{ normally suffices for } l \quad (11)$$

where *l* and each $p_1, \ldots, p_n$ $(1 \leq n)$ is a fluent literal. A *domain description* is a collection of initial, effect, sufficiency and defeasible sufficiency propositions. The initial propositions are referred to as *facts* and the other propositions are referred to as *causal laws.*

## 3  Transition functions of *ADC*

We take a slightly different approach in defining the semantics of the language *ADC.* Most of the action description languages [11; 6; 2] that are inspired by *A* have an independent characterization without involving any of the standard logical formalisms, such as logic programming, default logic, circumscription, classical logic, etc. But when we allow defeasible causality and keep open the possibility of hierarchies of such causality it is not clear if we can avoid using standard logical formalisms and have independent characterizations.

An inherent component of semantics of action description languages is the transition function *Res* from actions and states to power set of states, where a state is a set of fluents. Intuitively, if $s' \in Res(a, s)$ we say that execution of action a in a state s may take us to the state s'. When we have only deterministic actions we can treat *Res* as a function from actions and states to states.

In this section we define the transition function of *ADC* through a translation of the domain description into an extended logic program. (We use the answer set semantics of extended logic programs [4].) Given a set of causal

laws *D,* an action a, and a state .s, we construct an extended logic program $\Pi(D, a, s)$, and use its answer sets to define *Res(a,s)*. The extended logic program $\Pi(D, a, s)$ uses variables of three sorts: *situation* variables S, $S_1$,..., *fluent* variables F, $F_1$,..., and *action* variables $A, A_1$, —

For convenience, we use the following notations. The atom *Holds(f,s)* denotes the literal *holds(f,s)* if *f* is a positive fluent literal and the literal $\neg holds(|f|, s)$ if *f* is a negative fluent literal. Also, for any fluent *f*, $|f| = f$ and $|\neg f| = f$. By $\overline{Holds(f,s)}$ we denote the literal which is complementary to *Holds(f,s)*. In addition, we will have an special postfix function symbol ' that we will apply to fluents and fluent variables. It is used to distinguish negative from positive fluent literals occurring as predicate arguments.[6]

We now describe the construction of the program $\Pi(D, a, s)$, which consists of the translations of the individual propositions in *D* and certain other rules.

Translation 3.1 *[Constructing $\Pi(D, a, s)$]*

*1. Initial Database*
If *f* is true in *s* then $\Pi(D, a, s)$ contains

(1.1) $holds(f, s) \leftarrow$

Otherwise, $\Pi(D, a, s)$ contains

(1.2) $\neg holds(f, s) \leftarrow$

*2. Inertia Rule*
$\Pi(D, a, s)$ contains the following inertia rules:

(2.1) $holds(F, res(a, s)) \leftarrow holds(F, s),$
$\qquad not\ ab_0(F, a, s), not\ ab_1(F, a, s)$

(2.2) $\neg holds(F, res(a, s)) \leftarrow \neg holds(F, s),$
$\qquad not\ ab_0(F', a, s), not\ ab_1(F', a, s)$

This rule is motivated by the minimality consideration which states that only changes that happens to a state as a result of an action *a* are the ones dictated by the direct and indirect effects of the action. If a has a direct effect on a fluent F (resp. the negated fluent F'), we say that that F (resp. F') is abnormal with respect to *a* in *s*. Thus, we expect $ab_0$(F, a, *s*) (resp. $ab_0$(F', a, .s)) to be true (see item 3 below). Similarly, if If a has an indirect effect on a fluent F (resp. the negated fluent F'), we say that that F (resp. F') is abnormal with respect to *a* in *s*. Thus, we expect $ab_1$(F, a,s) (resp. $ab_1$(F',a,*)) to be true (see item 4 below).

*3. Translating Effect propositions*

For each effect proposition of the form (8) in *D* which describes the effect of the action a, $\Pi(D, a, s)$ contains the following rules:

$Holds(f, res(a, s)) \leftarrow Holds(p_1, s), \ldots, Holds(p_n, s)$

---
[6]We could have used $\neg$ instead of ' and write $holds(\neg(f), s)$, but the unary negation $\neg$ is already used as predicate negation in extended logic programs.

$ab_0(\overline{f}, a, s) \leftarrow Holds(p_1, s), \ldots, Holds(p_n, s)$

where $\overline{f} = l$ if $f = \neg l$; otherwise $\overline{f} = f'$.

*4. Translating non-defeasible sufficiency propositions*

For each non-defeasible sufficiency proposition of the form (10) in *D*, $\Pi(D, a, s)$ contains the following rules:

$Holds(f, res(a, s)) \leftarrow Holds(p_1, res(a, s)), \ldots,$
$\qquad\qquad Holds(p_n, res(a, s))$

$ab_0(\overline{f}, a, s) \leftarrow Holds(p_1, res(a, s)), \ldots,$
$\qquad\qquad Holds(p_n, res(a, s))$

*5. Translating defeasible sufficiency propositions*

For each defeasible sufficiency proposition of the form (11) in *D*, $\Pi(D, a, s)$ contains the following rules:

$Holds(f, res(a, s)) \leftarrow Holds(p_1, res(a, s)), \ldots,$
$\qquad\qquad Holds(p_n, res(a, s)), not\ ab_0(f, a, s)$

$ab_1(\overline{f}, a, s) \leftarrow Holds(p_1, res(a, s)), \ldots,$
$\qquad\qquad Holds(p_n, res(a, s)), not\ ab_0(\overline{f}, a, s)$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

**Definition 3.1** For a state *s*, and an action *a*, $Res(a, s)$ is defined with respect to the causal laws *D* through the program $\Pi(D, a, s)$ in the following way:

Let **A** be the set of answer sets of $\Pi(D, a, s)$.

If $\Pi(D, a, s)$ is inconsistent (i.e., only the set Lit is in **A**) then $Res(a, s)$ is undefined, meaning that the action *a* is not executable in the state *s*.

Otherwise,

$Res(a, s) = \{\{f : holds(f, res(a, s)) \in A\} \cup \{\neg f : \neg holds(f, res(a, s)) \in A\} | A \in \mathbf{A}\}$ $\quad\square$

Note that $Res(a, s)$ is defined if contradictory effect rules are absent and sufficiency propositions are signed[7]. In addition stratification of the sufficiency propositions guarantee $Res(a, s)$ to be a singleton set.

### 3.1 Examples

**Example 3.1** Let $D_1$ consist of the the causal laws (1), (2), (3), and (4). Let $s_1 = \{open, power\}$, and $s_2 = \{light\_on, power\}$. Intuitively, the action *turning_switch* should cause the transition from $s_1$ to $s_2$ and from $s_2$ to $s_1$. We now verify that this indeed is the case, by examining the answer sets of the programs $\Pi(D_1, turning\_switch, s_1)$ and $\Pi(D_1, turning\_switch, s_2)$. For $s_1$ the initial database S is:

$\{holds(open, s_1), holds(power, s_1), \neg holds(light\_on, s_1)\}$

The following four rules are the effect propositions:

---
[7]We use the term 'signed' loosely. What we mean is that the sufficiency propositions result in a signing of the rules obtained in step 4 and 5. Stratification of sufficiency propositions has a similar meaning.

$$holds(open, res(turning\_switch, s_1)) \leftarrow$$
$$\neg holds(open, s_1)$$
$$ab_0(open', turning\_switch, s_1) \leftarrow$$
$$\neg holds(open, s_1)$$
$$\neg holds(open, res(turning\_switch, s_1)) \leftarrow$$
$$holds(open, s_1)$$
$$ab_0(open, turning\_switch, s_1) \leftarrow holds(open, s_1)$$

**Non-defeasible rules:**

$$\neg holds(light\_on, res(A, s_1)) \leftarrow holds(open, res(A, s_1))$$
$$ab_0(light\_on', A, s_1) \leftarrow holds(open, res(A, s_1))$$
$$holds(light\_on, res(A, s_1)) \leftarrow \neg holds(open, res(A, s_1))$$
$$ab_0(light\_on', A, s_1) \leftarrow \neg holds(open, res(A, s_1))$$

**Inertia rules:**

$$holds(F, res(turning\_switch, s_1)) \leftarrow holds(F, s_1),$$
$$not\ ab_0(F, turning\_switch, s_1),$$
$$not\ ab_1(F, turning\_switch, s_1)$$
$$\neg holds(F, res(turning\_switch, s_1)) \leftarrow \neg holds(F, s_1),$$
$$not\ ab_0(F', turning\_switch, s_1),$$
$$not\ ab_1(F', turning\_switch, s_1)$$

The effect rules depend on the initial state $S$ which will be part of any answer set of the program. Hence, the set

$$E = \{\neg holds(open, res(turning\_switch, s_1)),$$
$$ab_0(open, turning\_switch, s_1)\}$$

will also be part of any answer set. Therefore, from the third and fourth non-defeasible rules we get that the set:

$$ND = \{holds(light\_on, res(turning\_switch, s_1)),$$
$$ab_0(light\_on', turning\_switch, s_1)\}$$

will be part of any answer set too. The inertia rules, based on the initial state, will try to make the fluent literals *power, -light_on* and *open* true in the situation *re8(turning switch, s\)*. We can ignore the $ab_1$ predicates since there are no defeasible rules. However, *abo* in *E* blocks *open* and *abo* in *ND* blocks *-light-on*. Hence, the inertia rules generate the set

$$I = \{holds(power, res(turning\_light, s_1))\}$$

No other *holds* literal with situation argument *res(turning_ight,8\)* or *abo* predicate with third argument $s_1$ will be in any answer set. Notice too, that the evaluation of the rules was stratified, first the effect rules, then the non-defeasible rules and finally the inertia rules. Hence, it is easy to see that $Res(turning\_switch, s_1) = s_2$ and it is completely determined by $E \cup ND \cup I$.

It can also be shown that $Res(turning switch, s_2) = s_1$. □

Example 3.2 If we add to the previous example the causal laws (5) and (6), we can see that the answer set of $\Pi(D_1, unplug, s_2)$ is inconsistent. This is because we will have the effect rule:

$$\neg holds(power, res(unplug, s_2)) \qquad (12)$$

that is unconditionally true, and the following two non-defeasible sufficiency rules:

$$\neg holds(light\_on, res(unplug, s_2)) \leftarrow$$
$$\neg holds(power, res(unplug, s_2))$$
$$holds(light\_on, res(unplug, s_2)) \leftarrow$$
$$\neg holds(open, res(unplug, s_2))$$

The first rule justifies $\neg holds(light\_on, res(unplug, s_2))$ because of the effect rule and the second rule will justify $holds(light\_on, res(unplug, s_2))$ based on the inertia rule:

$$\neg holds(open, res(unplug, s_2)) \leftarrow \neg holds(open, s_2),$$
$$not\ ab_0(open', unplug, s_2),$$
$$not\ ab_1(open', unplug, s_2)$$

The literal *-holds(open,$S_2$)* is part of the initial database, *open'* will never appear in the first argument of any *abo* in the translation and *ab\* is irrelevant in the example. □

Example 3.3 In this final example, we show the correct representation of the problem. We need the causal laws (1), (2), (5) and (6), plus the defeasible sufficiency propositions:

**open normally suffices for $\neg light\_on$**
**$\neg open$ normally suffices for $light\_on$**

In this case the second non-defeasible sufficiency rule from the previous example is replaced by the rule

$$holds(light\_on, res(unplug, s_2)) \leftarrow$$
$$\neg holds(open, res(unplug, s_2)),$$
$$not\ ab_0(light\_on, unplug, s_2)$$

and this rule will not be justified since the following rule is also part of the translation

$$ab_0(light\_on, unplug, s_2) \leftarrow \neg$$
$$holds(power, res(unplug, s_2))$$

which is justified by (12). O

Observe that the translation presented above is equivalent to a simpler version where, we replace both *abo* and $ab_1$ by *ab,* and in the <u>first rule of step 5</u> we replace *not $ab_0(f,A,s)$* by *not $Hold8(f, re8(A,S))$*.[8] This is similar to Turner's [11] use of

*not Holds(f, res(A,S))* instead of an abnormal predicate to define sufficiency propositions. However, with Turner's translation there is no obvious way to add defeasible causality to his formalism. Furthermore, the structure of our translation (compared to Gelfond's suggestion) allows us to easily add more levels of defeasible propositions; i.e., defeasible propositions that can be defeated by other defeasible propositions[9] which could

---

[8] This simplification was <u>suggested by Michael Gelfond</u>. The subtle use of the *not Holds(f,<u>res(A,</u> S))* instead of *Holda(f,res(A, S))* was first introduced in [5] to reason about actions in presence of incomplete information.

[9] Geffner in [3] also deals with hierarchy of defeasible propositions by assigning non-negative integer priorities to them and then characterizes them using an approach based on infinitesimal probabilities. We plan to compare our approaches in the full paper.

be also be defeated by other propositions. In general, when there are several levels of defeasibility the guard *not Holds(f,res(A,S))* can be used in the last (meaning highest) but one layer. All layers below that need an *ab* predicate. In the language *AVC* the lowest layer corresponds to inertia, the next layer corresponds to the defeasible causality and the highest layer is the effect propositions and non-defeasible causality propositions.

# 4 Semantics of *AVC* and a logic programming translation

In this section we use the formalization of *Res* to present a semantics of *AVC* which allows us to make inferences about the effect of a sequence of actions.

An *interpretation* is a partial function $\Psi$ from sequences of actions to states such that:

(1) The empty sequence [ ] belongs to the domain of $\Psi$ and

(2) $\Psi$ is prefix-closed [10].

$\Psi([\,])$ is called the initial state of $\Psi$.

Given a fluent $f$ and a state $\sigma$, we say that $f$ holds in $\sigma$ ($f$ is *true* in $\sigma$) if $f \in \sigma$; $\neg f$ holds in $\sigma$ ($f$ is *false* in $\sigma$) if $f \notin \sigma$.

**Definition 4.1** An interpretation $\Psi$ *satisfies* the causal laws of $D$ if for any sequence $\alpha \circ a$ from the language of $D$,

if $Res(a, \Psi(\alpha))$ is defined

then $\Psi(\alpha \circ a) \in Res(a, \Psi(\alpha))$,

otherwise $\Psi(\alpha \circ a)$ is undefined.

We say $\Psi$ is a *causal model* of $D$ if it satisfies the causal laws of $D$. $\square$

We say an initial proposition of the form **initially** $f$ is true in an interpretation $\Psi$ if $f$ is true in $\Psi([])$.

**Definition 4.2** An interpretation $\Psi$ will be called a *model* of a domain description $D$ in $\mathcal{ADC}$ if the following conditions are satisfied:

1. $\Psi$ is a causal model of $D$.

2. Initial propositions of $D$ are true in $\Psi([])$. $\square$

A query in $\mathcal{ADC}$ is an expression of the form

$$f \text{ after } a_1, \dots, a_n \qquad (13)$$

where, $f$ is a fluent literal and $a_i$'s are actions.

**Definition 4.3** We say a query of the form (13) is true in an interpretation $\Psi$ if $f$ is true in $\Psi([a_1, \dots, a_n])$.

We say $D \models f$ **after** $a_1, \dots, a_n$ if $f$ is true in all models of $\Psi([a_1, \dots, a_n])$. $\square$

---

[10] By "prefix closed" we mean that for any sequence of actions $\alpha$ and action $a$, if $\alpha \circ a$ is in the domain of $\Psi$ then so is $\alpha$.

## 4.1 Logic programming characterization

We can modify the translation $\Pi\{D, a, s\}$ to give a logic programming characterization of *AVC*. (Note that $\Pi(D, a, s)$ only characterizes the effect of a single action and to characterize *AVC* we need to be able to reason about sequences of actions.) We introduce a new predicate *p_holds* (which intuitively means 'possibly holds') and replace all occurrences of *holds(f,res(a,s))* and *-holds(f,res(a,s))* in the steps 1, 3, 4 and 5, and the head of the rules in step 2, by *p.holds(f, res(a, S))* and *p-holds(f,res(a, S))*. We then modify step 2 by replacing the constants *a* and *s*, by the variables *A* and *S*. Then, we add the following rules:[11]

$\neg reachable(S) \leftarrow p\_holds(F, S), p\_holds(F', S)$
$reachable(S) \leftarrow not \ \neg reachable(S)$
$holds(F, S) \leftarrow p\_holds(F, S), reachable(S)$
$\neg holds(F, S) \leftarrow p\_holds(F', S), reachable(S)$

The first rule verifies if the effects produced by actions when we are trying to reach the situation $S$ are contradictory. This indicates that $S$ should be unreachable. The second rule says that $S$ is reachable otherwise.[12] And, finally we add rules to make sure that sufficiency propositions are true w.r.t. the initial state $S_0$. For a non-defeasible sufficiency proposition of the form (10) we add the rules:

$p\_Holds(f, s_0) \leftarrow p\_Holds(p_1, s_0), \dots, p\_Holds(p_n, s_0).$
$ab_{init}(\bar{f}, s_0) \leftarrow p\_Holds(p_1, s_0), \dots, p\_Holds(p_n, s_0).$

For defeasible sufficiency propositions of the form (11) we add the rule:

$p\_Holds(f, s_0) \leftarrow p\_Holds(p_1, s_0), \dots, p\_Holds(p_n, s_0),$
$\qquad not \ ab_{init}(\bar{f}, s_0).$

Let us refer to the resultant program as $\Pi(D)$.

**Proposition 4.1** For any domain description $D$ with complete information about the initial state

$D \models f$ **after** $a_1, \dots, a_n$ iff $\Pi(D) \models holds(f, [a_1, \dots, a_n])$ $\square$

## 5 Defeasible effect propositions

Our new approach of defining the semantics of the action language as a logic program gives us the flexibility to directly add defeasible effect propositions to *AVC*. Consider the following effect proposition.

*open-faucet* causes *water_on_sink.*

However, in the rare occasion when the pipe is blocked the action will not have any effect[13]. A better description for the effect proposition would be:

---

[11] Similar to $Holds(f, s)$, the atom $p\_Holds(f, s)$ denotes the literal $p\_holds(f, s)$ if $f$ is a positive fluent literal and the literal $p\_holds(\bar{f}, s)$ if $f$ is a negative fluent literal.

[12] This indirect definition of *reachable* is to be able to represent $\forall(F)(p\_holds(F, S), p\_holds(F', S))$ in the logic rules.

[13] Note the similarity to McCarthy's example of the potato in the pipe.

*open-faucet* normally_causes *water_on _sink*.

The constraint about the case when the pipe is blocked can then be specified as:

*pipe blocked* suffices for $\neg water\_on\_sink$.

In general, we will call a defeasible effect proposition an expression of the form

$$a \text{ normally\_causes } f \text{ if } p_1, \ldots, p_n \qquad (14)$$

where $a$ is an action, and $l$ and and each of $p_1, \ldots, p_n$ $(0 \leq n)$ is a fluent literal.

To describe the semantics of a domain description $D$ with defeasible effect propositions we expand the program $\Pi(D, a, s)$ as follows.

### 6. Translating Defeasible Effect propositions

For each defeasible effect proposition of the form (14) in $D$ which describes the effect of the action a, $\Pi(D, a, s)$ contains the following rules:

$$Holds(f, res(a, s)) \leftarrow Holds(p_1, s), \ldots, Holds(p_n, s),$$
$$not \ pre\_ab(\overline{f}, s)$$
$$ab_0(\overline{f}, a, s) \leftarrow Holds(p_1, s), \ldots, Holds(p_n, s),$$
$$not \ pre\_ab(\overline{f}, s)$$

We also need to add to the two rules of the translation of non-defeasible sufficiency propositions the following rule (step 4 of the original transformation):

$$pre\_ab(\overline{f}, res(a, s)) \leftarrow Holds(p_1, s), \ldots, Holds(p_n, s).$$

To expand the logic programming characterization of *AVC* domains we merely need to replace the *Holds* in the rules of part (6) with the corresponding *p_holds* and proposition 4.1 holds for the new domain descriptions.

In the above formalization we assumed that defeasible effect propositions are overridden by non-defeasible causal rules and defeasible effect propositions do not contradict with non-defeasible effect propositions. A different formalization can be given when we modify the assumptions such as (i) allowing defeasible effect propositions to contradict with non-defeasible effect propositions and giving non-defeasible effect propositions a higher priority; (ii) giving effect propositions higher priority than sufficiency propositions. The new formalizations will mainly express the priority between the different abnormal predicates. We will discuss this in greater detail in the full paper.

## 6 Conclusion

We have showed how to formalize defeasible causality in action theories. We used logic programming to encode the defeasible causalities. It is not clear to us if use of a logical formalism can be avoided (in the spirit of *A* [5]) while encoding defeasibility, particularly when we allow several level of defeasibility. In the full paper we will add additional examples involving several levels of defeasibility and discuss how defeasibility can be incorporated in other specification languages. To make our point clear

we avoided additional constructs and features. Most of these can be easily added to the proposed language, and we will need default logic (as in [11]), instead of logic programming, to encode defeasibility when we allow fluent formulas instead of fluent literals.

## References

[1] C. Baral. Reasoning about Actions : Non-deterministic effects, Constraints and Qualification. In *Proc. of IJCAI 95,* pages 2017-2023, 1995.

[2] C. Baral, M. Gelfond, and A. Provetti. Representing Actions: Laws, Observations and Hypothesis. *Journal of Logic Programming,* 31(1-3):201-243, May 1997.

[3] H. Geffner. A qualitative model for Temporal reasoning with incomplete information. In *AAAI 96,* pages 1176-1181. 1996.

[4] M. Gelfond and V. Lifschitz. Classical negation in logic programs and disjunctive databases. *New Generation Computing,* pages 365-387, 1991.

[5] M. Gelfond and V. Lifschitz. Representing actions and change by logic programs. *Journal of Logic Programming,* 17(2,3,4):301 323, 1993.

[6] G. Kartha and V. Lifschitz. Actions with indirect effects: Preliminary report. In K*R 94* pages 341 350, 1994.

[7] F. Lin. Embracing causality in specifying the indirect effects of actions. In *Proc. of 1JCAI 95,* pages 1985-1993,95.

[8] N. McCain and M. Turner. A causal theory of ramifications and qualifications. In *Proc. of IJCAI 95,* pages 1978-1984, 95.

[9] R. Reiter. The frame problem in the situation calculus: A simple solution (sometimes) and a completeness result for goal regression. In V. Lifschitz, editor, *Artificial Intelligence and Mathematical Theory of Computation,* pages 359-380. Academic Press, 1991.

[10] E. Sandewall. The range of applicability of some non-monotonic logics for strict inertia. *Journal of Logic and Computation,* 4(5):581-616, October 1994.

[11] H. Turner. Representing actions in default logic: A situation calculus approach. In *Proceedings of the Symposium in honor of Michael Gelfond 's 50th birthday (also in Common Sense 96),* 1995.

[12] Y. Zhang. Compiling causality into action theories. In *Proc, of Common Sense 96,* pages 263-270,1996.