

# Prioritized Goals and Subgoals in a Logical Account of Goal Change – A Preliminary Report

Shakil M. Khan and Yves Lespérance

Department of Computer Science and Engineering  
York University, Toronto, ON, Canada  
{skhan,lesperan}@cse.yorku.ca

**Abstract.** Most previous logical accounts of goal change do not deal with prioritized goals and do not handle subgoals and their dynamics properly. Many are restricted to achievement goals. In this paper, we develop a logical account of goal change that addresses these deficiencies. In our account, we do not drop lower priority goals permanently when they become inconsistent with other goals and the agent's knowledge; rather, we make such goals inactive. We ensure that the agent's chosen goals/intentions are consistent with each other and the agent's knowledge. When the world changes, the agent recomputes her chosen goals and some inactive goals may become active again. This ensures that our agent maximizes her utility. We also propose an approach for handling subgoals and their dynamics. We prove that the proposed account has some intuitively desirable properties.

## 1 Introduction

There has been much work on modeling agent's mental states, beliefs, goals, and intentions, and how they interact and lead to rational decisions about action. As well, there has been a lot of work on modeling belief change. But the dynamics of motivational attitudes has received much less attention. Most formal models of goal and goal change [1,2,3,4,5,6] assume that all goals are equally important and many only deal with achievement goals (one exception to this is the model of prioritized goals in [7]). Moreover, most of these frameworks do not guarantee that an agent's goals will properly evolve when an action/event occurs, e.g. when the agent's beliefs/knowledge changes or a goal is adopted or dropped. Also, they do not model the dependencies between goals and the subgoals and plans adopted to achieve these goals. For instance, subgoals and plans adopted to bring about a goal should be dropped when the parent goal becomes impossible, is achieved, or is dropped. Dealing with these issues is important for developing effective models of rational agency. It is also important for work on BDI agent programming languages, where handling declarative goals is an active research topic [8,9].

In this paper, we present a formal model of prioritized goals and their dynamics that addresses some of these issues. Specifically, we propose a framework, where an agent can have multiple goals at different priority levels, possibly inconsistent with each other. We define intentions as the maximal set of highest priority goals that is consistent given the agent's knowledge. Our model supports the specification of general

temporally extended goals, not just achievement goals, and handles subgoals and their dynamics.

We start with a (possibly inconsistent) initial set of *prioritized goals* or desires that are totally ordered according to priority, and specify how these goals evolve when actions/events occur and the agent's knowledge changes. We define the agent's *chosen goals* or intentions in terms of this goal hierarchy. Our agents maximize their utility; they will abandon a chosen goal  $\phi$  if an opportunity to commit to a higher priority but inconsistent with  $\phi$  goal arises. To this end, we keep all prioritized goals in the goal base unless they are explicitly dropped. At every step, we compute an optimal set of chosen goals given the hierarchy of prioritized goals, preferring higher priority goals, such that chosen goals are consistent with each other and with the agent's knowledge. Thus at any given time, some goals in the hierarchy are active, i.e. chosen, while others are inactive. Some of these inactive goals may later become active, e.g. if a higher priority active goal that is currently blocking an inactive goal becomes impossible. We also show how the dependencies between goals and subgoals can be modeled. Finally, we prove some interesting properties about the dynamics of chosen goals.

As mentioned above, our formalization of prioritized goals ensures that the agent always tries to maximize her utility, and as such a limitation of our framework is that it displays an idealized form of rationality. In Section 5, we discuss how this relates to Bratman's theory of practical reasoning [10]. We use an action theory based on the situation calculus [11] along with our formalization of paths in the situation calculus as our base formalism.

The paper is organized as follows: in the next section, we outline our base framework. In Section 3, we formalize *paths* in the situation calculus to support modeling temporally extended goals. In Section 4, we present our model of prioritized goals. In section 5, we present our formalization of goal dynamics and discuss some of its properties. In Section 6, we discuss what it means for an agent to have a subgoal and how subgoals change as a result of changes to their parent goals. Then in the last section, we summarize our results, discuss previous work in this area, and point to possible future work.

## 2 Action and Knowledge

Our base framework for modeling goal change is the situation calculus [11] as formalized in [12]. In this framework, a possible state of the domain is represented by a situation. There is a set of initial situations corresponding to the ways the agents believe the domain might be initially, i.e. situations in which no actions have yet occurred.  $\text{Init}(s)$  means that  $s$  is an initial situation. The actual initial state is represented by a special constant  $S_0$ . There is a distinguished binary function symbol  $do$  where  $do(a, s)$  denotes the successor situation to  $s$  resulting from performing the action  $a$ . Thus the situations can be viewed as a set of trees, where the root of each tree is an initial situation and the arcs represent actions. Relations (and functions) whose truth values vary from situation to situation, are called relational (functional, resp.) fluents, and are denoted by predicate (function, resp.) symbols taking a situation term as their last argument. There is a special predicate  $\text{Poss}(a, s)$  used to state that action  $a$  is executable in situation  $s$ .

Our framework uses a theory  $D_{basic}$  that includes the following set of axioms:<sup>1</sup> (1) action precondition axioms, one per action  $a$  characterizing  $Poss(a, s)$ , (2) successor state axioms (SSA), one per fluent, that succinctly encode both effect and frame axioms and specify exactly when the fluent changes [12], (3) initial state axioms describing what is true initially including the mental states of the agents, (4) unique name axioms for actions, and (5) domain-independent foundational axioms describing the structure of situations [14].

Following [15,16], we model knowledge using a possible worlds account adapted to the situation calculus.  $K(s', s)$  is used to denote that in situation  $s$ , the agent thinks that she could be in situation  $s'$ . Using  $K$ , the knowledge of an agent is defined as:<sup>2</sup>  $Know(\Phi, s) \stackrel{\text{def}}{=} \forall s'. K(s', s) \supset \Phi(s')$ , i.e. the agent knows  $\Phi$  in  $s$  if  $\Phi$  holds in all of her  $K$ -accessible situations in  $s$ .  $K$  is constrained to be reflexive, transitive, and Euclidean in the initial situation to capture the fact that agents' knowledge is true, and that agents have positive and negative introspection. As shown in [16], these constraints then continue to hold after any sequence of actions since they are preserved by the successor state axiom for  $K$ . We also assume that all actions are public, i.e. whenever an action (including exogenous events) occurs, the agent learns that it has happened. Note that, we work with knowledge rather than belief. Although much of our formalization should extend to the latter, we leave this for future work.

### 3 Paths in the Situation Calculus

To support modeling temporally extended goals, we introduce a new sort of *paths*, with (possibly sub/super-scripted) variables  $p$  ranging over paths. A path is essentially an infinite sequence of situations, where each successor situation along the path can be reached by performing some *executable* action in the preceding situation. We introduce a predicate  $OnPath(p, s)$ , meaning that the situation  $s$  is on the path  $p$ . Also,  $Starts(p, s)$  means that  $s$  is the starting situation of path  $p$ . A path  $p$  starts with the situation  $s$  iff  $s$  is the earliest situation on  $p$ :<sup>3</sup>

#### Axiom 1

$$Starts(p, s) \equiv OnPath(p, s) \wedge \forall s'. OnPath(p, s') \supset s \leq s'.$$

In the standard situation calculus, paths are implicitly there, and a path can be viewed as a pair  $(s, F)$  consisting of a situation  $s$  representing the starting situation of the path, and a function  $F$  from situations to actions (called *Action Selection Functions* (ASF) or

<sup>1</sup> We will be quantifying over formulae, and thus we assume  $D_{basic}$  includes axioms for encoding formulae as first order terms, as in [13]. We will also be using lists of integers, and assume that  $D_{basic}$  includes axiomatizations of integers and lists.

<sup>2</sup> A state formula  $\Phi(s)$  takes a single situation as argument and is evaluated with respect to that situation.  $\Phi$  may contain a placeholder constant *now* that stands for the situation in which  $\Phi$  must hold.  $\Phi(s)$  is the formula that results from replacing *now* by  $s$ . Where the intended meaning is clear, we sometimes suppress the placeholder.

<sup>3</sup> In the following,  $s < s'$  means that  $s'$  can be reached from  $s$  by performing a sequence of executable actions.  $s \leq s'$  is an abbreviation for  $s < s' \vee s = s'$ .

strategies in [5]), such that from the starting situation  $s$ ,  $F$  defines an infinite sequence of situations by specifying an action for every situation starting from  $s$ . Thus, one way of axiomatizing paths is by making them correspond to such pairs  $(s, F)$ :

## Axiom 2

$$\forall p. \text{Starts}(p, s) \supset (\exists F. \text{Executable}(F, s) \wedge \forall s'. \text{OnPath}(p, s') \equiv \text{OnPathASF}(F, s, s')),$$

$$\forall F, s. \text{Executable}(F, s) \supset \exists p. \text{Starts}(p, s) \wedge \forall s'. \text{OnPathASF}(F, s, s') \equiv \text{OnPath}(p, s').$$

This says that for every path there is an executable ASF that produces exactly the sequence of situations on the path from its starting situation. Also, for every executable ASF and situation, there is a path that corresponds to the sequence of situations produced by the ASF starting from that situation.

$$\text{OnPathASF}(F, s, s') \stackrel{\text{def}}{=} s \leq s' \wedge \forall a, s^*. s < do(a, s^*) \leq s' \supset F(s^*) = a,$$

$$\text{Executable}(F, s) \stackrel{\text{def}}{=} \forall s'. \text{OnPathASF}(F, s, s') \supset \text{Poss}(F(s'), s').$$

Here,  $\text{OnPathASF}(F, s, s')$  [6] means that the situation sequence defined by  $(s, F)$  includes the situation  $s'$ . Also, the situation sequence encoded by a strategy  $F$  and a starting situation  $s$  is executable iff for all situations  $s'$  on this sequence, the action selected by  $F$  in  $s'$  is executable in  $s'$ .

In our framework, we will use both state and path formulae. A state formula is a formula that has a free situation variable in it, whereas a path formula is one that has a free path variable. State formulae are used in the context of knowledge while path formulae are used in that of goals. We use  $\Phi(s), \Psi(s), \dots$  and  $\phi(p), \psi(p), \dots$  possibly with decorations to represent state and path formulae, respectively. Note that, by incorporating infinite paths in our framework, we can evaluate goals over these and handle arbitrary temporally extended goals; thus, unlike some other situation calculus based accounts where goal formulae are evaluated w.r.t. finite paths (e.g. [7]), we can handle for example unbounded maintenance goals.

We next define some useful constructs. A state formula  $\Phi$  *eventually holds* over the path  $p$  if  $\Phi$  holds in some situation that is on  $p$ , i.e.  $\diamond\Phi(p) \stackrel{\text{def}}{=} \exists s'. \text{OnPath}(p, s') \wedge \Phi(s')$ . Other Temporal Logic operators can be defined similarly, e.g. always  $\Phi$ :  $\square\Phi(p)$ .

An agent *knows* in  $s$  that  $\phi$  has become *inevitable* if  $\phi$  holds over all paths that starts with a  $K$ -accessible situation in  $s$ :

$$\text{KInevitable}(\phi, s) \stackrel{\text{def}}{=} \forall p. \text{Starts}(p, s') \wedge K(s', s) \supset \phi(p).$$

An agent knows in  $s$  that  $\phi$  is impossible if she knows that  $\neg\phi$  is inevitable in  $s$ , i.e.  $\text{KImpossible}(\phi, s) \stackrel{\text{def}}{=} \text{KInevitable}(\neg\phi, s)$ .

Thirdly, we define what it means for a path  $p'$  to be a suffix of another path  $p$  w.r.t. a situation  $s$ :

$$\text{Suffix}(p', p, s) \stackrel{\text{def}}{=} \text{OnPath}(p, s) \wedge \text{Starts}(p', s)$$

$$\wedge \forall s'. s \leq s' \supset \text{OnPath}(p, s') \equiv \text{OnPath}(p', s').$$

That is, a path  $p'$  is a suffix of another path  $p$  w.r.t. a situation  $s$  iff  $s$  is on  $p$ , and  $p'$ , which starts with  $s$ , is exactly the same as the subpath of  $p$  that starts with  $s$ .

Fourthly,  $\text{SameHistory}(s_1, s_2)$  means that the situations  $s_1$  and  $s_2$  share the same history of actions, but perhaps starting from different initial situations:

### Axiom 3

$$\begin{aligned} \text{SameHistory}(s_1, s_2) &\equiv (\text{Init}(s_1) \wedge \text{Init}(s_2)) \\ &\vee (\exists a, s'_1, s'_2. s_1 = \text{do}(a, s'_1) \wedge s_2 = \text{do}(a, s'_2) \wedge \text{SameHistory}(s'_1, s'_2)). \end{aligned}$$

Thus, if  $s_1$  can be reached from some initial situation by performing a sequence of actions  $\sigma$ , then  $s_2$  can be reached from a (possibly different) initial situation by executing  $\sigma$ .

Finally, we say that  $\phi$  has become *inevitable* in  $s$  if  $\phi$  holds over all paths that starts with a situation that has the same action history as  $s$ :

$$\text{Inevitable}(\phi, s) \stackrel{\text{def}}{=} \forall p, s'. \text{Starts}(p, s') \wedge \text{SameHistory}(s', s) \supset \phi(p).$$

## 4 Prioritized Goals

Most work on formalizing goals only deals with static goal semantics and not their dynamics. There are two main categories of motivational attitudes, namely goal [1,17] (AKA choice [2], wish [10] or preference), and intention. While goals are sometimes allowed to be inconsistent [10], intentions are mostly required to be consistent. Another difference is that agents are committed to their intentions, but not necessarily to their goals [10]. Intention is sometimes primitive [17,3] and sometimes a defined concept, specified in terms of goals [1,2,4]. In this section, we formalize goals or desires with different priorities, which we call *prioritized goals* (p-goals, henceforth). These p-goals are not required to be mutually consistent and need not be actively pursued by the agent. In terms of these, we define the consistent set of *chosen goals* or intentions (c-goals, henceforth) that the agent is committed to. In the next section, we formalize goal dynamics by providing a SSA for p-goals. The agent's c-goals are automatically updated when her p-goals change. We deal with subgoals and their dynamics in Section 6.

Not all of the agent's goals are equally important to her. Thus, it is useful to support a priority ordering over goals. This information can be used to decide which of the agent's c-goals should no longer be actively pursued in case they become mutually inconsistent. Following [6], we specify each p-goal by its own accessibility relation/fluents  $G$ . A path  $p$  is  $G$ -accessible at priority level  $n$  in situation  $s$  (denoted by  $G(p, n, s)$ ) if all the goals of the agent at level  $n$  are satisfied over this path and if it starts with a situation that has the same history (in terms of the actions performed so far) as  $s$ . The latter requirement ensures that the agent's p-goal-accessible paths reflect the actions that have been performed so far. A smaller  $n$  represents higher priority, and the highest priority level is 0. Thus in this framework, we assume that the set of p-goals are totally ordered according to priority. We say that an agent has the p-goal that  $\phi$  at level  $n$  in situation  $s$  iff  $\phi$  holds over all paths that are  $G$ -accessible at  $n$  in  $s$ :

$$\text{PGoal}(\phi, n, s) \stackrel{\text{def}}{=} \forall p. G(p, n, s) \supset \phi(p).$$

To be able to refer to all the p-goals of the agent at some given priority level, we also define *only p-goals*.

$$\text{OPGoal}(\phi, n, s) \stackrel{\text{def}}{=} \text{PGoal}(\phi, n, s) \wedge (\forall p. \phi(p) \supset G(p, n, s)).$$

An agent has the *only p-goal* that  $\phi$  at level  $n$  in situation  $s$  iff  $\phi$  is a p-goal at  $n$  in  $s$ , and any path over which  $\phi$  holds is  $G$ -accessible at  $n$  in  $s$ .

A domain theory for our framework  $D$  includes the axioms of a theory  $D_{\text{basic}}$  as in the previous section, the axiomatization of paths i.e. axioms 1-3, domain dependent initial goal axioms (see below), the domain independent axioms 4-7 and the definitions that appear in this section and the next. The modeler must provide initial goal axioms of the following form:

### INITIAL GOAL AXIOMS

- (a)  $\text{Init}(s) \supset ((G(p, 0, s) \equiv \text{Starts}(p, s') \wedge \text{Init}(s') \wedge \phi_0(p))$   
 $\wedge (G(p, 1, s) \equiv \text{Starts}(p, s') \wedge \text{Init}(s') \wedge \phi_1(p)) \wedge \dots$   
 $\wedge (G(p, k-1, s) \equiv \text{Starts}(p, s') \wedge \text{Init}(s') \wedge \phi_{k-1}(p)))$ ,
- (b)  $\forall n, p, s. \text{Init}(s) \wedge n \geq k \supset (G(p, n, s) \equiv \text{Starts}(p, s') \wedge \text{Init}(s'))$ ,
- (c)  $\text{Init}(s) \supset \text{NPGoals}(s) = k$ .

The p-goals  $\phi_0, \phi_1, \dots, \phi_{k-1}$  (from highest to lowest priority) of the agent in the initial situations are specified by the Initial Goal Axiom (a); each of them defines a set of initial goal paths for a given priority level, and must be consistent. We assume that the agent has a finite number  $k$  of initial p-goals. For  $n \geq k$ , we make  $G(p, n, s)$  true for every path  $p$  that starts with an initial situation in (b). Thus at levels  $n \geq k$ , the agent has the trivial p-goal that she be in an initial situation. We also have a distinguished functional fluent  $\text{NPGoals}(s)$  that represents the number of prioritized goals that the agent has (i.e. the location of the first empty slot after the last p-goal). Initially  $\text{NPGoals}$  is set to  $k$  in (c). Later, we will specify the dynamics of p-goals by giving SSAs for  $G$  and  $\text{NPGoals}$ .

We use the following as a running example. We have an agent who initially has the following three p-goals:  $\phi_0 = \Box\text{BeRich}$ ,  $\phi_1 = \Diamond\text{GetPhD}$ , and  $\phi_2 = \Box\text{BeHappy}$  at level 0, 1, and 2, respectively (see second column of Table 1). Assume that while initially the agent knows that all of her p-goals are individually achievable, she knows that her p-goal  $\Diamond\text{GetPhD}$  is inconsistent with her highest priority p-goal  $\Box\text{BeRich}$  as well as with her p-goal  $\Box\text{BeHappy}$ , while the latter are consistent with each other. It is straightforward to specify a domain action theory such that it entails this. Thus in our example, we have  $\text{OPGoal}(\phi_i(p) \wedge \text{Starts}(p, s) \wedge \text{Init}(s), i, S_0)$ , for  $i = 0, 1, 2$ . Also, for any  $n \geq 3$ , we have  $\text{OPGoal}(\text{Starts}(p, s) \wedge \text{Init}(s), n, S_0)$ .

**Table 1.** Example of an Agent's PGoals and their Dynamics

G-Level	$S_0, S'_1$	$S_1$	$S_2$	$S_3$
4	TRUE	TRUE	$\Box\text{BeRich} \wedge \Box\text{WorkHard} \wedge \Box\text{BeEnergetic}$	TRUE
3	TRUE	$\Box\text{BeRich} \wedge \Box\text{WorkHard}$	$\Box\text{BeRich} \wedge \Box\text{WorkHard}$	TRUE
2	$\Box\text{BeHappy}$	$\Box\text{BeHappy}$	$\Box\text{BeHappy}$	$\Box\text{BeHappy}$
1	$\Diamond\text{GetPhD}$	$\Diamond\text{GetPhD}$	$\Diamond\text{GetPhD}$	$\Diamond\text{GetPhD}$
0	$\Box\text{BeRich}$	$\Box\text{BeRich}$	$\Box\text{BeRich}$	$\Box\text{BeRich}$

While p-goals or desires are allowed to be known to be impossible to achieve, an agent's c-goals or intentions must be realistic. Not all of the  $G$ -accessible paths are realistic in the sense that they start with a  $K$ -accessible situation. To filter these out, we define *realistic* p-goal accessible paths:

$$G_R(p, n, s) \stackrel{\text{def}}{=} G(p, n, s) \wedge \text{Starts}(p, s') \wedge K(s', s),$$

i.e., a path  $p$  is  $G_R$ -accessible at level  $n$  in situation  $s$  if it is  $G$ -accessible at  $n$  in  $s$ , and if  $p$  starts with a situation that is  $K$ -accessible in  $s$ . Thus  $G_R$  prunes out the paths from  $G$  that are known to be impossible, and since we define c-goals in terms of realistic p-goals, this ensures that c-goals are realistic. We say that an agent has the *realistic p-goal* that  $\phi$  at level  $n$  in situation  $s$  iff  $\phi$  holds over all paths that are  $G_R$ -accessible at  $n$  in  $s$ :

$$\text{RPGoal}(\phi, n, s) \stackrel{\text{def}}{=} \forall p. G_R(p, n, s) \supset \phi(p).$$

Using realistic p-goals, we next define c-goals. The idea of how we compute c-goal-accessible paths is as follows: the set of  $G_R$ -accessibility relations represents a set of prioritized temporal propositions that are candidates for the agent's c-goals. Given  $G_R$ , in each situation we want to compute the agent's c-goals such that it is the *maximal consistent* set of higher priority realistic p-goals. We do this iteratively starting with the set of all possible paths (i.e. paths that starts with a  $K$ -accessible situation). At each iteration we compute the intersection of this set with the next highest priority set of  $G_R$ -accessible paths. If the intersection is not empty, we thus obtain a new chosen set of paths at level  $i$ . We call a p-goal chosen by this process an *active* p-goal. If on the other hand, the intersection is empty, then it must be the case that the p-goal represented by this level is either in conflict with another active higher priority p-goal/a combination of two or more active higher priority p-goals, or is known to be impossible. In that case, that p-goal is ignored (i.e. marked as inactive), and the chosen set of paths at level  $i$  is the same as at level  $i - 1$ . We repeat this until we reach  $i = \text{NPGoals}$ . Axiom 4 "computes" this intersection:<sup>4</sup>

#### Axiom 4

$$\begin{aligned} G_{\cap}(p, n, s) \equiv & \text{if } (n = 0) \text{ then} \\ & \text{if } \exists p'. G_R(p', n, s) \text{ then } G_R(p, n, s) \\ & \text{else } \text{Starts}(p, s') \wedge K(s', s) \\ & \text{else} \\ & \text{if } \exists p'. (G_R(p', n, s) \wedge G_{\cap}(p', n - 1, s)) \\ & \text{then } (G_R(p, n, s) \wedge G_{\cap}(p, n - 1, s)) \\ & \text{else } G_{\cap}(p, n - 1, s). \end{aligned}$$

C-goal accessible paths are the result of this intersection after all priority levels have been considered:

$$G_C(p, s) \stackrel{\text{def}}{=} G_{\cap}(p, \text{NPGoals}(s) - 1, s).$$

<sup>4</sup> **if**  $\phi$  **then**  $\psi_1$  **else**  $\psi_2$  is an abbreviation for  $(\phi \supset \psi_1) \wedge (\neg\phi \supset \psi_2)$ .

We define an agent's c-goals in terms of the  $G_C$ -accessible paths:

$$\text{CGoal}(\phi, s) \stackrel{\text{def}}{=} \forall p. G_C(p, s) \supset \phi(p),$$

i.e., the agent has the c-goal that  $\phi$  if  $\phi$  holds over all of her  $G_C$ -accessible paths.

We also define what it means for an agent to have a c-goal at some level  $n$ :

$$\text{CGoal}(\phi, n, s) \stackrel{\text{def}}{=} \forall p. G_{\cap}(p, n, s) \supset \phi(p),$$

i.e. an agent has the c-goal at level  $n$  that  $\phi$  if  $\phi$  holds over all paths that are in the prioritized intersection of the set of  $G_R$ -accessible paths up to level  $n$ .

In our example, the agent's realistic p-goals are  $\square\text{BeRich}$ ,  $\diamond\text{GetPhD}$ , and  $\square\text{BeHappy}$  in order of priority. The  $G_{\cap}$ -accessible paths at level 0 in  $S_0$  are the ones that start with a  $K$ -accessible situation and where  $\square\text{BeRich}$  holds. The  $G_{\cap}$ -accessible paths at level 1 in  $S_0$  are the same as at level 0, since there are no  $K$ -accessible paths over which both  $\diamond\text{GetPhD}$  and  $\square\text{BeRich}$  hold. Finally, the  $G_{\cap}$ -accessible paths at level 2 in  $S_0$  and hence the  $G_C$ -accessible paths are those that start with a  $K$ -accessible situation and over which  $\square\text{BeRich} \wedge \square\text{BeHappy}$  holds. Also, it can be shown that initially our example agent has the c-goals that  $\square\text{BeRich}$  and  $\square\text{BeHappy}$ , but not  $\diamond\text{GetPhD}$ .

Note that by our definition of c-goals, the agent can have a c-goal that  $\phi$  in situation  $s$  for various reasons: 1)  $\phi$  is known to be inevitable in  $s$ ; 2)  $\phi$  is an active p-goal at some level  $n$  in  $s$ ; 3)  $\phi$  is a consequence of two or more active p-goals at different levels in  $s$ . To be able to refer to c-goals for which the agent has a primitive motivation, i.e. c-goals that result from a single active p-goal at some priority level  $n$ , in contrast to those that hold as a consequence of two or more active p-goals at different priority levels, we define *primary* c-goals:

$$\text{PrimCGoal}(\phi, s) \stackrel{\text{def}}{=} \exists n. \text{PGoal}(\phi, n, s) \wedge \exists p. G(p, n, s) \wedge G_{\cap}(p, n, s).$$

That is, an agent has the primary c-goal that  $\phi$  in situation  $s$ , if  $\phi$  is a p-goal at some level  $n$  in  $s$ , and if there is a  $G$ -accessible path  $p$  at  $n$  in  $s$  that is also in the prioritized intersection of  $G_R$ -accessible paths upto  $n$  in  $s$ . The last two conjuncts are required to ensure that  $n$  is an active level. Thus if an agent has a primary c-goal that  $\phi$ , then she also has the c-goal that  $\phi$ , but not necessarily vice-versa. It can be shown that initially our example agent has the primary c-goals that  $\square\text{BeRich}$  and  $\square\text{BeHappy}$ , but not their conjunction. This shows that (strictly speaking) primary c-goals are not closed under logical consequence.

## 5 Goal Dynamics

An agent's goals change when her knowledge changes as a result of the occurrence of an action (including exogenous events), or when she adopts or drops a goal. We formalize this by specifying how p-goals change. C-goals are then computed using realistic p-goals in every new situation as explained above.

We introduce two actions for adopting and dropping a p-goal,  $\text{adopt}(\phi)$  and  $\text{drop}(\phi)$ , and a third for adopting a subgoal  $\psi$  w.r.t. a supergoal  $\phi$ ,  $\text{adopt}(\psi, \phi)$ . The action precondition axioms for these are as follows:



**Axiom 5**

$$\begin{aligned} \text{Poss}(\text{adopt}(\phi), s) &\equiv \neg \exists n. \text{PGoal}(\phi, n, s), \\ \text{Poss}(\text{adopt}(\psi, \phi), s) &\equiv \neg \exists n. \text{PGoal}(\psi, n, s) \wedge \exists n'. \text{PGoal}(\phi, n', s), \\ \text{Poss}(\text{drop}(\phi), s) &\equiv \exists n. \text{PGoal}(\phi, n, s). \end{aligned}$$

That is, an agent can adopt the p-goal that  $\phi$ , if she does not already have  $\phi$  as her p-goal at some level. An agent can adopt a subgoal  $\psi$  w.r.t. the parent goal that  $\phi$  if she does not already have the p-goal that  $\psi$  at some level, and if at some level she currently has the parent goal that  $\phi$ . The  $\text{drop}(\phi)$  action is possible in  $s$  if  $\phi$  is a p-goal at some level  $n$  in  $s$ .

In the following, we specify the dynamics of p-goals by giving the SSA for  $G$  and then discuss each case, one at a time:

**Axiom 6 (SSA for  $G$ )**

$$\begin{aligned} G(p, n, \text{do}(a, s)) &\equiv \\ &\forall \phi, \psi. (a \neq \text{adopt}(\phi) \wedge a \neq \text{adopt}(\psi, \phi) \wedge a \neq \text{drop}(\phi) \wedge \text{Progressed}(p, n, a, s)) \\ &\vee \exists \phi. (a = \text{adopt}(\phi) \wedge \text{Adopted}(p, n, a, s, \phi)) \\ &\vee \exists \phi, \psi. (a = \text{adopt}(\psi, \phi) \wedge \text{SubGoalAdopted}(p, n, a, s, \psi, \phi)) \\ &\vee \exists \phi. (a = \text{drop}(\phi) \wedge \text{Dropped}(p, n, a, s, \phi)). \end{aligned}$$

The overall idea of the SSA for  $G$  is as follows. First of all, to handle the occurrence of a non-adopt/drop (i.e. regular) action  $a$ , we progress all  $G$ -accessible paths to reflect the fact that this action has just happened; this is done using the  $\text{Progressed}(p, n, a, s)$  construct, which replaces each  $G$ -accessible path  $p'$  with starting situation  $s'$ , by its suffix  $p$  provided that it starts with  $\text{do}(a, s')$ :

$$\text{Progressed}(p, n, a, s) \stackrel{\text{def}}{=} \exists p'. G(p', n, s) \wedge \text{Starts}(p', s') \wedge \text{Suffix}(p, p', \text{do}(a, s')).$$

Any path over which the next action performed is not  $a$  is eliminated from the respective  $G$  accessibility level.

Secondly, to handle adoption of a p-goal  $\phi$ , we add a new proposition containing the p-goal to the agent's goal hierarchy. We assume that the newly adopted p-goal  $\phi$  has the lowest priority. Thus in addition to progressing the  $G$ -accessible paths at all levels as above, we eliminate the paths over which  $\phi$  does not hold from the  $\text{NPGoals}(s)$ -th  $G$ -accessibility level, and the agent acquires the p-goal that  $\phi$  at level  $\text{NPGoals}(s)$ :

$$\begin{aligned} \text{Adopted}(p, n, a, s, \phi) &\stackrel{\text{def}}{=} \mathbf{if} (n = \text{NPGoals}(s)) \mathbf{then} (\text{Progressed}(p, n, a, s) \wedge \phi(p)) \\ &\quad \mathbf{else} \text{Progressed}(p, n, a, s). \end{aligned}$$

The third case of subgoal adoption is discussed in the next section.

Finally, to handle dropping of a p-goal  $\phi$ , we replace the propositions that imply the dropped goal in the agent's goal hierarchy by the "trivial" proposition that the history of actions in the current situation has occurred. Thus in addition to progressing all  $G$ -accessible paths as above, we add back all paths that share the same history with  $\text{do}(a, s)$

to the existing  $G$ -accessibility levels where the agent has the p-goal that  $\phi$ , and thus these  $G$ -accessibility levels now amount to the “trivial” p-goal that  $\text{CorrectHist}(s, \text{path})$ .<sup>5</sup>

$$\begin{aligned} \text{Dropped}(p, n, a, s, \phi) &\stackrel{\text{def}}{=} \\ &\mathbf{if} \text{ PGoal}(\phi, n, s) \mathbf{then} \exists s'. \text{Starts}(p, s') \wedge \text{SameHistory}(s', \text{do}(a, s)) \\ &\mathbf{else} \text{ Progressed}(p, n, a, s). \end{aligned}$$

The SSA for  $\text{NPGoals}(s)$  is as follows:

**Axiom 7 (SSA for  $\text{NPGoals}(s)$ )**

$$\begin{aligned} \text{NPGoals}(\text{do}(a, s)) = k &\equiv \\ \neg(\exists \phi. a = \text{adopt}(\phi)) \wedge \neg(\exists \psi. \phi. a = \text{adopt}(\psi, \phi)) \wedge \text{NPGoals}(s) = k &\vee \\ \exists \phi. a = \text{adopt}(\phi) \wedge \text{NPGoals}(s) + 1 = k &\vee \\ \exists \psi, \phi. a = \text{adopt}(\psi, \phi) \wedge \text{AdjustSubGoalAdopt}(\phi, s) = k. & \end{aligned}$$

That is, when the agent adopts a p-goal, her current  $\text{NPGoals}$  is incremented by one. We discuss the adjustment of  $\text{NPGoals}$  required for subgoal adoption in the next section. Finally,  $\text{NPGoals}$  is not affected by any other action.

Returning to our example, recall that our agent has the c-goals/active p-goals in  $S_0$  that  $\Box\text{BeRich}$  and  $\Box\text{BeHappy}$ , but not  $\Diamond\text{GetPhD}$ , since the latter is inconsistent with her higher priority p-goal  $\Box\text{BeRich}$ . Assume that, after the action  $\text{goBankrupt}$  happens in  $S_0$ , the p-goal  $\Box\text{BeRich}$  becomes impossible. Then in  $S'_1 = \text{do}(\text{goBankrupt}, S_0)$ , the agent has the c-goal that  $\Diamond\text{GetPhD}$ , but not  $\Box\text{BeRich}$  nor  $\Box\text{BeHappy}$ ;  $\Box\text{BeRich}$  is excluded from the set of c-goals since it has become impossible to achieve (i.e. unrealistic). Also, since her higher priority p-goal  $\Diamond\text{GetPhD}$  is inconsistent with her p-goal  $\Box\text{BeHappy}$ , the agent will make  $\Box\text{BeHappy}$  inactive.

Note that, while it might be reasonable to drop a p-goal (e.g.  $\Diamond\text{GetPhD}$ ) that is in conflict with another higher priority active p-goal (e.g.  $\Box\text{BeRich}$ ), in our framework we keep such p-goals around. The reason for this is that although  $\Box\text{BeRich}$  is currently inconsistent with  $\Diamond\text{GetPhD}$ , the agent might later learn that  $\Box\text{BeRich}$  has become impossible to bring about (e.g. after  $\text{goBankrupt}$  occurs), and then might want to pursue  $\Diamond\text{GetPhD}$ . Thus, it is useful to keep these inactive p-goals since this allows the agent to maximize her utility (that of her chosen goals) by taking advantage of such opportunities. As mentioned earlier, c-goals are our analogue to intentions. Recall that Bratman’s model of intentions limits the agent’s practical reasoning – agents do not always optimize their utility and don’t always reconsider all available options in order to allocate their reasoning effort wisely. In contrast to this, our c-goals are defined in terms of the p-goals, and at every step, we ensure that the agent’s c-goals maximizes her utility so that these are the set of highest priority goals that are consistent given the agent’s knowledge. Thus, our notion of c-goals is not as persistent as Bratman’s notion of intention [10]. For instance as mentioned above, after the action  $\text{goBankrupt}$  happens in  $S_0$ , the agent will lose the c-goal that  $\Box\text{BeHappy}$ , although she did not drop it and it did not

<sup>5</sup>  $\text{CorrectHist}(s, \text{path})$  is defined as  $\text{Starts}(\text{path}, s') \wedge \text{SameHistory}(s', s)$ ; here  $\text{path}$  is a placeholder that stands for a path and  $s$  represents the current situation.

become impossible or achieved. In this sense, our model is that of an idealized agent. There is a tradeoff between optimizing the agent's chosen set of prioritized goals and being committed to chosen goals. In our framework, chosen goals behave like intentions with an automatic filter-override mechanism [10] that forces the agent to drop her chosen goals when opportunities to commit to other higher priority goals arise. In the future, it would be interesting to develop a logical model that captures the pragmatics of intention reconsideration by supporting control over it.

We now show that our formalization of prioritized goals has some desirable properties. Some of these (e.g. Proposition 3a) are analogues of the AGM postulates; others (e.g. adopting logically equivalent goals has the same result, etc.) were left out for space reasons. First we show that c-goals are consistent:

**Proposition 1 (Consistency)**

$$D \models \forall s. \neg \text{CGoal}(\text{False}, s).$$

Thus, the agent cannot have both  $\phi$  and  $\neg\phi$  as c-goals in a situation  $s$  and there is a path that is  $G_C$ -accessible in  $s$ . Even if all of the agent's p-goals become known to be impossible, the set of  $G_C$ -accessible paths will be precisely those that starts with a  $K$ -accessible situation, and thus the agent will only choose the propositions that are known to be inevitable.

We also have the property of realism [1], i.e. if an agent knows that something has become inevitable, then she has this as a c-goal:

**Proposition 2 (Realism)**

$$D \models \forall \phi, s. \text{KInevitable}(\phi, s) \supset \text{CGoal}(\phi, s).$$

Note that this is not necessarily true for p-goals and primary c-goals – an agent may know that something has become inevitable and not have it as her p-goal/primary c-goal, which is intuitive. In fact, this is the reason why we define p-goals in terms of  $G$ -accessible paths rather than  $G_R$ . While the property of realism is often criticized, one should view these inevitable goals as something that holds in the worlds that the agent intends to bring about, rather than something that the agent is actively pursuing.

A consequence of Proposition 1 and 2 is that an agent does not have a c-goal that is known to be impossible, i.e.  $D \models \forall \phi, s. \text{CGoal}(\phi, s) \supset \neg \text{KImpossible}(\phi, s)$ .

We next discuss some properties of the framework w.r.t. goal change. Proposition 3 says that (a) an agent acquires the p-goal that  $\phi$  at some level  $n$  after she adopts it, and (b) that she acquires the primary c-goal (and c-goal) that  $\phi$  after she adopts it in  $s$ , provided that she does not have the c-goal in  $s$  that  $\neg\phi$  next.

**Proposition 3 (Adoption)**

- (a)  $D \models \exists n. \text{PGoal}(\phi, n, \text{do}(\text{adopt}(\phi), s))$ ,
- (b)  $D \models \neg \text{CGoal}(\neg \exists s', p'. \text{Starts}(s') \wedge \text{Suffix}(p', \text{do}(\text{adopt}(\phi), s')) \wedge \phi(p'), s) \supset \text{PrimCGoal}(\phi, \text{do}(\text{adopt}(\phi), s))$ .

We can also show that after dropping the p-goal that  $\phi$  at  $n$  in  $s$ , an agent does not have the p-goal (and thus the primary c-goal) that the progression of  $\phi$  at  $n$ , i.e.  $\text{ProgressionOf}(\phi, \text{drop}(\phi), s)$ , provided that  $\text{ProgressionOf}(\phi, \text{drop}(\phi), s)$  is not inevitable in  $\text{do}(\text{drop}(\phi), s)$ .

**Proposition 4 (Drop)**

$$\begin{aligned} D \models & \text{PGoal}(\phi, n, s) \\ & \wedge \neg \text{Inevitable}(\text{ProgressionOf}(\phi, \text{drop}(\phi), s), \text{do}(\text{drop}(\phi), s)) \\ & \supset \neg \text{PGoal}(\text{ProgressionOf}(\phi, \text{drop}(\phi), s), n, \text{do}(\text{drop}(\phi), s)), \end{aligned}$$

where,

$$\text{ProgressionOf}(\phi, a, s) \stackrel{\text{def}}{=} \exists p', s'. \text{Starts}(p', s') \wedge \text{Suffix}(p', \text{do}(a, s')) \wedge \phi(p').$$

Note that, this does not hold for  $\text{CGoal}$ , as  $\phi$  could still be a consequence of two or more of her remaining primary c-goals.

The next few properties concern the persistence of these motivational attitudes. First, we have a persistence property for achievement realistic p-goals:

**Proposition 5 (Persistence of Achievement RPGoals)**

$$D \models \text{RPGoal}(\diamond\Phi, n, s) \wedge \text{Know}(\neg\Phi, s) \wedge \forall\psi. a \neq \text{drop}(\psi) \supset \text{RPGoal}(\diamond\Phi, n, \text{do}(a, s)).$$

This says that if an agent has a realistic p-goal that  $\diamond\Phi$  in  $s$ , then she will retain this realistic p-goal after some action  $a$  has been performed in  $s$ , provided that she knows that  $\Phi$  has not yet been achieved, and  $a$  is not the action of dropping a p-goal. Note that, we do not need to ensure that  $\diamond\Phi$  is still known to be possible or consistent with higher priority active p-goals, since the SSA for  $G$  does not automatically drop such incompatible p-goals from the goal hierarchy.

For achievement chosen goals we have the following:

**Proposition 6 (Persistence of Achievement Chosen Goals)**

$$\begin{aligned} D \models & \text{OPGoal}(\diamond\Phi \wedge \text{CorrectHist}(s), n, s) \wedge \text{CGoal}(\diamond\Phi, s) \\ & \wedge \text{Know}(\neg\Phi, s) \wedge \forall\psi. a \neq \text{drop}(\psi) \wedge \neg \text{CGoal}(\neg\diamond\Phi, n-1, \text{do}(a, s)) \\ & \supset \text{CGoal}(\diamond\Phi, n, \text{do}(a, s)). \end{aligned}$$

Thus, in situation  $s$ , if an agent has the only p-goal at level  $n$  that  $\diamond\Phi$  and the correct history of actions in  $s$  has been performed, and if  $\diamond\Phi$  is also a chosen goal in  $s$  (and thus she has the primary c-goal that  $\diamond\Phi$ ), then she will retain the c-goal that  $\diamond\Phi$  at level  $n$  after some action  $a$  has been performed in  $s$ , provided that:

- she knows in  $s$  that  $\Phi$  has not yet been achieved,
- that  $a$  is not the action of dropping a p-goal,
- and that at level  $n-1$  the agent does not have the c-goal in  $\text{do}(a, s)$  that  $\neg\diamond\Phi$ , i.e.  $\diamond\Phi$  is consistent with higher priority c-goals after  $a$  has been performed in  $s$ .

Note that this property also follows if we replace the consequent with  $\text{CGoal}(\diamond\Phi, \text{do}(a, s))$ , and thus it deals with the persistence of c-goals. Note however that, it does not hold if we replace the  $\text{OPGoal}$  in the antecedent with  $\text{PGoal}$ ; the reason for this is that the agent might have a p-goal at level  $n$  in  $s$  that  $\phi$  and the c-goal in  $s$  that  $\phi$ , but not have  $\phi$  as a primary c-goal in  $s$ , e.g.  $n$  might be an inactive level because another p-goal at  $n$  has become impossible, and  $\phi$  could be a c-goal in  $s$  because it is a consequence of two other primary c-goals. Thus even if  $\neg\phi$  is not a c-goal after  $a$  has been performed in  $s$ , there is no guarantee that the level  $n$  will be active in  $\text{do}(a, s)$  or that all the active p-goals that contributed to  $\phi$  in  $s$  are still active.

We believe that the dropping of an unrelated p-goal will not affect persistence, and hence it should be possible to strengthen Proposition 5 and 6. Also, in the future we would like to generalize these two propositions to deal with arbitrary temporally extended goals.

## 6 Handling Subgoals

In this section, we deal with the dynamics of subgoals. As mentioned earlier, a subgoal must be dropped when the parent goal is dropped or becomes impossible. When adopting a subgoal  $\psi$  with respect to a supergoal  $\phi$ , in addition to recording the newly adopted goal  $\psi$ , we need to model the fact that  $\psi$  is a subgoal of  $\phi$ . This information can later be used to drop the subgoal when the parent goal is dropped. One way of modeling this is to ensure that the adoption of a subgoal  $\psi$  w.r.t. a parent goal  $\phi$  adds new p-goals that contain *both this subgoal and this parent goal i.e.  $\psi \wedge \phi$  at a lower priority than the parent goal  $\phi$* . This ensures that when the parent goal is dropped, the subgoal is also dropped. To see this, recall that to handle the dropping of a goal  $\phi$ , we drop the p-goals at all  $G$ -accessibility levels that imply  $\phi$ . Thus, if we drop the parent goal  $\phi$ , it will also drop all of its subgoals including  $\psi$ , since the  $G$ -accessibility levels where the parent goal  $\phi$  holds include the  $G$ -accessibility levels where the subgoal  $\psi$  holds. Note that, if there are more than one level where the supergoal  $\phi$  is a p-goal, then we copy all these levels, i.e. for each level  $n$  where  $\phi$  is a p-goal, we add a (lower priority) level to the goal hierarchy. As we will see, this ensures that the sub-subgoals and sub-sub-subgoals etc. are also properly dropped when the supergoal is dropped. Also, this means that dropping a subgoal does not necessarily drop the supergoal.

Before going over the formal details, let us mention some useful bookkeeping tools that we will use:  $\text{Length}(l)$  returns the number of elements in list  $l$ ;  $\text{Nth}(l, i)$  returns the  $i$ -th element in list  $l$ , and -1 if  $i > \text{Length}(l)$ ;  $\text{Sort}(l)$  returns a sorted version of list  $l$ . The part of the SSA for  $G$  that handles subgoal adoption is defined as follows:

$$\begin{aligned} \text{SubGoalAdopted}(p, n, a, s, \psi, \phi) \stackrel{\text{def}}{=} & (n < \text{NPGoals}(s) \wedge \text{Progressed}(p, n, a, s)) \vee \\ & (\text{NPGoals}(s) \leq n < \text{NPGoals}(s) + \text{Length}(\text{AddList}(\phi, s)) \\ & \wedge \exists i, m. (n = \text{NPGoals}(s) + i \wedge m = \text{Nth}(\text{AddList}(\phi, s), i) \\ & \wedge \text{Progressed}(p, m, a, s) \wedge \psi(p))) \vee \\ & (n \geq \text{NPGoals}(s) + \text{Length}(\text{AddList}(\phi, s)) \wedge \text{Progressed}(p, n, a, s)). \end{aligned}$$

That is, if the action involves the adoption of a subgoal  $\psi$  w.r.t. a supergoal  $\phi$ , we adjust  $G$  to incorporate (possibly several) new p-goals. We will discuss each case in turn. First,

note that the existing p-goals are just carried over by progressing them; this is handled by the first disjunct.

Secondly, we adjust  $G$  starting at level  $NPGoals(s)$ . We add a number of new levels that include the conjunction of the only p-goal and the subgoal at a lower priority for all the current only p-goals that imply the parent goal  $\phi$ . For example, say at level  $i$  we have an OPGoal that  $\phi_i$  and it implies the parent goal that  $\phi$ ; then we add at a lower priority the conjoined goal of the progressed version of  $\phi_i$  and the subgoal  $\psi$ . Our formalization of this uses the abbreviation  $AddList(\phi, s)$  which is a sorted list of levels such that the parent goal is implied by the only p-goal at this level.  $AddList$  is defined as:  $AddList(\phi, s) \stackrel{\text{def}}{=} \text{Sort}([n \mid PGoal(\phi, n, s)])$ . The length of this list indicates the number of lower priority goals that needs to be added. As discussed above, this ensures that the agent will drop the subgoal when the parent goal is dropped (but not necessarily vice-versa). Note that if this process adds two or more new p-goals to the agent's goal hierarchy, we maintain the original ordering; e.g. suppose that the agent adopted  $\psi$  w.r.t.  $\phi$ , that there are two  $G$ -accessibility levels  $m$  and  $n$  such that the agent has the only p-goal that  $\phi_m$  at  $m$  and  $\phi_n$  at  $n$ , that  $\phi_m$  implies  $\phi$  and  $\phi_n$  implies  $\phi$ , and that  $n > m$ . In that case, the SSA for  $G$  will add the p-goal  $\phi_m \wedge \phi$  at level  $NPGoals(s)$  and the p-goal  $\phi_n \wedge \phi$  at level  $NPGoals(s) + 1$ .

Finally, all the remaining levels involving trivially true goals are just carried over by progressing them.

The part of the SSA for  $NPGoals$  that handles subgoal adoption is defined as follows:

$$\text{AdjustSubGoalAdopt}(\phi, s) \stackrel{\text{def}}{=} NPGoals(s) + \text{Length}(AddList(\phi, s)).$$

That is, when the agent adopts a subgoal w.r.t. a parent goal, her current  $NPGoals$  is incremented by the number of new p-goals adopted in this process.

Let us go back to our example. Suppose that the agent knows that one way of always being rich is to always work hard, which in turns can be fulfilled by always being energetic. Assume that with this in mind, our agent adopts the subgoal that  $\Box\text{WorkHard}$  w.r.t. the p-goal that  $\Box\text{BeRich}$ , and then adopts the sub-subgoal that  $\Box\text{BeEnergetic}$  w.r.t. the subgoal that  $\Box\text{WorkHard}$  starting in  $S_0$ . Then the agent's goal hierarchy in  $S_1 = do(adopt(\Box\text{WorkHard}, \Box\text{BeRich}), S_0)$  should include the p-goal that  $\Box\text{WorkHard}$  and in  $S_2 = do(adopt(\Box\text{BeEnergetic}, \Box\text{WorkHard}), S_1)$  should also include the p-goal that  $\Box\text{BeEnergetic}$ . According to the SSA for  $G$ , our agent's goal hierarchy in  $S_1$  and in  $S_2$  will be as in Table 1.<sup>6</sup> In  $S_0$ , the supergoal  $\Box\text{BeRich}$  holds at level 0 and thus  $AddList(\Box\text{BeRich}, S_0) = [0]$ . Similarly in  $S_1$ , the supergoal  $\Box\text{WorkHard}$  holds at level 3 and thus  $AddList(\Box\text{WorkHard}, S_1) = [3]$ . Now, suppose that in  $S_2$  the agent wants to drop the p-goal that  $\Box\text{WorkHard}$ . Then in  $S_3 = do(drop(\Box\text{WorkHard}), S_2)$ , she should no longer have  $\Box\text{BeEnergetic}$  as a p-goal, but should retain the supergoal that  $\Box\text{BeRich}$ . After the agent drops the p-goal that  $\Box\text{WorkHard}$ , by the SSA for  $G$  we can see that all the  $G$ -accessible levels where  $\Box\text{WorkHard}$  holds will be replaced by the only p-goal that  $\text{CorrectHist}(S_2, path)$  (see  $S_3$  in Table 1). This shows that dropping  $\Box\text{WorkHard}$  results in the dropping of all of its subgoals (in this case  $\Box\text{BeEnergetic}$ ), but that its parent goal  $\Box\text{BeRich}$  is retained.

<sup>6</sup> For simplicity in Table 1, we only show the agent's relevant p-goals rather than its only p-goals (which in addition reflect the actions that have been performed so far, i.e.  $\text{CorrectHist}(s)$ ).

We define the SubGoal relation as follows:

$$\text{SubGoal}(\psi, \phi, s) \stackrel{\text{def}}{=} \exists n. \text{PGoal}(\phi, n, s) \wedge \neg \text{PGoal}(\psi, n, s) \\ \wedge \forall n. \text{PGoal}(\psi, n, s) \supset \text{PGoal}(\phi, n, s).$$

This says that  $\psi$  is a subgoal of  $\phi$  in situation  $s$  iff there exists an  $G$ -accessibility level  $n$  in  $s$  such that  $\phi$  is a p-goal at  $n$  while  $\psi$  is not, and for all  $G$ -accessibility levels in  $s$  where  $\psi$  is a p-goal,  $\phi$  is also a p-goal. Note that, while our formalization of subgoal dynamics allows a subgoal to have multiple parents, in this definition we assume that a subgoal can't have more than one parent. In the future, we will work on relaxing this constraint.

We now discuss some properties concerning the dynamics of subgoals and the dependencies between a subgoal and its parent goal. Proposition 7 states that (a) an agent acquires the p-goal that  $\psi$  after she adopts it as a subgoal of another goal  $\phi$  in  $s$ , provided that she has the p-goal at some level in  $s$  that  $\phi$ , and (b) she also acquires the primary c-goal that  $\psi$  after she adopts it as a subgoal of  $\phi$  in  $s$ , provided that she has the primary c-goal in  $s$  that  $\phi$ , and that she does not have the c-goal in  $s$  that  $\neg\psi$  next.

**Proposition 7 (Subgoal Adoption)**

- (a)  $D \models \exists m. \text{PGoal}(\phi, m, s) \supset \exists n. \text{PGoal}(\psi, n, \text{do}(\text{adopt}(\psi, \phi), s))$ ,
- (b)  $D \models \text{PrimCGoal}(\phi, s) \\ \wedge \neg \text{CGoal}(\neg \exists s', p'. \text{Starts}(s') \wedge \text{Suffix}(p', \text{do}(\text{adopt}(\psi, \phi), s')) \wedge \psi(p'), s) \\ \supset \text{PrimCGoal}(\psi, \text{do}(\text{adopt}(\psi, \phi), s))$ .

The next property says that after dropping the p-goal that  $\phi$  in  $s$ , an agent does not have the p-goal (and thus the primary c-goal) that the progression of  $\psi$ , provided that  $\psi$  is a subgoal of  $\phi$  in  $s$ , and that the progression of  $\psi$  is not inevitable in  $\text{do}(\text{drop}(\phi), s)$ .

**Proposition 8 (Supergoal Drop)**

$$D \models \text{SubGoal}(\psi, \phi, s) \wedge \neg \text{Inevitable}(\text{ProgressionOf}(\psi, \text{drop}(\phi), s), \text{do}(\text{drop}(\phi), s)) \\ \supset \neg \exists n. \text{PGoal}(\text{ProgressionOf}(\psi, \text{drop}(\phi), s), n, \text{do}(\text{drop}(\phi), s)).$$

As with Proposition 4, this does not hold if we replace PGoal in the consequence with CGoal since  $\psi$  could be a consequence of a combination of other active p-goals.

The next two properties say that dropping a subgoal does not effect the parent goal.

**Proposition 9 (Subgoal Drop)**

- (a)  $D \models \text{SubGoal}(\psi, \phi, s) \\ \supset \exists n. \text{PGoal}(\text{ProgressionOf}(\phi, \text{drop}(\psi), s), n, \text{do}(\text{drop}(\psi), s))$ ,
- (b)  $D \models \text{SubGoal}(\psi, \phi, s) \wedge \text{PrimCGoal}(\phi, s) \\ \supset \text{PrimCGoal}(\text{ProgressionOf}(\phi, \text{drop}(\psi), s), \text{do}(\text{drop}(\psi), s))$ .

That is, (a) an agent retains the p-goal that the progression of  $\phi$  after she drops a subgoal  $\psi$  of  $\phi$ , and (b) she also retains the primary c-goal that the progression of  $\phi$  after she drops a subgoal  $\psi$  of  $\phi$  in  $s$ , provided that she has the primary c-goal that  $\phi$  in  $s$ .

Finally, it can be shown that the SubGoal relation is transitive, i.e. if  $\psi_1$  is a subgoal of  $\phi$  in  $s$ , and if  $\psi_2$  is a subgoal of  $\psi_1$  in  $s$ , then  $\psi_2$  must also be a subgoal of  $\phi$  in  $s$ .

## 7 Discussion and Future Work

In this paper, we presented a formalization of prioritized goals, subgoals, and their dynamics. Our formalization ensures that an agent's chosen goals are always consistent and that goals and subgoals properly evolve as a result of regular actions as well as of adopting and dropping goals. Although we made some simplifying assumptions, in this paper we have focused on developing an expressive framework that captures an idealized form of rationality without worrying about tractability. It would be desirable to study restricted fragments of the logic where reasoning is tractable. Also, before defining more limited forms of rationality, one should have a clear specification of what ideal rationality really is so that one understands what compromises are being made.

While in our account chosen goals are closed under logical consequence, primary c-goals are not. Thus, our formalization of primary c-goals is related to the non-normal modal formalizations of intentions found in the literature [3], and as such it does not suffer from the side-effect problem [1]. For instance, in our framework an agent can have the primary c-goal to get her teeth fixed and know that this always involves pain, but not have the primary c-goal to have pain.

Also, since we are using the situation calculus, we can easily represent procedural goals/plans, e.g. the goal to do  $a_1$  and then  $a_2$  can be written as:  $\text{PGoal}(\exists s, s_1. \text{Starts}(s_1) \wedge \text{OnPath}(s) \wedge s = \text{do}(a_2, \text{do}(a_1, s_1)), 0, S_0)$ . Golog [12] can be used to represent complex plans/programs. So we can model the adoption of plans as subgoals.

Recently, there have been a few proposals that deal with goal change. Shapiro *et al.* [18] present a situation calculus based framework where an agent adopts a goal when she is requested to do so, and remains committed to this goal unless the requester cancels this request; a goal is retained even if the agent learns that it has become impossible, and in this case the agent's goals become inconsistent. Shapiro and Brewka [7] modify this framework to ensure that goals are dropped when they are believed to be impossible or when they are achieved. Their account is similar to ours in the sense that they also assume a priority ordering over the set of (in their case, requested) goals, and in every situation they compute chosen goals by computing a maximal consistent goal set that is also compatible with the agent's beliefs. However, their model has some unintuitive properties: the agent's chosen goals in  $\text{do}(a, s)$  may be quite different from her chosen goals in  $s$ , although  $a$  did not make any of her goals in  $s$  impossible or inconsistent with higher priority goals, because inconsistencies between goals at the same priority level are resolved differently. In their framework, this can happen because goals are only partially ordered. Note that, while one might argue that a partial order over goals might be more general, allowing this means that additional control information is required to obtain a single goal state after the agent's goals change. Also, we provide a more expressive formalization of prioritized goals – we model goals using infinite paths, and thus can model many types of goals that they cannot. Finally they model prioritized goals by treating the agent's p-goals as an arbitrary set of temporal formulae, and then defining the set of c-goals as a subset of the p-goals. But our possible world semantics has some advantages over this: it clearly defines when goals are consistent with each other and with what is known. One can easily specify how goals change when an action  $a$  occurs, e.g. the goal to do  $a$  next and then do  $b$  becomes the goal to do  $b$  next, the goal that  $\diamond\Phi \vee \diamond\Psi$  becomes the goal that  $\diamond\Psi$  if  $a$  makes achieving  $\Phi$  impossible, etc.



There has been much work on agent programming languages with declarative goals where the dynamics of goals and intentions and the dependencies between goals and subgoals are modeled (e.g. [19,20,9] and the references therein). However, most of these are not based on a formal theory of agency, and to the best of our knowledge, none maintains the consistency of (chosen) goals (e.g. when adopting a plan to achieve a goal, these frameworks do not ensure that this plan is consistent with the agent's other concurrent goals/plans). Also, most of these do not deal with temporally extended goals, and as a result they often need to accommodate inconsistent goal-bases to allow the agent to achieve conflicting states at different time points (e.g. the default logic based framework in [21]); chosen goals are required to be consistent. In [22], the authors formalized two semantics for representing conflicting goals, using propositional and default logic; they argued that even logically consistent goals can be conflicting, e.g. when multiple goals/plans are chosen to fulfill the same (super)goal. Unlike us however, they do not address how an agent chooses the goals that she will actively pursue. In [6], the authors present a situation calculus based agent programming language where the agent executes a program while maximizing the achievement of a set of prioritized goals. However, they do not formalize goal dynamics.

One limitation of our account is that we assume that the agent's p-goals are totally ordered in terms of priority. Also, newly adopted p-goals are assigned the lowest priority. A consequence of this is that an agent's c-goals depend on the adoption order of her p-goals. For instance, given a fixed starting situation, an agent can end up with two different sets of c-goals by adopting  $\phi$  followed by  $\psi$ , and by adopting  $\psi$  followed by  $\phi$ . This has very different results when  $\phi$  and  $\psi$  conflict with each other. We would like to address this by incorporating the priority of the p-goal as an argument to the *adopt* action, and handling this in the framework. Finally, one could argue that our agent wastes resources trying to optimize her c-goals at every step. In the future, we would like to develop an account where the agent is strongly committed to her chosen goals, and where the filter override mechanism is only triggered under specific conditions.

## References

1. Cohen, P.R., Levesque, H.J.: Intention is Choice with Commitment. *Artificial Intelligence* 42(2-3), 213–361 (1990)
2. Sadek, M.D.: A Study in the Logic of Intention. In: *Third Intl. Conf. on Principles of Knowledge Representation and Reasoning (KR&R 1992)*, Cambridge, MA, pp. 462–473 (1992)
3. Konolige, K., Pollack, M.E.: A Representationalist Theory of Intention. In: *Thirteenth Intl. J. Conf. on Artificial Intelligence (IJCAI 1993)*, Chambéry, France, pp. 390–395 (1993)
4. Singh, M.P.: *Multiagent Systems – A Theoretical Framework for Intentions, Know-How, and Communications*. LNCS (LNAI), vol. 799. Springer, Heidelberg (1994)
5. Shapiro, S., Lespérance, Y., Levesque, H.J.: Goals and Rational Action in the Situation Calculus - A Preliminary Report. In: *Working Notes of the AAAI Fall Symposium on Rational Agency: Concepts, Theories, Models, and Applications*, Cambridge, MA, November 1995, pp. 117–122 (1995)
6. Sardina, S., Shapiro, S.: Rational Action in Agent Programs with Prioritized Goals. In: *Second Intl. J. Conf. on Autonomous Agents and Multi-Agent Sys (AAMAS 2003)*, Melbourne, Australia, pp. 417–424 (2003)
7. Shapiro, S., Brewka, G.: Dynamic Interactions Between Goals and Beliefs. In: *Twentieth Intl. J. Conf. on Artificial Intelligence (IJCAI 2007)*, India, pp. 2625–2630 (2007)

8. Winikoff, M., Padgham, L., Harland, J., Thangarajah, J.: Declarative and Procedural Goals in Intelligent Agent Systems. In: Eighth Intl. Conf. on Principles and Knowledge Representation and Reasoning (KR&R 2002), Toulouse, France, pp. 470–481 (2002)
9. Bordini, R.H., Dastani, M., Dix, J., Fallah-Seghrouchni, A.E. (eds.): *Multi-Agent Programming: Languages, Platforms and Applications*. Springer, Heidelberg (2005)
10. Bratman, M.E.: *Intentions, Plans, and Practical Reason*. Harvard University Press, Cambridge (1987)
11. McCarthy, J., Hayes, P.J.: Some Philosophical Problems from the Standpoint of Artificial Intelligence. *Machine Intelligence* 4, 463–502 (1969)
12. Reiter, R.: *Knowledge in Action. Logical Foundations for Specifying and Implementing Dynamical Systems*. MIT Press, Cambridge (2001)
13. DeGiacomo, G., Lespérance, Y., Levesque, H.J.: ConGolog, a Concurrent Programming Language Based on the Situation Calculus. *Artificial Intelligence* 121, 109–169 (2000)
14. Levesque, H.J., Pirri, F., Reiter, R.: Foundations for a Calculus of Situations. *Electronic Transactions of AI (ETAI)* 2(3-4), 159–178 (1998)
15. Moore, R.C.: A Formal Theory of Knowledge and Action. In: Hobbs, J.R., Moore, R.C. (eds.) *Formal Theories of the Commonsense World*, pp. 319–358. Ablex, Greenwich (1985)
16. Scherl, R., Levesque, H.: Knowledge, Action, and the Frame Problem. *Artificial Intelligence* 144(1-2), 1–39 (2003)
17. Rao, A.S., Georgeff, M.P.: Modeling Rational Agents with a BDI-Architecture. In: Fikes, R., Sandewall, E. (eds.) *Second Intl. Conf. on Principles of Knowledge Representation and Reasoning (KR&R 1991)*, San Mateo, CA, pp. 473–484. Morgan Kaufmann Publishers, San Francisco (1991)
18. Shapiro, S., Lespérance, Y., Levesque, H.J.: Goal Change in the Situation Calculus. *J. of Logic and Computation* 17(5), 983–1018 (2007)
19. Sardina, S., deSilva, L., Padgham, L.: Hierarchical Planning in BDI Agent Programming Languages: A Formal Approach. In: Fifth Intl. J. Conf. on Autonomous Agents and Multi-Agent Systems (AAMAS 2006), Hakodate, Japan, pp. 1001–1008 (2006)
20. van Riemsdijk, M.B., Dastani, M., Dignum, F., Meyer, J.J.C.: Dynamics of Declarative Goals in Agent Programming. In: Leite, J., Omicini, A., Torroni, P., Yolum, p. (eds.) *DALT 2004. LNCS (LNAI)*, vol. 3476, pp. 1–18. Springer, Heidelberg (2005)
21. van Riemsdijk, M.B., Dastani, M., Meyer, J.J.: Ch.: Semantics of Declarative Goals in Agent Programming. In: Fourth Int'l J. Conf. on Autonomous Agents and Multiagent Sys. (AAMAS 2005), pp. 133–140 (2005)
22. van Riemsdijk, M.B., Dastani, M., Meyer, J.J.Ch.: Goals in Conflict: Semantic Foundations of Goals in Agent Programming. *International Journal of Autonomous Agents and Multi-Agent Systems (JAAMAS)* 18(3), 471–500 (2009)