

# From Actions to Programs as Abstract Actual Causes

Bitá Banihashemi<sup>\*1</sup>, Shakil M. Khan<sup>\*2</sup>, Mikhail Soutchanski<sup>3</sup>

<sup>1</sup> York University, 4700 Keele St, Toronto, ON M3J 1P3, Canada

<sup>2</sup> University of Regina, 3737 Wascana Pkwy, Regina, SK S4S 0A2, Canada

<sup>3</sup> (Former) Ryerson University, 350 Victoria St, Toronto, ON M5B 2K3, Canada

bita@eecs.yorku.ca, shakil.khan@uregina.ca, mes@cs.ryerson.ca

## Abstract

Causality plays a central role in reasoning about observations. In many cases, it might be useful to define the conditions under which a non-deterministic program can be called an actual cause of an effect in a setting where a sequence of programs are executed one after another. There can be two perspectives, one where at least one execution of the program leads to the effect, and another where all executions do so. The former captures a “weak” notion of causation and is more general than the latter stronger notion. In this paper, we give a definition of weak potential causes. Our analysis is performed within the situation calculus basic action theories and we consider programs formulated in the logic programming language ConGolog. Within this setting, we show how one can utilize a recently developed abstraction framework to relate causes at various levels of abstraction, which facilitates reasoning about programs as causes.

## Introduction

Actual or token causation is concerned with identifying the events or actions in a trace that can be considered as causes of an observed effect. The seminal work of Pearl (2000) provided the foundations and served as inspiration for research on actual causes in AI. This research culminated in the book (Halpern 2016) that summarized a number of previously developed definitions concerning when an event can be considered as an actual cause of an effect. These definitions are developed within the framework of structural equations models (SEM), where a simple event is understood as assigning a value to an endogenous variable.

However, this perspective does not facilitate the study of causation for more complex activities such as control flow in programs. It can be interesting and important to define when a non-deterministic program is an actual cause of an effect in a setting where a sequence of programs are executed one after another. This immediately leads to the question when can one intuitively say that a program is an actual cause? One perspective can be that a non-deterministic program is a *weak potential cause*, if at least one execution of the program leads to a situation where the effect holds. Another perspective is that a program is a *strong potential cause* if

all executions of the program produce the effect. Note that a strong potential cause is also weak, but not vice versa. Also, in both cases we talk about *potential* causes, since they can manifest only in some of the situations that are produced by the execution of the program sequence.

As an example, imagine *Suzy* buying a lottery ticket that later wins a reward. If one conceptualizes the complex actions of purchasing the ticket as a highly non-deterministic program, then it is reasonable to say that this program was a weak potential cause of the fact that *Suzy* won, since there is an execution of this program that leads to a situation where the effect holds and another to a situation where it doesn't.

Again, imagine a computer system that involves multiple interacting agents. The typical examples of such systems arise in computer security contexts where the behaviours of the agents are specified by non-deterministic protocols due to versatility of possible agent interactions. In this context, one might be interested in determining if all of the executions of a protocol led to the successful handling of a security leak. This corresponds to the case of a strong potential cause.

In this paper, we give a definition of the more inclusive notion of weak cause. We consider programs formulated in the high-level logic programming language ConGolog (De Giacomo, Lespérance, and Levesque 2000), which is based on action theories specified in the situation calculus (SC) (McCarthy and Hayes 1969; Reiter 2001). We build on a recently proposed definition of actual cause in the SC (Batusov and Soutchanski 2018), which only considers primitive actions as causes. Since we focus on programs as causes, a natural question that arises then is how these two notions can be related. The programs can be complex, but often they can be conceptualized at some abstract high-level (HL) as actions. It turns out that the abstraction framework proposed in (Banihashemi, De Giacomo, and Lespérance 2017) that can relate programs with primitive actions is also useful for relating a subclass of weak potential causes (in particular, weak causes that are also strong) at different levels of abstraction.

On the semantic level, models of programs can be very complicated, but reasoning about effects of actions that serve as their abstractions can be easier since essential details are encapsulated in a simpler HL model. We argue that HL and low-level (LL) causes can be related in a kind of commutative diagram. Namely, if an HL action is found to be a cause

<sup>\*</sup>These authors contributed equally.

of an effect, this action is associated to a program  $\delta$  defined over an LL theory that implements it, and this effect is an abstraction of an LL formula  $\phi$  (i.e.,  $\phi$  is a refinement of the effect), then at the LL,  $\delta$  must be a cause of  $\phi$ . This result is one of our main contributions. We focus here on semantics and leave computational issues to future work.

The rest of this paper is organized as follows. In the next section, we outline the SC. Then, we discuss previous work on actual cause and abstraction. Subsequently, we define what it means for a program to be a cause. Finally, we study how abstraction can be utilized to reason about abstract causes. We conclude with some discussion of previous work and avenues for future research.

## Preliminaries

Our base framework for this is the situation calculus (SC) (McCarthy and Hayes 1969) as formalized in (Reiter 2001). We assume that there is a *finite number of action types*  $\mathcal{A}$ . Moreover, we assume that the terms of object sort are a countably infinite set  $\mathcal{N}$  of standard names for which we have the unique names assumption and domain closure. For simplicity, and w.l.o.g., we assume that there are no functions other than constants and no non-fluent predicates.

A basic action theory (BAT)  $\mathcal{D}$  is the union of the following disjoint sets: the foundational, domain independent, (second-order, or SO) axioms of the SC; (first-order, or FO) precondition axioms; (FO) successor state axioms (SSAs) describing how fluents change between situations; (FO) unique names axioms for actions and (FO) domain closure on action types; (SO) unique name axioms and domain closure for object constants; and (FO) axioms describing the initial configuration of the world. A special predicate  $Poss(a, s)$  is used to state that action  $a$  is executable in situation  $s$ ; precondition axioms characterize this predicate. The abbreviation  $Executable(s)$  means that every action performed in reaching situation  $s$  was possible in the situation in which it occurred. The binary relation  $\sqsubset$  defines precedence on situations; thus  $s \sqsubset s'$  indicates  $s$  is a sub-history of  $s'$ . Note that  $s \sqsubseteq s'$  is an abbreviation for  $s \sqsubset s' \vee s = s'$ . Also,  $s \leq s'$  states that  $s'$  is a successor situation of  $s$  and that every action between  $s$  and  $s'$  is in fact executable. We write  $do([a_1, a_2, \dots, a_{n-1}, a_n], s)$  as an abbreviation for  $do(a_n, do(a_{n-1}, \dots, do(a_2, do(a_1, s)) \dots))$ ; for an action sequence  $\vec{a}$ , we often write  $do(\vec{a}, s)$  for  $do([\vec{a}], s)$ .

An SC formula is *uniform* in  $s$  iff it does not mention  $Poss$ ,  $\sqsubset$ , or equality on situations, it does not quantify over situations, and whenever it mentions a term of sort situation then that term is  $s$ . Also, we use upper-case Greek letters for situation-suppressed SC formulae and we denote by  $\Phi[s]$  the formula obtained from  $\Phi$  by restoring the situation argument  $s$  into all fluents in  $\Phi$ . To represent and reason about complex actions, various *high-level programming languages* have been defined. Here we concentrate on (a fragment of) ConGolog (De Giacomo, Lespérance, and Levesque 2000) that includes the following constructs:

$$\delta ::= nil \mid \alpha \mid \Phi? \mid (\delta_1; \delta_2) \mid (\delta_1 | \delta_2) \mid (\pi x. \delta(x)) \mid \delta^* \mid (\delta_1 || \delta_2).$$

Thus, complex actions can be composed using constructs

that include the *empty program* ( $nil$ ), primitive actions ( $\alpha$ ), waiting for a condition ( $\Phi?$ ), sequence ( $\delta_1; \delta_2$ ), nondeterministic branch ( $\delta_1 | \delta_2$ ), nondeterministic choice of arguments ( $\pi x. \delta(x)$ ), nondeterministic iteration ( $\delta^*$ ), and concurrent execution ( $\delta_1 || \delta_2$ ). Intuitively,  $\pi x. \delta(x)$  nondeterministically picks a binding for the variable  $x$  and performs the program  $\delta$  for this binding of  $x$ .

The semantics of ConGolog is specified in terms of single-step transitions, using the following two predicates (De Giacomo, Lespérance, and Levesque 2000): (i)  $Trans(\delta, s, \delta', s')$ , which holds if one step of program  $\delta$  in situation  $s$  may lead to situation  $s'$  with  $\delta'$  remaining; and (ii)  $Final(\delta, s)$ , which holds if program  $\delta$  may legally terminate in situation  $s$ . The definitions of  $Trans$  and  $Final$  we use are as in (De Giacomo, Lespérance, and Levesque 2000), except that the test construct  $\Phi?$  does not yield any transition, but is final when satisfied. The predicate  $Do(\delta, s, s')$  means that program  $\delta$ , when executed starting in situation  $s$ , has  $s'$  as its legal terminating situation. It is defined as  $Do(\delta, s, s') \doteq \exists \delta'. Trans^*(\delta, s, \delta', s') \wedge Final(\delta', s')$  where  $Trans^*$  denotes the reflexive transitive closure of  $Trans$ . We use  $\mathcal{K}$  to denote the axioms defining ConGolog.

Following (De Giacomo, Lespérance, and Muise 2012), we say that a ConGolog program  $\delta$  is *situation-determined* (SD) in  $s$  if for every sequence of transitions, the remaining program is determined by the resulting situation, i.e.,

$$SituationDetermined(\delta, s) \doteq \forall s', \delta', \delta''. Trans^*(\delta, s, \delta', s') \wedge Trans^*(\delta, s, \delta'', s') \supset \delta' = \delta''.$$

**Example.** Our running example involves a simple rescue robot  $Rob$  that is designed to aid first responders. Initially  $Rob$  is at the *Station* but as an emergency at location  $L1$  exists,  $Rob$  is expected to go to  $L1$  and assist in the rescue operations (by removing rubble or by evacuating people). Action  $go_{LL}(r, l)$  takes robot  $r$  to location  $l$ , and is executable if  $r$  is not already at that location. Action  $removeRubble(r, l)$  (resp.  $evacuate(r, l)$ ) can be performed by robot  $r$  at location  $l$  to remove rubble (resp. evacuate people); these actions are executable if  $r$  is at location  $l$ . Fluent  $At_{LL}(r, l, s)$  indicates  $r$ 's location to be  $l$  at situation  $s$ . Fluents  $Cleared(r, l, s)$  and  $Evacuated(r, l, s)$  evaluate to true when the robot  $r$  has removed rubble and evacuated people at location  $l$  respectively. Robot  $r$  is also able to update the software packages it uses by performing action  $updateSW_{LL}(r, v)$ , where  $v$  indicates the version of the software. Fluent  $UpdatedSW_{LL}(r, v, s)$  indicates if software has been updated to version  $v$ . Initially, we assume a new version  $V2021$  is available.

The BAT for this domain  $\mathcal{D}_i^{ex}$  includes the following action precondition axioms (throughout, we assume that free variables are universally quantified from the outside):

$$\begin{aligned} Poss(go_{LL}(r, l), s) &\equiv \neg At_{LL}(r, l, s), \\ Poss(updateSW_{LL}(r, v), s) &\equiv \neg UpdatedSW_{LL}(r, v, s), \\ Poss(removeRubble(r, l), s) &\equiv At_{LL}(r, l, s), \\ Poss(evacuate(r, l), s) &\equiv At_{LL}(r, l, s). \end{aligned}$$

Moreover,  $\mathcal{D}_i^{ex}$  includes the following SSAs:

$$\begin{aligned} UpdatedSW_{LL}(r, v, do(a, s)) &\equiv \\ a = updateSW_{LL}(r, v) \vee UpdatedSW_{LL}(r, v, s), \end{aligned}$$

$$\begin{aligned}
At_{LL}(r, l, do(a, s)) &\equiv a = go_{LL}(r, l) \vee \\
&(At_{LL}(r, l, s) \wedge \neg \exists l'. l' \neq l \wedge a = go_{LL}(r, l')), \\
Cleared(r, l, do(a, s)) &\equiv \\
a = removeRubble(r, l) &\vee Cleared(r, l, s), \\
Evacuated(r, l, do(a, s)) &\equiv \\
a = evacuate(r, l) &\vee Evacuated(r, l, s).
\end{aligned}$$

Thus, e.g.,  $r$  will be located at  $l$  in  $do(a, s)$  iff  $a$  refers to  $r$  going to  $l$ , or  $r$  was already at  $l$  in  $s$  and  $a$  is not the action of  $r$  going to a different location  $l'$ .

$\mathcal{D}_l^{ex}$  also includes the following initial state axioms:  
 $At_{LL}(Rob, Station, S_0), \neg UpdatedSW_{LL}(Rob, V2021, S_0),$   
 $\neg Evacuated(Rob, L1, S_0), \neg Cleared(Rob, L1, S_0).$   $\triangleleft$

## Theoretical Foundations

### Actual Cause

Given a trace of events, *actual achievement causes* are the events that are behind achieving an effect. In this section, we review previous work on achievement causality in the SC (Batusov and Soutchanski 2018). An effect here is an SC formula  $\Phi[s]$  that is *uniform in  $s$*  and that may include quantifiers over object variables. Given an effect  $\Phi$ , the actual causes are defined relative to a *causal setting* that includes a BAT  $\mathcal{D}$  representing the domain dynamics, and a ground situation  $\sigma$ , representing the “narrative” (i.e., trace of events) where the effect was observed.

**Definition 1 (Causal Setting)** *A causal setting is a tuple  $\langle \mathcal{D}, \sigma, \Phi[s] \rangle$ , where  $\mathcal{D}$  is a BAT,  $\sigma$  is a ground situation term of the form  $do([\alpha_1, \dots, \alpha_n], S_0)$  with ground action functions  $\alpha_1, \dots, \alpha_n$  such that  $\mathcal{D} \models Executable(\sigma)$ , and  $\Phi[s]$  is an SC formula uniform in  $s$  such that  $\mathcal{D} \models \neg \Phi[S_0] \wedge \Phi[\sigma]$ .*

Since the theory  $\mathcal{D}$  does not change, when referring to a causal setting we will often suppress  $\mathcal{D}$  and simply write  $\langle \sigma, \Phi \rangle$ . Also, here  $\Phi$  is required to hold by the end of the narrative  $\sigma$ , and thus we ignore the cases where  $\Phi$  is not achieved by the actions in  $\sigma$ , since in that case, the achievement cause truly does not exist.

As all changes in the SC result from actions, the achievement causes of an effect are contained within a set of ground action terms occurring in  $\sigma$ . However, since  $\sigma$  might include multiple occurrences of the same action, one also needs to identify the situations where those actions were executed.

According to (Batusov and Soutchanski 2018), if some action  $\alpha$  of the action sequence in  $\sigma$  triggers the formula  $\Phi$  to change its truth value from false to true relative to  $\mathcal{D}$ , and if there are no actions in  $\sigma$  after  $\alpha$  that change the value of  $\Phi$  back to false, then  $\alpha$  is an actual cause of achieving  $\Phi$  in  $\sigma$ . They showed that using the single-step regression operator  $\rho$  (i.e., one-step version of the regression operator defined in (Reiter 2001)), in addition to the primary action that actually brings about the effect of interest, one can recursively compute the chain of actions that build up to the primary achievement cause. The following inductive definition formalizes this intuition. Let  $\Pi_{apa}(\alpha, \sigma)$  be the r.h.s. of the precondition axiom for  $\alpha$  in  $\sigma$ .

**Definition 2 (Achievement Cause)** *A causal setting  $\mathcal{C} = \langle \mathcal{D}, \sigma, \Phi[s] \rangle$  satisfies the achievement condition of  $\Phi$  via the*

*situation term  $do(\alpha^*, \sigma^*) \sqsubseteq \sigma$  iff there is an action  $\alpha'$  and situation  $\sigma'$  such that*

$$\mathcal{D} \models \neg \Phi[\sigma'] \wedge \forall s. do(\alpha', \sigma') \sqsubseteq s \sqsubseteq \sigma \supset \Phi[s],$$

*and either  $\alpha^* = \alpha'$  and  $\sigma^* = \sigma'$ , or the associated causal setting  $\langle \sigma', \rho[\Phi[s], \alpha'] \wedge \Pi_{apa}(\alpha', \sigma') \rangle$  satisfies its achievement condition via the situation term  $do(\alpha^*, \sigma^*)$ . Whenever a causal setting  $\mathcal{C}$  satisfies the achievement condition via situation  $do(\alpha^*, \sigma^*)$ , the action  $\alpha^*$  executed in situation  $\sigma^*$  is said to be an achievement cause of  $\mathcal{C}$ .*

According to (Batusov and Soutchanski 2018), the achievement causes of  $\mathcal{C}$  form a finite sequence of situation-action pairs, which is called the *achievement causal chain* of  $\mathcal{C}$ .

**Example (Cont’d).** Consider causal setting  $\mathcal{C}_{ex} = \langle \mathcal{D}_l^{ex}, \sigma_{ex1}, \Phi_{ex1} \rangle$ , where  $\Phi_{ex1} = \exists r, l. Cleared(r, l)$  and  $\sigma_{ex1} = do([updateSW_{LL}(Rob, V2021), go_{LL}(Rob, L1), removeRubble(Rob, L1)], S_0)$ . Then by Definition 2, the action  $removeRubble(Rob, L1)$  performed in situation  $S_2 = do([updateSW_{LL}(Rob, V2021), go_{LL}(Rob, L1)], S_0)$  is an achievement cause of  $\mathcal{C}_{ex}$ . This is the case since  $removeRubble(Rob, L1)$  is the first action after which the effect  $\Phi_{ex1}$  becomes true. Moreover, we can show that  $go_{LL}(Rob, L1)$  executed in  $S_1 = do(updateSW_{LL}(Rob, V2021), S_0)$  is another achievement cause of  $\mathcal{C}_{ex}$ , since the causal setting  $\langle \mathcal{D}_l^{ex}, \Phi', S_2 \rangle$  satisfies the achievement condition  $\Phi'$  via the situation term  $do(go_{LL}(Rob, L1), S_1)$ , where  $\Phi' = \rho[\Phi_{ex1}, removeRubble(Rob, L1)] \wedge \Pi_{apa}(removeRubble(Rob, L1), S_2)$ . Finally, these are all the causes, and in particular  $updateSW_{LL}(Rob, V2021)$  executed in  $S_0$  is not an achievement cause of  $\mathcal{C}_{ex}$ .  $\triangleleft$

### Abstraction

We will use the abstraction framework of (Banihashemi, De Giacomo, and Lespérance 2017) for reasoning about abstract causes. In this framework, there is a high-level (HL) or abstract action theory  $\mathcal{D}_h$  and a low-level (LL) or concrete action theory  $\mathcal{D}_l$  representing the dynamics of the domain at different levels of detail.  $\mathcal{D}_h$  (resp.  $\mathcal{D}_l$ ) involves a finite set of primitive action types  $\mathcal{A}_h$  (resp.  $\mathcal{A}_l$ ) and a finite set of primitive fluent predicates  $\mathcal{F}_h$  (resp.  $\mathcal{F}_l$ ). Also,  $\mathcal{D}_h$  and  $\mathcal{D}_l$  are assumed to share no domain specific symbols except for standard names for objects in  $\mathcal{N}$ .

**Definition 3 (Refinement Mapping)** *A refinement mapping  $m$  is a function that associates each HL primitive action type  $A$  in  $\mathcal{A}_h$  to a SD ConGolog program  $\delta_A$  defined over the LL theory that implements the action, i.e.,  $m(A(\vec{x})) = \delta_A(\vec{x})$ . Moreover,  $m$  maps each situation-suppressed HL fluent  $F(\vec{x})$  in  $\mathcal{F}_h$  to a situation-suppressed formula  $\Phi_F(\vec{x})$  defined over the LL theory that characterizes the concrete conditions under which  $F(\vec{x})$  holds in a situation.*

We extend the notation so that  $m(\Phi)$  stands for the result of substituting every fluent  $F(\vec{x})$  in situation-suppressed formula  $\Phi$  by  $m(F(\vec{x}))$ . Also, we apply  $m$  to action sequences with  $m(\alpha_1, \dots, \alpha_n) \doteq m(\alpha_1); \dots; m(\alpha_n)$  for  $n \geq 1$  and  $m(\epsilon) \doteq nil$ , where  $\epsilon$  is the empty sequence of actions.

To relate the HL and LL models/theories, a variant of bisimulation (Milner 1989) is defined as follows.

**Definition 4 (*m*-Bisimulation)** Given  $M_h$  a model of  $\mathcal{D}_h$ , and  $M_l$  a model of  $\mathcal{D}_l \cup \mathcal{K}$ , a relation  $B \subseteq \Delta_S^{M_h} \times \Delta_S^{M_l}$  (where  $\Delta_S^M$  stands for the situation domain of  $M$ ) is an *m*-bisimulation relation between  $M_h$  and  $M_l$  if  $\langle s_h, s_l \rangle \in B$  implies that: (i)  $s_h$  evaluates each HL primitive fluent the same as the evaluation of the refinement of the fluent in  $s_l$ ; (ii) for every HL primitive action type  $A$  in  $\mathcal{A}_h$ , if there exists  $s'_h$  s.t.  $M_h \models \text{Poss}(A(\vec{x}), s_h) \wedge s'_h = \text{do}(A(\vec{x}), s_h)$ , then there exists  $s'_l$  s.t.  $M_l \models \text{Do}(m(A(\vec{x})), s_l, s'_l)$  and  $\langle s'_h, s'_l \rangle \in B$ ; and (iii) for every HL primitive action type  $A$  in  $\mathcal{A}_h$ , if there exists  $s'_l$  s.t.  $M_l \models \text{Do}(m(A(\vec{x})), s_l, s'_l)$ , then there exists  $s'_h$  s.t.  $M_h \models \text{Poss}(A(\vec{x}), s_h) \wedge s'_h = \text{do}(A(\vec{x}), s_h)$  and  $\langle s'_h, s'_l \rangle \in B$ .

$M_h$  is *m*-bisimilar to  $M_l$ , written  $M_h \sim_m M_l$ , iff there exists an *m*-bisimulation relation  $B$  between  $M_h$  and  $M_l$  such that  $(S_0^{M_h}, S_0^{M_l}) \in B$ .

**Definition 5 (Sound abstraction)**  $\mathcal{D}_h$  is a sound abstraction of  $\mathcal{D}_l$  relative to refinement mapping  $m$  iff for all models  $M_l$  of  $\mathcal{D}_l \cup \mathcal{K}$ , there exists a model  $M_h$  of  $\mathcal{D}_h$  s.t.  $M_h \sim_m M_l$ .

With a sound abstraction, if the HL theory entails that a sequence of actions is executable and achieves a condition, then the LL must also entail that there exists an executable refinement of the sequence such that the “translated” condition holds afterwards. Also, if the LL theory considers the executability of a refinement of a sequence of HL actions to be satisfiable and a refinement of an HL condition to hold afterwards, then the HL must also consider the executability of the sequence of HL actions satisfiable after which the condition must hold as well.

**Definition 6 (Complete abstraction)**  $\mathcal{D}_h$  is a complete abstraction of  $\mathcal{D}_l$  relative to refinement mapping  $m$  iff for all models  $M_h$  of  $\mathcal{D}_h$ , there exists a model  $M_l$  of  $\mathcal{D}_l \cup \mathcal{K}$  s.t.  $M_l \sim_m M_h$ .

With a complete abstraction, if the LL theory entails that some refinement of a sequence of HL actions is executable and achieves a “translated” HL condition, then the HL also entails that the action sequence is executable and the condition holds afterwards. Also, if the HL theory considers the executability of a sequence of actions to be satisfiable and a condition to hold after that, then the LL must also consider the executability of the refinement of the sequence of HL actions satisfiable after which a “translated” condition must hold as well.

Note that this approach supports the use of ConGolog programs to specify the possible dynamics of the domain at both the HL and LL; this is done by following (De Giacomo et al. 2016) and “compiling” the program into the BAT  $\mathcal{D}$  to get a new BAT  $\mathcal{D}'$  whose executable situations are exactly those that can be reached by executing the program.

**Example (Cont’d).** In our example, we define an HL BAT  $\mathcal{D}_h^{ex}$  that abstracts over some details of  $\mathcal{D}_l^{ex}$ . At the HL, we abstract over details of rescue actions. Action  $\text{rescue}(r, l)$  abstracts over the process of either clearing rubble or evacuating people. The fluent  $\text{AidedInRescue}(r, l, s)$  indicates if robot  $r$  has aided in rescue at location  $l$ . For simplicity, actions  $\text{updateSW}_{HL}(r, v)$  and  $\text{go}_{HL}(r, l)$  are defined similar to  $\text{updateSW}_{LL}(r, v)$  and  $\text{go}_{LL}(r, l)$  respectively.

$\mathcal{D}_h^{ex}$  includes the following precondition axioms:

$$\begin{aligned} \text{Poss}(\text{updateSW}_{HL}(r, v), s) &\equiv \neg \text{UpdatedSW}_{HL}(r, v, s), \\ \text{Poss}(\text{go}_{HL}(r, l), s) &\equiv \neg \text{At}_{HL}(r, l, s), \\ \text{Poss}(\text{rescue}(r, l), s) &\equiv \text{At}_{HL}(r, l, s). \end{aligned}$$

The HL BAT also includes the following SSAs:

$$\begin{aligned} \text{AidedInRescue}(r, l, \text{do}(a, s)) &\equiv \\ a = \text{rescue}(r, l) \vee \text{AidedInRescue}(r, l, s). \end{aligned}$$

$\text{At}_{HL}$  and  $\text{UpdatedSW}_{HL}$  have SSAs similar to their LL counterparts respectively.

$\mathcal{D}_h^{ex}$  contains the following initial state axioms:

$$\begin{aligned} \text{At}_{HL}(\text{Rob}, \text{Station}, S_0), \neg \text{UpdatedSW}_{HL}(\text{Rob}, V2021, S_0), \\ \neg \text{AidedInRescue}(\text{Rob}, L1, S_0). \end{aligned}$$

**Refinement Mapping  $m^{ex}$**  We specify the relationship between the HL and LL BATs through a refinement mapping  $m^{ex}$  which is defined as follows:

$$\begin{aligned} m^{ex}(\text{go}_{HL}(r, l)) &= \text{go}_{LL}(r, l), \\ m^{ex}(\text{updateSW}_{HL}(r, v)) &= \text{updateSW}_{LL}(r, v), \\ m^{ex}(\text{rescue}(r, l)) &= \text{evacuate}(r, l) \mid \text{removeRubble}(r, l), \\ m^{ex}(\text{At}_{HL}(r, l)) &= \text{At}_{LL}(r, l), \\ m^{ex}(\text{UpdatedSW}_{HL}(r, v)) &= \text{UpdatedSW}_{LL}(r, v), \\ m^{ex}(\text{AidedInRescue}(r, l)) &= \text{Cleared}(r, l) \vee \text{Evacuated}(r, l). \end{aligned}$$

By using Theorem 9 in (Banihashemi, De Giacomo, and Lespérance 2017), it can be confirmed that  $\mathcal{D}_h^{ex}$  is a sound abstraction of  $\mathcal{D}_l^{ex}$  relative to the mapping  $m^{ex}$ .  $\triangleleft$

## Programs as Actual Causes

We now return to our discussion of abstract causes. As seen in the previous section, Definition 2 appeals to regression, a syntactic notion, and this requires the effect formula  $\Phi[s]$  to be uniform in  $s$ . However, this is too restrictive for us as it is hard to adapt for abstract causes. Specifically, it is hard to define regression over programs; recall Reiter defined regression over primitive actions.<sup>1</sup> Therefore, we start by introducing the notion of *dynamic effect formulae* in the SC, which is motivated by the notion of *epistemic dynamic formulae* (Khan and Lespérance 2021).

**Definition 7 (Dynamic Effect Formula)** Let  $\vec{x}$  and  $\theta_{\vec{a}}$  respectively range over object terms and a sequence of action terms. The class of situation-suppressed dynamic effect formulae  $\psi$  is defined inductively using the following grammar:

$$\psi ::= P(\vec{x}) \mid \text{ExecSeq}(\theta_{\vec{a}}) \mid \text{After}(\theta_{\vec{a}}, \psi) \mid \neg\psi \mid \psi_1 \wedge \psi_2 \mid \exists \vec{x}. \psi.$$

That is, a dynamic effect formula can be a situation-suppressed fluent, a formula that says that some sequence of actions  $\theta_{\vec{a}}$  is executable, a formula that indicates some dynamic effect formula holds after some sequence of actions has occurred, or one that can be built from other dynamic effect formulae using the usual connectives. Note that  $\psi$  can have quantification over object variables, but must not include quantification over situations or the precedence operator  $\sqsubset$ . We use lower-case Greek letters for dynamic effect

<sup>1</sup>Note that, previously (Mo, Li, and Liu 2016) has proposed an extension of regression for programs; investigating whether their definition can be adapted for our purpose is future work.

formulae.  $\psi[s]$  is the formula obtained from  $\psi$  by restoring the appropriate situation argument into all fluents in  $\psi$ .

**Definition 8**

$$\psi[s] \doteq \begin{cases} P(\vec{x}, s) & \text{if } \psi \text{ is } P(\vec{x}) \\ \exists s'. Do(\theta_{\vec{a}}, s, s') & \text{if } \psi \text{ is } ExecSeq(\theta_{\vec{a}}) \\ \psi'[do([\theta_{\vec{a}}], s)] & \text{if } \psi \text{ is } After(\theta_{\vec{a}}, \psi') \\ \neg(\psi'[s]) & \text{if } \psi \text{ is } (\neg\psi') \\ \psi_1[s] \wedge \psi_2[s] & \text{if } \psi \text{ is } (\psi_1 \wedge \psi_2) \\ \exists \vec{y}. (\psi'[s]) & \text{if } \psi \text{ is } (\exists \vec{y}. \psi') \end{cases}$$

We generalize causal settings by allowing effects in our framework to be any dynamic effect formula  $\psi$ , i.e., we no longer require the effect to be uniform in  $s$ . Also, we do not require the trace to be a ground situation term, so it can now include arbitrary (non-ground) action terms. This generalization allows for the modeling of abstract causes.

**Definition 9 (Generalized Causal Setting)** A *generalized causal setting* is a tuple  $\langle \mathcal{D}, \delta, \psi \rangle$ , where  $\mathcal{D}$  is a BAT,  $\delta$  is a ConGolog program, and  $\psi$  is a dynamic effect formula s.t.:

$$\mathcal{D} \cup \mathcal{K} \models \neg\psi[S_0] \wedge \exists s'. Do(\delta, S_0, s') \wedge \psi[s'].$$

Thus, there is at least one execution of the program  $\delta$  starting in the initial situation  $S_0$  after which the effect  $\psi$  holds.

As discussed in the previous section, the definition of actual achievement cause given by (Batusov and Soutchanski 2018) only deals with narratives that are linear sequences of actions. Consequently, their causes are actions (or more precisely, action-situation pairs).<sup>2</sup> To facilitate the modeling of abstract causes, we extend this by allowing narratives to be linear sequences of ConGolog programs. This allows programs to be identified as causes of observed effects. In the following, we progressively define what it means for a ConGolog program to be a weak potential cause, starting with primary causes. Note that, given a generalized causal setting there can be more than one primary potential cause of the effect as the program can have multiple possible executions.

**Definition 10** Given a generalized causal setting  $\mathcal{C} = \langle \mathcal{D}, (\delta_1; \dots; \delta_n), \psi \rangle$  and a model  $M$  of  $\mathcal{D} \cup \mathcal{K}$ , a program  $\delta_{i+1} \in \{\delta_1, \dots, \delta_n\}$  is called a *primary weak potential cause of  $\psi$  relative to  $\mathcal{C}$  and  $M$  if and only if*:

$$\begin{aligned} M \models & \exists s_i, s_{i+1}, s_n. Do((\delta_1; \dots; \delta_i), S_0, s_i) \wedge \neg\psi[s_i] \\ & \wedge Do(\delta_{i+1}, s_i, s_{i+1}) \wedge Do((\delta_{i+2}; \dots; \delta_n), s_{i+1}, s_n) \\ & \wedge \forall s'. s_{i+1} \leq s' \leq s_n \supset \psi[s']. \end{aligned}$$

The triple  $(s_i, s_{i+1}, \psi)$  is called a *witness for this*.

That is, a program  $\delta_{i+1}$  in the scenario  $(\delta_1; \dots; \delta_n)$  is a *primary weak potential cause relative to a model  $M$  of theory  $\mathcal{D} \cup \mathcal{K}$  and causal setting  $\mathcal{C}$*  if and only if there is an execution of the prefix  $(\delta_1; \dots; \delta_i)$  that ends in situation  $s_i$  in which  $\psi$  is false, situation  $s_{i+1}$  can be reached by executing  $\delta_{i+1}$  starting from  $s_i$ , situation  $s_n$  can be reached by

<sup>2</sup>Here and in the sequel, for brevity, we omit the terms *actual* and *achievement* when we talk about causes, since we exclusively consider actual achievement causes in this paper.

executing the remaining programs starting from  $s_{i+1}$ , and  $\psi$  holds in all situations from  $s_{i+1}$  up to  $s_n$ . Essentially, this is a straightforward generalization of the base case of Definition 2 and ensures that there is an execution of the scenario over which  $\psi$  was achieved by some action in  $\delta_{i+1}$  and  $\psi$  persisted till the end of the trace, i.e., it was not later made false by a subsequent action.

Moreover, we define what it means for a program to be a primary weak potential cause relative to a causal setting.

**Definition 11 (Primary Weak Potential Cause)** Given a generalized causal setting  $\mathcal{C} = \langle \mathcal{D}, (\delta_1; \dots; \delta_n), \psi \rangle$ , a program  $\delta_i \in \{\delta_1, \dots, \delta_n\}$  is called a *primary weak potential cause relative to  $\mathcal{C}$  if and only if for all models  $M$  of  $\mathcal{D} \cup \mathcal{K}$ ,  $\delta_i$  is a primary weak potential cause of  $\psi$  relative to  $\mathcal{C}$  and  $M$ .*

Next, we define weak potential causes in general. These include both primary and non-primary causes reflecting both base and inductive cases of Definition 2.

**Definition 12** Given a generalized causal setting  $\mathcal{C} = \langle \mathcal{D}, (\delta_1; \dots; \delta_n), \psi \rangle$  and a model  $M$  of  $\mathcal{D} \cup \mathcal{K}$ , a program  $\delta_i \in \{\delta_1, \dots, \delta_n\}$  is called a *weak potential cause of  $\psi$  relative to  $\mathcal{C}$  and  $M$  if and only if*:

1.  $\delta_i$  is a primary weak potential cause wrt  $\mathcal{C}$  and  $M$  with witness  $(s', s'', \psi')$ , where  $\psi' = \psi$ , or
2.  $\delta_j$  (where  $i < j \leq n$ ) is a weak potential cause relative to setting  $\mathcal{C}$  and  $M$  with witness  $(s_{j-1}, do([\vec{a}_j], s_{j-1}), \psi_j)$ , and  $\delta_i$  is a primary weak potential cause relative to the setting  $\langle \mathcal{D}, (\delta_1; \dots; \delta_{j-1}), \psi' \rangle$  and model  $M$  with witness  $(s', s'', \psi')$ , where  $\psi' = ExecSeq(\vec{a}_j) \wedge After(\vec{a}_j, \psi_j)$ .

We call  $(s', s'', \psi')$  a *witness for  $\delta_i$  being a weak potential cause wrt  $\mathcal{C}$  and  $M$ .*

Thus,  $\delta_i$  is a *weak potential cause relative to model  $M$  and generalized causal setting  $\mathcal{C}$*  if and only if it is either a primary weak potential cause wrt  $\mathcal{C}$  and  $M$ , or it is a primary weak potential cause of another weak potential cause  $\delta_j$ , i.e., it enables  $\delta_j$  by ensuring that the appropriate execution path  $\vec{a}_j$  of  $\delta_j$  that brought about  $\delta_j$ 's own effect  $\psi_j$  is executable (i.e., that  $ExecSeq(\vec{a}_j)$ ) and by fulfilling the conditions under which the execution of  $\vec{a}_j$  achieved  $\psi_j$  (i.e., that  $After(\vec{a}_j, \psi_j)$ ).

**Definition 13 (Weak Potential Cause)** Given a generalized causal setting  $\mathcal{C} = \langle \mathcal{D}, (\delta_1; \dots; \delta_n), \psi \rangle$ , a program  $\delta_i \in \{\delta_1, \dots, \delta_n\}$  is called a *weak potential cause of  $\psi$  relative to  $\mathcal{C}$  if and only if for all models  $M$  of  $\mathcal{D} \cup \mathcal{K}$ ,  $\delta_i$  is a weak potential cause of  $\psi$  relative to  $\mathcal{C}$  and  $M$ .*

Moreover, if  $\mathcal{D}$  is initially completely specified, there is only one model; in that case, we call  $\psi'$  from the witness  $(s', s'', \psi')$  in Definition 12 a *witness to the fact that  $\delta_i$  is a weak potential cause of  $\psi$  relative to  $\mathcal{C}$ .*

Thus, we only call a program a *weak potential cause relative to a generalized causal setting* if it is a weak potential cause in all models of the theory.

**Example (Cont'd).** Consider the generalized causal setting  $\langle \mathcal{D}_i^{ex}, (updateSW_{LL}(Rob, V2021); go_{LL}(Rob, L1));$

$\delta_{rescue}(Rob, L1), \exists r, l. Cleared(r, l)$ , where  $\delta_{rescue}(r, l) = m^{ex}(rescue(r, l))$ . Then according to our definitions,  $\delta_{rescue}(Rob, L1)$  is the primary weak potential cause relative to the above setting, as in all models,  $\mathcal{D}_l^{ex} \cup \mathcal{K} \models \exists s_2, s_3. Do((updateSW_{LL}(Rob, V2021); go_{LL}(Rob, L1)), S_0, s_2) \wedge \neg \exists r, l. Cleared(r, l)[s_2] \wedge Do(\delta_{rescue}(Rob, L1), s_2, s_3) \wedge \exists r, l. Cleared(r, l)[s_3]$ , and by the SSA for *Cleared*, the effect persists until the end of scenario.

Note that, as only in some executions of the scenario  $\exists r, l. Cleared(r, l)$  is true,  $\delta_{rescue}(Rob, L1)$  is considered *weak*. If we instead consider the effect  $\exists r, l. Cleared(r, l) \vee Evacuated(r, l)$ , then  $\delta_{rescue}(Rob, L1)$  can be considered as the primary strong potential cause in the sense that in all executions of the scenario  $\delta_{rescue}$  achieves the effect.

Moreover, we can also show that  $go_{LL}(Rob, L1)$  is a weak potential cause, since it is a primary weak potential cause wrt the setting  $\langle \mathcal{D}_l^{ex}, (updateSW_{LL}(Rob, V2021); go_{LL}(Rob, L1)), ExecSeq(removeRubble(Rob, L1)) \wedge After(removeRubble(Rob, L1), \exists r, l. Cleared(r, l)) \rangle$ .

On the other hand,  $updateSW_{LL}(Rob, V2021)$  cannot be shown to be a weak potential cause.  $\triangleleft$

Our notion of programs as actual cause above is a weak and more inclusive one. We consider a program as a cause if there is at least one execution where the program is a cause. In some cases, it might be useful to consider a stronger version, where a program is considered to be a cause if it is a cause according to all executions of the program. A thorough investigation of such a variant is future work.

When the program  $\delta$  is finite, terminating, and composed of ground actions only, one can show that the intermediate effects (i.e.,  $[ExecSeq(\vec{a}) \wedge After(\vec{a}, \phi)]$ ) can be straightforwardly computed using Reiter's regression. Also, and in particular, when  $\delta$  is a finite sequence of ground actions, the causes computed using our definition and Batusov and Soutchanski's (2018) definition are the same.

**Theorem 14** *Let  $\delta = \alpha_1; \dots; \alpha_n$  be a finite sequence of ground actions. Then  $(\alpha_i, do([\alpha_1, \dots, \alpha_{i-1}], S_0))$  is a cause relative to the causal setting  $\langle \mathcal{D}, do([\delta], S_0), \Phi \rangle$  according to Definition 2 iff  $\alpha_i$  is a potential cause relative to the generalized causal setting  $\langle \mathcal{D}, \delta, \Phi \rangle$  according to Definition 13.<sup>3</sup>*

*Proof.* (Sketch) By Definitions 2, 8, 10, 11, and the definition of *Do*, the primary cause computed by both definitions is the same. Moreover, by induction on the number of causes and using Definitions 2, 8, 10, 11, 12, 13, the definition of *Do*, and Reiter's (2001) regression theorem, it can be shown that the intermediate effects in these two definitions (i.e., formulae  $\rho[\Phi[s], \alpha'] \wedge \Pi_{apa}(\alpha', \sigma')$  in Definition 2 and  $ExecSeq(\vec{a}_j) \wedge After(\vec{a}_j, \psi)$  in Definition 12) are equivalent, and thus it can be shown that the causes computed in both cases are the same.  $\square$

Given the above, it is easy to see that when  $\delta$  is a finite sequence of ground actions, all properties shown for (Batusov and Soutchanski 2018)'s framework also hold in ours. These include the proper handling of preemption and switches.

<sup>3</sup>Note that, an SC formula  $\Phi$  is also a dynamic effect formula.

## Reasoning about Abstract Causes

We now focus on investigating how reasoning about abstract causes can be simplified. In particular, we will show that under some conditions, a subclass of weak potential causes at various levels of abstraction can be related. This subclass involves weak causes that are also strong in the sense that all executions of the cause achieve the effect (see the corollary below). This reduces reasoning about abstract causes (i.e., programs) at the LL to that of actions as causes at the HL when said conditions are met.

We start by formalizing some of these conditions. First, we assume that every HL action  $\alpha_i$  is mapped via  $m$  to an LL program  $\delta_i$  that may take part in an LL scenario; in this way, any abstract scenario can be refined by a concrete one.

Moreover, we assume only action sequences that refine some HL action sequence are executed in the LL BAT:

**Assumption 1 (All LL behaviours refine HL actions)**

$\mathcal{D}_l \cup \mathcal{K} \models \forall s. Executable(s) \supset \exists \delta. Trans^*(ANYSEQHL, S_0, \delta, s)$ ,  
where  $ANYSEQHL \doteq (\bigwedge_{A_i \in \mathcal{A}_h} \pi \vec{x}. m(A_i(\vec{x})))^*$ ,  
i.e., do any sequence of refinements of HL actions.

Furthermore, we require that LL effects are non-transient wrt HL actions:

**Assumption 2 (Non-transiency of LL Effects)** *Suppose the set  $F_R^{F_i}$  includes all the fluent literals in a refinement of an HL fluent  $F_i$ . We assume that:*

$$\begin{aligned} \mathcal{D}_l \cup \mathcal{K} \models \forall s. Do(ANYSEQHL, S_0, s) \supset \\ \bigwedge_{A_i \in \mathcal{A}_h} \bigwedge_{F_i \in \mathcal{F}_h} \bigwedge_{F_L \in F_R^{F_i}} \forall s', s'', \vec{x}, \vec{y}. \\ F_L(\vec{y})[s] \wedge Do(m(A_i(\vec{x})), s, s') \wedge F_L(\vec{y})[s'] \\ \wedge s < s'' < s' \supset F_L(\vec{y})[s''] \end{aligned}$$

The above essentially requires the LL theory to entail that if a fluent literal  $F_L$  that is in a refinement of an HL fluent  $F_i$  is true in both the situations before and after the execution of the refinement of an HL action  $A_i(\vec{x})$ , then it should remain true in all intermediate situations of execution of the refinement of  $A_i(\vec{x})$  as well. This condition must hold after any sequence of refinements of HL actions, i.e.,  $Do(ANYSEQHL, S_0, s)$ . To see why this is necessary, consider the following example. Suppose that at the HL, we have the generalized causal setting  $\langle \mathcal{D}_h, (\alpha; \beta), F_{hl} \rangle$  in which  $\alpha$  is the only primary weak potential cause. Assume the following mapping:  $m(\alpha) = a$  and  $m(\beta) = b_1; b_2$ , and  $m(F_{hl}) = F_{ll}$ . At the LL, after performing  $a$ ,  $F_{ll}$  becomes true, and after performing  $b_1$  and  $b_2$ ,  $F_{ll}$  becomes false and true respectively. Hence, in the setting  $\langle \mathcal{D}_l, m(\alpha; \beta), m(F_{hl}) \rangle$ ,  $m(\beta)$  is considered the only primary weak potential cause if the analysis is done at the LL using Definition 13. To achieve correspondence of potential causes between HL and LL, we need to rule out such cases.

To investigate how causes at the abstract and concrete levels are related, we first consider *sound abstractions*. For this, we first show that if at the HL, an effect  $\Phi$  which is not true in the initial situation, holds after the execution of sequence of actions  $\vec{a}$ , then at the LL, the refinement of  $\Phi$  is false in the initial situation, and holds after the execution of a refinement of  $\vec{a}$ . Moreover, if at the HL, an action  $\alpha_k$  (in  $\vec{a}$ ) is

considered a primary weak potential cause wrt the generalized casual setting  $\mathcal{C}_h = \langle \mathcal{D}_h, (\vec{\alpha}), \Phi \rangle$ , then a refinement of  $\alpha_k$  can be considered a primary weak potential cause wrt the setting  $\mathcal{C}_m = \langle \mathcal{D}_l, m(\vec{\alpha}), m(\Phi) \rangle$  at the LL.<sup>4</sup>

**Theorem 15** *Suppose that  $\mathcal{D}_h$  is a sound abstraction of  $\mathcal{D}_l$  wrt some refinement mapping  $m$ , and that Assumptions 1 and 2 hold. Then for any ground HL action sequence  $\vec{\alpha}$  and for any HL situation suppressed formula  $\Phi$  such that  $\mathcal{D}_h \models Executable(do(\vec{\alpha}, S_0)) \wedge \neg \Phi[S_0] \wedge \Phi[do([\vec{\alpha}], S_0)]$ , we have:*

1.  $\mathcal{D}_l \cup \mathcal{K} \models \neg m(\Phi)[S_0] \wedge \exists s. Do(m(\vec{\alpha}), S_0, s) \wedge m(\Phi)[s]$ .
2. If  $\vec{\alpha} = \vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1}$  and  $\alpha_k$  is the primary weak potential cause wrt the generalized causal setting  $\mathcal{C}_h = \langle \mathcal{D}_h, (\vec{\alpha}), \Phi \rangle$ , then  $m(\alpha_k)$  is the unique primary weak potential cause wrt the generalized causal setting  $\mathcal{C}_m = \langle \mathcal{D}_l, m(\vec{\alpha}), m(\Phi) \rangle$ .

*Proof.* (Sketch) Part 1 follows from Theorem 4 in (Banihashemi, De Giacomo, and Lespérance 2017), as  $\mathcal{D}_h$  is a sound abstraction of  $\mathcal{D}_l$  wrt mapping  $m$ . Based on part 1, we conclude that there is a primary weak potential cause at the low-level. Part 2 is shown by contradiction. Suppose  $\alpha_k$  is the primary weak potential cause relative to  $\mathcal{C}_h$ , but  $m(\alpha_k)$  is not the unique primary weak potential cause relative to  $\mathcal{C}_m$ . We know by Assumption 1 that all low-level actions (that may be executed) are part of some refinement of some high-level action. Hence, we can assume  $m(\vec{\alpha}) = m(\vec{\alpha}_{j-1} \alpha_j \vec{\alpha}_{j+1})$  where  $j \neq k$ , and  $m(\alpha_j)$  is the primary weak potential cause in the casual setting  $\mathcal{C}_m$ . This means that we can take an arbitrary model  $M_l$  of  $\mathcal{D}_l \cup \mathcal{K}$  such that  $m(\alpha_j)$  is the primary weak potential cause relative to  $\mathcal{C}_m$  and  $M_l$ . Since  $\mathcal{D}_h$  is a sound abstraction of  $\mathcal{D}_l$  wrt to  $m$ , there exists a model  $M_h$  of  $\mathcal{D}_h$  such that  $M_h$  and  $M_l$  are bisimilar, and that  $\alpha_j$  is considered the primary weak potential cause relative to  $\mathcal{C}_h$  and  $M_h$ . But this contradicts the fact that  $\mathcal{D}_h$  entails that  $\alpha_k$  is the primary weak potential cause wrt  $\mathcal{C}_h$ . Assumption 2 ensures that the effect achieved by  $m(\alpha_k)$  is persistent until the end of the trace.  $\square$

We now focus on showing the correspondence of non-primary (indirect) causes at the abstract and concrete levels. As the witnessing formulae may be different for each low-level model, we assume complete information (a single model) at the LL.

**Theorem 16** *Suppose that  $\mathcal{D}_h$  is a sound abstraction of  $\mathcal{D}_l$  wrt some refinement mapping  $m$ , and that Assumptions 1 and 2 hold. Then for any ground HL action sequence  $\vec{\alpha}$  and for any HL situation suppressed formula  $\Phi$  such that  $\mathcal{D}_h \models Executable(do(\vec{\alpha}, S_0)) \wedge \neg \Phi[S_0] \wedge \Phi[do([\vec{\alpha}], S_0)]$ , we have:*

- If  $\vec{\alpha} = \vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1} \alpha_j \vec{\alpha}_{j+1}$ ,  $\alpha_j$  is a weak potential cause wrt the generalized causal setting  $\mathcal{C}_h = \langle \mathcal{D}_h, (\vec{\alpha}), \Phi \rangle$  with witness  $\Phi_j$ , and,
- $\alpha_k$  is the primary weak potential cause wrt the setting  $\mathcal{C}'_h = \langle \mathcal{D}_h, (\vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1}), ExecSeq(\alpha_j) \wedge After(\alpha_j, \Phi_j) \rangle$ , and,

<sup>4</sup>In the following, we will quantify over action sequences and so we need to encode sequences as first-order terms as in (De Giacomo, Lespérance, and Levesque 2000). For notational simplicity, we suppress this encoding and use sequences as terms directly.

- $\mathcal{D}_l$  is initially completely specified, and,
- $m(\alpha_j)$  is a weak potential cause wrt the causal setting  $\mathcal{C}_m = \langle \mathcal{D}_l, m(\vec{\alpha}), m(\Phi) \rangle$  with witness  $m(\Phi_j)$ ,

then,  $m(\alpha_k)$  is the unique primary weak potential cause wrt the generalized causal setting  $\mathcal{C}'_m = \langle \mathcal{D}_l, m(\vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1}), \phi'_L \rangle$ , where  $\mathcal{D}_l \cup \mathcal{K} \models \exists s^*, \vec{a}_j. Do(m(\vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1}), S_0, s^*) \wedge Do(m(\alpha_j), s^*, do([\vec{a}_j], s^*))$  and  $\phi'_L = ExecSeq(\vec{a}_j) \wedge After(\vec{a}_j, m(\Phi_j))$ .

*Proof.* (Sketch) The proof is by contradiction. Suppose  $\alpha_k$  is the primary weak potential cause relative to  $\mathcal{C}'_h$ , but  $m(\alpha_k)$  is not the unique primary weak potential cause relative to  $\mathcal{C}'_m$ . We know by Assumption 1 that all low-level actions (that may be executed) are part of some refinement of some high-level action. Hence, we can assume another action  $\alpha_p$  where  $p \neq k$  is in the sequence of actions  $\vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1}$ , and  $m(\alpha_p)$  is the primary weak potential cause in the casual setting  $\mathcal{C}'_m$ . This means that we can take an arbitrary model  $M_l$  of  $\mathcal{D}_l \cup \mathcal{K}$  such that  $m(\alpha_p)$  is the primary weak potential cause relative to  $\mathcal{C}'_m$  and  $M_l$ . Since  $\mathcal{D}_h$  is a sound abstraction of  $\mathcal{D}_l$  wrt to  $m$ , there exists a model  $M_h$  of  $\mathcal{D}_h$  such that  $M_h$  and  $M_l$  are bisimilar, and that  $\alpha_p$  is considered the primary weak potential cause relative to  $\mathcal{C}'_h$  and  $M_h$ . But this contradicts the fact that  $\mathcal{D}_h$  entails that  $\alpha_k$  is the primary weak potential cause wrt  $\mathcal{C}'_h$ . Assumption 2 ensures that the effect achieved by  $m(\alpha_k)$  is persistent until the end of the trace, and having a single model ensures that there is a sequence of actions  $\vec{a}_j$  such that  $ExecSeq(\vec{a}_j) \wedge After(\vec{a}_j, m(\Phi_j))$ .  $\square$

**Example (Cont'd).** Consider the HL setting  $\mathcal{C}_h = \langle \mathcal{D}_h^{ex}, (\vec{\alpha}), \phi_e \rangle$ , where  $\vec{\alpha} = [updateSW_{HL}(Rob, V2021), go_{HL}(Rob, L1), rescue(Rob, L1)]$  and  $\phi_e = \exists r, l. AidedInRescue(r, l)$ . Using similar reasoning as before, we can show that  $rescue(Rob, L1)$  is the primary weak potential cause relative to  $\mathcal{C}_h$ . Moreover,  $go_{HL}(Rob, L1)$  is another cause relative to  $\mathcal{C}_h$ .

By Theorem 15, we have that  $m^{ex}(rescue(Rob, L1)) = \delta_{rescue}(Rob, L1)$  is the primary weak potential cause relative to setting  $\mathcal{C}_m = \langle \mathcal{D}_l^{ex}, m^{ex}(\vec{\alpha}), m^{ex}(\phi_e) \rangle$ , where  $m^{ex}(\phi_e) = \exists r, l. Cleared(r, l) \vee Evacuated(r, l)$  and  $m^{ex}(\vec{\alpha}) = updateSW_{LL}(Rob, V2021); go_{LL}(Rob, L1); \delta_{rescue}(Rob, L1)$ . Moreover, by Theorem 16, the action  $go_{LL}(Rob, L1)$  is considered another weak potential cause relative to  $\mathcal{C}_m$ .  $\triangleleft$

Notice that since the number of actions and fluents that a reasoner needs to consider are typically higher at the LL, Theorems 15 and 16 can yield important efficiency benefits.

Corollary 5 of (Banihashemi, De Giacomo, and Lespérance 2017) ([BDL17]), showed that if  $\mathcal{D}_h$  is a sound abstraction of  $\mathcal{D}_l$  wrt  $m$ , then the different sequences of LL actions that are refinements of a given HL action sequence all have the same effects on the HL fluents, and more generally on HL situation-suppressed formulae, i.e., from the HL perspective they are deterministic:

**Corollary 17 (from BDL17)** *If  $\mathcal{D}_h$  is a sound abstraction of  $\mathcal{D}_l$  wrt  $m$ , then for any sequence of ground HL actions  $\vec{\alpha}$  and for any HL situation-suppressed formula  $\phi$ , we have:*

$$\mathcal{D}_l \cup \mathcal{K} \models \forall s, s'. Do(m(\vec{\alpha}), S_0, s) \wedge Do(m(\vec{\alpha}), S_0, s') \supset (m(\phi)[s] \equiv m(\phi)[s']).$$

This indicates that the weak potential causes in Theorems 15 and 16 are in fact strong in the sense that in all executions of the program, the effect is achieved.

With *complete abstractions*, we can show that if at the LL, the refinement of an effect  $\Phi$  which is not true in the initial situation, holds after the execution of a refinement of sequence of actions  $\vec{\alpha}$ , then at the HL, the effect  $\Phi$  is false in the initial situation, and holds after the execution of  $\vec{\alpha}$ . Moreover, if at the LL, the refinement of an action  $\alpha_k$  (in  $\vec{\alpha}$ ) is considered the primary weak potential cause wrt the generalized casual setting  $\mathcal{C}_m = \langle \mathcal{D}_l, m(\vec{\alpha}), m(\Phi) \rangle$ , then  $\alpha_k$  can be considered a primary weak potential cause wrt the setting  $\mathcal{C}_h = \langle \mathcal{D}_h, (\vec{\alpha}), \Phi \rangle$  at the HL.

**Theorem 18** *Suppose that  $\mathcal{D}_h$  is a complete abstraction of  $\mathcal{D}_l$  wrt some refinement mapping  $m$ . Then for any ground HL action sequence  $\vec{\alpha}$  and for any HL situation suppressed formula  $\Phi$  such that  $\mathcal{D}_l \cup \mathcal{K} \models \neg m(\Phi)[S_0] \wedge \exists s. Do(m(\vec{\alpha}), S_0, s) \wedge m(\Phi)[s]$ , we have that:*

1.  $\mathcal{D}_h \models \neg \Phi[S_0] \wedge Executable(do([\vec{\alpha}], S_0)) \wedge \Phi(do([\vec{\alpha}], S_0))$ .
2. If  $\vec{\alpha} = \vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1}$  and  $m(\alpha_k)$  is the primary weak potential cause wrt the generalized causal setting  $\mathcal{C}_m = \langle \mathcal{D}_l, m(\vec{\alpha}), m(\Phi) \rangle$  then  $\alpha_k$  is the unique primary weak potential cause wrt the setting  $\mathcal{C}_h = \langle \mathcal{D}_h, (\vec{\alpha}), \Phi \rangle$ .

When considering non-primary (indirect) causes at the abstract and concrete levels, similar to sound abstraction, we need to assume complete information (a single model) at the LL, since the witnessing formulae may be different for each LL model. The following theorem shows the correspondence between indirect causes at the concrete and abstract levels.

**Theorem 19** *Suppose that  $\mathcal{D}_h$  is a complete abstraction of  $\mathcal{D}_l$  wrt some refinement mapping  $m$ . Then for any ground HL action sequence  $\vec{\alpha}$  and for any HL situation suppressed formula  $\Phi$  such that  $\mathcal{D}_l \cup \mathcal{K} \models \neg m(\Phi)[S_0] \wedge \exists s. Do(m(\vec{\alpha}), S_0, s) \wedge m(\Phi)[s]$ , we have that:*

- If  $\mathcal{D}_l$  is initially completely specified, and,
- $\vec{\alpha} = \vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1} \alpha_j \vec{\alpha}_{j+1}$ ,  $m(\alpha_j)$  is a weak potential cause wrt the generalized causal setting  $\mathcal{C}_m = \langle \mathcal{D}_l, m(\vec{\alpha}), m(\Phi) \rangle$  with witness  $m(\Phi_j)$ , and,
- $m(\alpha_k)$  is the unique primary weak potential cause wrt the generalized causal setting  $\mathcal{C}'_m = \langle \mathcal{D}_l, m(\vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1}), \phi'_L \rangle$ , where  $\mathcal{D}_l \cup \mathcal{K} \models \exists s^*. \vec{a}_j. Do(m(\vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1}), S_0, s^*) \wedge Do(m(\alpha_j), s^*, do([\vec{a}_j], s^*))$  and  $\phi'_L = ExecSeq(\vec{a}_j) \wedge After(\vec{a}_j, m(\Phi_j))$ , and,
- $\alpha_j$  is a weak potential cause wrt the causal setting  $\mathcal{C}_h = \langle \mathcal{D}_h, \vec{\alpha}, \Phi \rangle$ ,

then,  $\alpha_k$  is the unique primary weak potential cause wrt the setting  $\mathcal{C}'_h = \langle \mathcal{D}_h, (\vec{\alpha}_{k-1} \alpha_k \vec{\alpha}_{k+1}), ExecSeq(\alpha_j) \wedge After(\alpha_j, \Phi_j) \rangle$ .

Proofs of Theorems 18 and 19 follow a similar reasoning as Theorems 15 and 16, respectively.

**Example (Cont'd).** Let  $\vec{\alpha} = [updateSW_{HL}(Rob, V2021), go_{HL}(Rob, L1), rescue(Rob, L1)]$  and  $\phi_e = \exists r, l. AidedInRescue(r, l)$ . Suppose at the LL,  $\mathcal{D}_l^{ex} \cup \mathcal{K} \models$

$\neg m^{ex}(\phi_e)[S_0] \wedge \exists s. Do(m^{ex}(\vec{\alpha}), S_0, s) \wedge m^{ex}(\phi_e)[s]$ . Moreover, suppose that  $\delta_{rescue}(Rob, L1)$  is the primary weak potential cause wrt setting  $\mathcal{C}_m = \langle \mathcal{D}_l^{ex}, m^{ex}(\vec{\alpha}), m^{ex}(\phi_e) \rangle$ , which brings about the effect  $m^{ex}(\phi_e) = \exists r, l. Cleared(r, l) \vee Evacuated(r, l)$ . Also,  $go_{LL}(Rob, L1)$  is another cause wrt  $\mathcal{C}_m$ .

Then by Theorem 18, we have that  $rescue(Rob, L1)$  is the primary weak potential cause wrt the setting  $\mathcal{C}_h = \langle \mathcal{D}_h^{ex}, (\vec{\alpha}), \phi_e \rangle$ , which brings about the effect  $\phi_e$ . Moreover, by Theorem 19 the action  $go_{HL}(Rob, L1)$  can be considered another weak potential cause wrt  $\mathcal{C}_h$ .  $\triangleleft$

Depending on requirements of the domain, a modeler can decide among sound, complete, or sound and complete abstractions, each providing efficiency benefits.

## Discussion

While there has been a lot of work on actual causation, to the best of our knowledge, our account is the first and the only proposal that investigates programs as actual causes. Perhaps the closest to our work is the one by (Datta et al. 2015), who identified a subset of actions (program steps) of a set of interacting programs as an actual cause for a violation of specific properties in a security domain. Our approach however, focuses on formalizing abstract actual causes as programs in the settings where the actions that led to the observed effect are only incompletely specified. Our framework is based on an expressive logical language for representing and reasoning about dynamic domains. In addition to non-deterministic programs, we allow for incomplete information that is represented through multiple models of a BAT. Furthermore, we investigate how abstraction may be used to facilitate representation and reasoning.

In this paper, we do not study how one can obtain an abstraction given ConGolog programs. Instead, we study causal reasoning that can be accomplished if we are given a sound and/or a complete abstraction of our causal theory. (Luo et al. 2020) proposed forgetting (of LL fluent and action symbols) to obtain a sound and complete abstraction of an LL BAT for a given mapping. Also, (Banihashemi, De Giacomo, and Lespérance 2017) identified the necessary and sufficient conditions for a (given) HL BAT to be a sound abstraction of an LL BAT under a mapping.

For simplicity, we focused on a single layer of abstraction, but the framework supports extending the hierarchy to several levels. In future work, we plan to investigate methodologies for designing abstract theories and refinement mappings with respect to given observed effects, as well as automated synthesis techniques to support this. Extending the current framework to support probabilistic actions (Belle and Levesque 2020) and approximate abstractions, and how such extensions facilitate reasoning about causality are important avenues for future research.

## References

Banihashemi, B.; De Giacomo, G.; and Lespérance, Y. 2017. Abstraction in Situation Calculus Action Theories. In Singh, S. P.; and Markovitch, S., eds., *Proceedings of the Thirty-*



- First AAAI Conference on Artificial Intelligence, 1048–1055. AAAI Press.
- Batusov, V.; and Soutchanski, M. 2018. Situation Calculus Semantics for Actual Causality. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, 1744–1752. AAAI Press.
- Belle, V.; and Levesque, H. J. 2020. Regression and progression in stochastic domains. *Artificial Intelligence*, 281: 103247.
- Datta, A.; Garg, D.; Kaynar, D. K.; Sharma, D.; and Sinha, A. 2015. Program Actions as Actual Causes: A Building Block for Accountability. In Fournet, C.; Hicks, M. W.; and Viganò, L., eds., *IEEE 28th Computer Security Foundations Symposium, CSF 2015, Verona, Italy, 13-17 July, 2015*, 261–275. IEEE Computer Society.
- De Giacomo, G.; Lespérance, Y.; and Levesque, H. J. 2000. ConGolog, a concurrent programming language based on the situation calculus. *Artif. Intell.*, 121(1-2): 109–169.
- De Giacomo, G.; Lespérance, Y.; and Muise, C. J. 2012. On supervising agents in situation-determined ConGolog. In *International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2012*, 1031–1038. IFAAMAS.
- De Giacomo, G.; Lespérance, Y.; Patrizi, F.; and Sardiña, S. 2016. Verifying ConGolog Programs on Bounded Situation Calculus Theories. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, 950–956. AAAI Press.
- Halpern, J. Y. 2016. *Actual Causality*. MIT Press. ISBN 978-0-262-03502-6.
- Khan, S. M.; and Lespérance, Y. 2021. Knowing Why - On the Dynamics of Knowledge about Actual Causes in the Situation Calculus. In Dignum, F.; Lomuscio, A.; Endriss, U.; and Nowé, A., eds., *AAMAS '21: 20th International Conference on Autonomous Agents and Multiagent Systems, Virtual Event, United Kingdom, May 3-7, 2021*, 701–709. ACM.
- Luo, K.; Liu, Y.; Lespérance, Y.; and Lin, Z. 2020. Agent Abstraction via Forgetting in the Situation Calculus. In *ECAI 2020 - 24th European Conference on Artificial Intelligence*, volume 325 of *Frontiers in Artificial Intelligence and Applications*, 809–816. IOS Press.
- McCarthy, J.; and Hayes, P. J. 1969. Some Philosophical Problems From the Standpoint of Artificial Intelligence. *Machine Intelligence*, 4: 463–502.
- Milner, R. 1989. *Communication and concurrency*. PHI Series in computer science. Prentice Hall.
- Mo, P.; Li, N.; and Liu, Y. 2016. Automatic Verification of Golog Programs via Predicate Abstraction. In Kaminka, G. A.; Fox, M.; Bouquet, P.; Hüllermeier, E.; Dignum, V.; Dignum, F.; and van Harmelen, F., eds., *ECAI 2016 - 22nd European Conference on Artificial Intelligence, 29 August-2 September 2016, The Hague, The Netherlands - Including Prestigious Applications of Artificial Intelligence (PAIS 2016)*, volume 285 of *Frontiers in Artificial Intelligence and Applications*, 760–768. IOS Press.
- Pearl, J. 2000. *Causality: Models, Reasoning, and Inference*. Cambridge University Press.
- Reiter, R. 2001. *Knowledge in Action. Logical Foundations for Specifying and Implementing Dynamical Systems*. The MIT Press.