

1 Explaining Actual Causation via Reasoning about 2 Actions and Change

3 **Emily C. LeBlanc**

4 College of Computing and Informatics
5 Drexel University
6 [Philadelphia, PA]
7 leblanc@drexel.edu

8 — Abstract —

9 In causality, an actual cause is often defined as an event responsible for bringing about a given
10 outcome in a scenario. In practice, however, identifying this event alone is not always sufficient
11 to provide a satisfactory explanation of how the outcome came to be. In this paper, we motivate
12 this claim using well-known examples and present a novel framework for reasoning more deeply
13 about actual causation. The framework reasons over a scenario and domain knowledge to identify
14 additional events that helped to “set the stage” for the outcome. By leveraging techniques from
15 Reasoning about Actions and Change, the approach supports reasoning over domains in which
16 the evolution of the state of the world over time plays a critical role and enables one to identify
17 and explain the circumstances that led to an outcome of interest. We utilize action language
18 \mathcal{AC} for defining the constructs of the framework. This language lends itself quite naturally to an
19 automated translation to Answer Set Programming, using which, reasoning tasks of considerable
20 complexity can be specified and executed. We speculate that a similar approach can also lead to
21 the development of algorithms for our framework.

22 **2012 ACM Subject Classification** Knowledge representation and reasoning, Causal reasoning
23 and diagnostics, Temporal reasoning

24 **Keywords and phrases** Actual Cause, Explanation, Reasoning about Actions and Change, Action
25 Language, Answer Set Programming, Knowledge Representation and Reasoning

26 **Digital Object Identifier** 10.4230/OASISs.ICLP.2018.16

27 **1 Introduction and Problem Description**

28 The comprehensive goal of this research has been to design, evaluate, and implement a novel
29 causal reasoning framework to discover causal explanations that are in closer agreement
30 with what common sense might lead one to conclude. Identifying actual causation concerns
31 determining how a specified consequence came to be in a given scenario and has long been
32 studied in a diversity of fields, including law, philosophy, and, more recently, computer science.
33 Also referred to as *causation in fact*, actual causation is a broad term that encompasses all
34 possible antecedents that have played a meaningful role in producing the consequence [5].
35 Consider the well-known Yale Shooting problem [16]:

36 *Shooting a turkey with a loaded gun will kill it. Suzy loads the gun and then shoots*
37 *the turkey. Why is the turkey dead?*

38 Intuition tells us that Suzy’s shooting of the turkey is the *actual cause* of its death. However,
39 if we know for certain that the gun was not loaded at the start of the story, then it is also
40 important to recognize that Suzy’s loading the gun played a key role in producing this
41 consequence. On the other hand, if the gun was loaded from the start, then this point may
42 not be as significant. Moreover, if we build upon this example to say that Tommy handed



© E. C. LeBlanc;
licensed under Creative Commons License CC-BY

Technical Communications of the 34th International Conference on Logic Programming (ICLP 2018).

Editors: Alessandro Dal Palu’, Paul Tarau, Neda Saeedloei, and Paul Fodor; Article No. 16; pp. 16:1–16:10

OpenAccess Series in Informatics



OASIS Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

43 Suzy the gun at the start of the scenario, then surely we want to identify Tommy’s action as
 44 a contributory cause of the turkey’s death. Hall [11] gives another classic example of actual
 45 causation in which two actors have each thrown a rock at a bottle and we wish to determine
 46 which actor’s throw caused the bottle to break. It is easy to imagine similar extensions to the
 47 example that require deeper reasoning about causation to properly explain how the bottle
 48 broke – for example, did a third actor instruct the original two to throw their rocks in the first
 49 place? Literature examples aside, sophisticated actual causal reasoning has been prevalent in
 50 human society and continues to have an undeniable impact on the advancement of science,
 51 technology, medicine, and other important fields. From the development of ancient tools to
 52 modern root cause analysis in business and industry, reasoning about causal influence in a
 53 historical sequence of events enables us to diagnose the cause of an outcome of interest and
 54 gives us insight into how to bring about, or even prevent, similar outcomes in future scenarios.
 55 Consider problems such as explaining the occurrence of a set of suspicious observations in a
 56 monitoring system, reasoning about the efficiency actions taken in an emergency evacuation
 57 scenario, or verifying how an automatically generated workflow produces the expected results.
 58 It is easy to imagine that in cases such as these, determining surface-level causation (e.g.,
 59 Suzy shot the turkey) may not be sufficient to provide a satisfactory explanation of how an
 60 outcome of interest to be.

61 In this dissertation work, we claim that reasoning about actual causation in complex
 62 scenarios requires the ability to identify more than the existence of a causal relationship.
 63 We may want a deeper understanding of the causal mechanism – was the outcome caused
 64 directly or indirectly? Did previously occurring events somehow *support* the causing event or
 65 the outcome’s ability to be caused? To this end, the overall goal of the dissertation work is to
 66 investigate and demonstrate the suitability of action language and answer set programming
 67 to design and realize a novel approach to automated reasoning about actual causation as
 68 described above. The framework leverages techniques from Reasoning about Actions and
 69 Change (RAC) to support reasoning over domains that change over time in response to a
 70 sequence of events, as well as to answer queries for detailed causal explanations of an *outcome*
 71 *of interest* in a specific scenario. The language of choice for the formalization of knowledge
 72 is action language \mathcal{AL} [2] which enables us to represent our knowledge of the direct and
 73 indirect effects of actions in a domain.

74 In the remainder of this summary, we present background on the action language \mathcal{AL}
 75 and its semantics, provide an overview of the framework and its behavior on a novel actual
 76 causation scenario, survey existing literature, and finally discuss open issues and expected
 77 achievements for the dissertation.

78 **2 Preliminaries**

79 As we have already described, this work leverages techniques from Reasoning about Actions
 80 and Change [20] to support reasoning over domains that change over time. We assume that
 81 knowledge of a domain exists as a set of causal laws called an *action description* describing
 82 direct and indirect effects of actions using the action language \mathcal{AL} [2]. These causal laws
 83 embody a transition diagram describing all possible world states of the domain and the
 84 events that trigger transitions between them. In the thesis investigation, we assume the
 85 existence of knowledge in this form, and while the work describes the formalization of the
 86 domain descriptions, the matter of the origin of knowledge is beyond the scope of the thesis.

87 The syntax of \mathcal{AL} builds upon an alphabet consisting of a set \mathcal{F} of symbols for *fluents* and

88 a set \mathcal{E} of symbols for *events*¹. The \mathcal{AL} is centered around a discrete-state-based representation
89 of the evolution of the domain.

90 Fluents are boolean properties of the domain whose truth value may change over time. A
91 (*fluent*) *literal* is a fluent f or its negation $\neg f$. Additionally, we define $\bar{f} = \neg f$ and $\overline{\neg f} = f$.

92 A statement of the form

$$93 \quad e \text{ causes } l_0 \text{ if } l_1, l_2, \dots, l_n \quad (1)$$

94 is called *dynamic causal law*, and intuitively states that, if event e in \mathcal{E} occurs in a state
95 in which literals l_1, \dots, l_n hold, then l_0 , the *consequence* of the law, will hold in the next
96 state. A statement

$$97 \quad l_0 \text{ if } l_1, \dots, l_n \quad (2)$$

98 is called *state constraint* and says that, in any state in which l_1, \dots, l_n hold, l_0 also holds.
99 This second kind of statement allows for an elegant and concise representation of indirect
100 effects, which increases the flexibility of the language. Finally, an *executability condition* is a
101 statement of the form:

$$102 \quad e \text{ impossible_if } l_1, \dots, l_n \quad (3)$$

103 where e and l_1, \dots, l_n are as above. (3) states that e cannot occur if l_1, \dots, l_n hold. A set
104 of statements of \mathcal{AL} is called an *action description*. The semantics of an action description
105 AD is defined by its *transition diagram* $\tau(AD)$, a directed graph $\langle N, E \rangle$ such that:

- 106 1. N is the collection of all states of AD ;
- 107 2. E is the set of all triples $\langle \sigma, e, \sigma' \rangle$ where σ, σ' are states, e is an event executable in σ ,
108 and σ, e, σ' satisfy the *successor state equation* [17]:

$$109 \quad \sigma' = Cn_Z(E(e, \sigma) \cup (\sigma \cap \sigma')) \quad (4)$$

110 where Z is the set of all state constraints of AD .

111 The argument of Cn_Z in (4) is the union of the set of direct effects $E(e, \sigma)$ of e , with the
112 set $\sigma \cap \sigma'$ of the facts “preserved by inertia”. The application of Cn_Z adds the “indirect effects”
113 to this union. A triple $\langle \sigma, e, \sigma' \rangle \in E$ is called a *transition* of $\tau(AD)$ and σ' is a *successor*
114 *state of* σ (under e). A sequence $\langle \sigma_1, \alpha_1, \sigma_2, \dots, \alpha_k, \sigma_{k+1} \rangle$ is a *path of* $\tau(D)$ of length k if
115 every $\langle \sigma_i, \alpha_i, \sigma_{i+1} \rangle$ is a transition in $\tau(D)$. We refer to state σ_1 of a path p as the *initial*
116 *state of* p . A path of length 0 contains only an initial state. In the next section, we build
117 upon this formalization to define a query to our framework for representing and reasoning
118 about actual cause.

119 3 Framework Overview and Foundational Example

120 In this section, we provide an overview of the causal reasoning framework alongside a novel
121 foundational example that showcases the reasoning capabilities and explanatory power of
122 the framework. It is a straightforward scenario in which an outcome of interest, say θ_E , is
123 not satisfied at the start of the scenario. After the occurrence of three events, say e_1, e_2 ,
124 and e_3 , the outcome has been caused. Given the outcome of interest, the sequence of events,
125 and knowledge of the domain in which they have occurred, our framework identifies causal
126 explanations for how θ_E may have come to be. In order to explain actual causation, we will
127 aim to characterize *transition events* which tell us the primary cause of an outcome and
128 whether or not it was caused directly or indirectly, as well as *outcome* and *supporting events*
129 which tell us which prior occurring events have contributed to causing the outcome.

¹ For convenience and compatibility with the terminology from RAC, in this paper we use *action* and *event* as synonyms.

130 **Query**

131 A query consists of an action description, a sequence of events, and the outcome of interest.
 132 The sequence of three scenario events and the outcome of interest for our example are
 133 represented by $v_E = \langle e_1, e_2, e_3 \rangle$, and $\theta_E = \{A, B, C, D, E, F\}$, respectively. The following
 134 action description AD_E characterizes events in the scenario's domain:

$$\begin{array}{l}
 \left\{ \begin{array}{l}
 e_1 \text{ impossible_if } A \\
 e_1 \text{ causes } E \text{ if } \neg E \\
 e_2 \text{ causes } D \text{ if } \neg D \\
 e_3 \text{ causes } A \text{ if } \neg A \\
 e_3 \text{ causes } C \text{ if } \neg C \\
 e_3 \text{ impossible_if } \neg E \\
 e_3 \text{ impossible_if } \neg F \\
 B \text{ if } C
 \end{array} \right. \quad \begin{array}{l}
 (5) \\
 (6) \\
 (7) \\
 (8) \\
 (9) \\
 (10) \\
 (11) \\
 (12)
 \end{array}
 \end{array}$$

136 Laws (5) and (6) describe event e_1 , telling us that e_1 can only occur when A does not
 137 hold and e_1 will cause E if it does not already hold. Law (7) states that e_2 will cause D to
 138 hold if it does not already hold. Similar to causal laws (6) and (7), laws (8) and (9) tell us
 139 that e_3 will cause A and C to hold if they do not hold. The executability conditions (10) and
 140 (11) state that e_3 can only occur when both E and F hold. Finally, the state constraint (12)
 141 tells us that B holds whenever C holds. Given the action description AD_E , the sequence of
 142 events v_E , and the outcome of interest θ_E , the triple $\mathcal{Q}_E = \langle AD_E, v_E, \theta_E \rangle$ is the *query* for
 143 our example. Next, we introduce the concept of a *scenario path*, a unique mapping of the
 144 scenario described by a query to a representation of how the state of the world has changed
 145 in response to the events.

146 **Scenario Path**

147 Scenario paths represent a unique unfolding of a scenario and provide a convenient represent-
 148 ation of how the domain changes over time in response to the events of the scenario. We
 149 reason over these paths to explain actual causation.

150 ► **Definition 1.** Given a query $\mathcal{Q} = \langle AD, v, \theta \rangle$, a *scenario path* is a path $\rho = \langle \sigma_1, \alpha_1, \sigma_2, \dots, \alpha_k,$
 151 $\sigma_{k+1} \rangle$ of $\tau(AD)$ satisfying the following conditions:

- 152 1. $\forall i, 1 \leq i \leq k, \alpha_i = e_i$
- 153 2. $\theta \not\subseteq \sigma_1$
- 154 3. $\exists i, 1 < i \leq k + 1, \theta \subseteq \sigma_i$

155 Condition 1 requires that the events in ρ correspond to the events of v , capturing the
 156 idea that each event of v represents a transition between states in ρ . Condition 2 requires
 157 that the set of fluent literals θ is not satisfied by the initial state of ρ , ensuring that the
 158 outcome has not already been caused prior to the known events of the story. Condition 3
 159 requires that θ is satisfied in at least one state after the initial state in ρ . Conditions 2 and 3
 160 together ensure that at least one event is responsible for causing θ to hold in ρ . The successor
 161 state equation (4) tells us some event in the scenario path must have directly or indirectly
 162 caused θ to be satisfied at some point after the initial state. The set of all scenario paths
 163 with respect to the query \mathcal{Q} is denoted by $P(\mathcal{Q}) = \{\rho_1, \rho_2, \dots, \rho_m\}$.

164 It is clear that there are multiple valid scenario paths in the set $P(\mathcal{Q}_E)$, each representing
 165 a valid evolution of state in response to the scenario's events in the domain given by AD_E .

■ **Table 1** Tabular representation of the scenario path $\rho_E \in P(\mathcal{Q}_E)$.

| State | Event | State Affecting Law(s) |
|--|------------------|--|
| $\sigma_1 = \{\neg A, \neg B, \neg C, \neg D, \neg E, F\}$ | $\alpha_1 = e_1$ | e_1 causes E if $\neg E$ |
| $\sigma_2 = \{\neg A, \neg B, \neg C, \neg D, E, F\}$ | $\alpha_2 = e_2$ | e_2 causes D if $\neg D$ |
| $\sigma_3 = \{\neg A, \neg B, \neg C, D, E, F\}$ | $\alpha_3 = e_3$ | e_3 causes A if $\neg A$, e_3 causes C if $\neg C$, B if C |
| $\sigma_4 = \{A, B, C, D, E, F\}$ | – | – |

166 For the purposes of this discussion, we choose a path with a complex causal mechanism that
 167 will exercise the causal reasoning framework. We will refer to this path as ρ_E . Table 1 shows
 168 the evolution of state in ρ_E in response to the events of v_E . The first column lists each state
 169 σ_i of ρ_E , and the second column gives the event α_i that caused a transition to the state
 170 σ_{i+1} . It is easy to see that ρ_E satisfies the conditions of Definition 1 with respect to AD_E ,
 171 v_E , and θ_E .

172 Transition Event

173 A *transition event* is an event in a scenario path that causes a transition from a state of the
 174 world where the outcome θ is not satisfied to a state of the world where θ is satisfied. In this
 175 section, we identify transition events and their direct and indirect effects on the outcome.

176 ► **Definition 2.** Given a scenario path $\rho = \langle \sigma_1, \alpha_1, \sigma_2, \dots, \alpha_k, \sigma_{k+1} \rangle$ and an outcome θ , event
 177 α_j , where $1 \leq j \leq k$, is a *transition event* of θ in ρ if the following conditions are satisfied by
 178 the transition $\langle \sigma_j, \alpha_j, \sigma_{j+1} \rangle$ of ρ :

- 179 1. $\theta \not\subseteq \sigma_j$
- 180 2. $\theta \subseteq \sigma_{j+1}$

181 Intuitively, event α_j is a transition event of outcome θ if the outcome was not satisfied
 182 when α_j occurred but *was* satisfied after its occurrence. Note that we have defined transition
 183 events in such a way that there can be multiple transition events for θ in ρ . Using Table 1, it
 184 is straightforward to verify that event e_3 is the only transition event of θ_E in the example
 185 scenario path ρ_E , clearly satisfying Conditions 1 and 2 of Definition 2.

186 Given a query $\mathcal{Q} = \langle AD, v, \theta \rangle$, a scenario path $\rho = \langle \sigma_1, \alpha_1, \sigma_2, \dots, \alpha_k, \sigma_{k+1} \rangle$ in $P(\mathcal{Q})$,
 187 and a transition event α_j for θ , the set of *direct effects* of α_j in θ is $d_\theta(\alpha_j, \rho) = \theta \cap E(\alpha_j, \sigma_j)$.
 188 Recall that $E(\alpha_j, \sigma_j)$ is the set of all direct effects of event α_j given that it occurs in state
 189 σ_j . The set of all direct effects of e_3 with respect to σ_3 , then, is $E(e_3, \sigma_3) = \{A, C\}$, in
 190 accordance with laws (8) and (9) in AD_E . The direct effects of e_3 in θ_E , then, is given by
 191 $d_{\theta_E}(e_3, \rho_E) = \theta_E \cap E(e_3, \sigma_3) = \{A, B, C, D, E, F\} \cap \{A, C\} = \{A, C\}$.

192 To determine the indirect effects of an event with respect to the outcome, first let
 193 $S = E(\alpha_j, \sigma_j) \cup (\sigma_j \cap \sigma_{j+1})$ represent the set of all literals directly caused by the transition
 194 event α_j and those preserved by inertia. Given a query $\mathcal{Q} = \langle AD, v, \theta \rangle$, a scenario path
 195 $\rho = \langle \sigma_1, \alpha_1, \sigma_2, \dots, \alpha_k, \sigma_{k+1} \rangle$ in $P(\mathcal{Q})$, and a transition event α_j for θ , the set of *indirect*
 196 *effects* of α_j in θ is $i_\theta(\alpha_j, \rho) = \theta \cap (\sigma_{j+1} \setminus S)$. Given the set $S_E = E(e_3, \sigma_3) \cup (\sigma_3 \cap \sigma_4) =$
 197 $\{A, C\} \cup \{D, E, F\} = \{A, C, D, E, F\}$ representing the direct effects of e_3 and the literals
 198 preserved by inertia, the indirect effects of e_3 in θ_E is

$$\begin{aligned}
 199 \quad i_{\theta_E}(e_3, \rho_E) &= \theta_E \cap (\sigma_4 \setminus S_E) \\
 200 &= \{A, B, C, D, E, F\} \cap (\{A, B, C, D, E, F\} \setminus \{A, C, D, E, F\}) \\
 201 &= \{A, B, C, D, E, F\} \cap \{B\} \\
 202 &= \{B\}
 \end{aligned}$$

204 This result is intuitive because e_3 directly caused C to hold by law (9) and we know from
 205 law (12) that whenever C holds in a certain state, then B holds. We claim that under these
 206 conditions, it must be the case the e_3 caused B indirectly.

207 First Causal Explanation

208 Both the knowledge of the transition event and its effects on the outcome are represented
 209 by the *first causal explanation*. Given the query $\mathcal{Q}_E = \langle AD_E, v_E, \theta_E \rangle$, the scenario path
 210 $\rho_E \in P(\mathcal{Q}_E)$, the transition event e_3 in ρ_E , and e_3 's direct and indirect effects, $d_{\theta_E}(\rho_E, \theta_E)$
 211 and $i_{\theta_E}(\rho_E, \theta_E)$, respectively, the *first causal explanation* for θ_E in ρ_E is the tuple

$$212 \quad C_E^1 = \langle \rho_E, e_3, d_{\theta_E}(\rho_E, \theta_E), i_{\theta_E}(\rho_E, \theta_E) \rangle$$

$$213 \quad = \langle \rho_E, e_3, \{A, C\}, \{B\} \rangle$$

215 Explanation C_E^1 summarizes our initial findings – the event e_3 caused a transition from a
 216 state where the outcome $\{A, B, C, D, E, F\}$ did not hold to a state where it did hold in the
 217 scenario path ρ_E . Specifically, literals A and C were direct effects of e_3 's occurrence while e_3
 218 caused B indirectly.

219 While C_E^1 tells us how the set of literals $\{A, B, C\}$ of θ_E were made to hold in scenario
 220 path ρ_E , we are still missing information about which, if any, events prior to e_3 caused the
 221 remaining literals $\{D, E, F\}$ to hold in state σ_4 . We also do not know if any prior occurring
 222 events influenced e_3 's ability to be a transition event of θ_E . In this work, supporting events
 223 are events that have occurred prior to a transition event α_j that enable α_j to be a transition
 224 event for the outcome θ . We identify two types of supporting events, *outcome supporting*
 225 *event* (OSEs) and *transition supporting events* (TSEs), both which are presented in the
 226 following sections. In order to identify both OSEs and TSEs in a scenario path ρ , we must
 227 first introduce the notion that an event α_i *ensures* that a literal l will hold in a specified
 228 state σ_j if it is the most recent transition event for l .

229 ► **Definition 3.** Given a scenario path $\rho = \langle \sigma_1, \alpha_1, \sigma_2, \dots, \alpha_k, \alpha_{k+1} \rangle$, event α_i is an *ensuring*
 230 *event* of $l \in \sigma_j$ in ρ if:

- 231 1. α_i is a transition event of $\{l\}$ in ρ
- 232 2. $i < j$
- 233 3. $j - i$ is minimal

234 Condition 1 leverages Definition 2 to require that event α_i responsible for l holding in
 235 some state of ρ . Condition 2 requires that α_i occurs before α_j in ρ . Condition 3 requires
 236 that α_i is the most recent transition event of l in ρ . We claim that if no event ensures $l \in \sigma_j$
 237 for a path ρ , this implies that l holds in every state of ρ because there exists no transition
 238 $\langle \sigma_i, \alpha_i, \sigma_{i+1} \rangle$ in the path such that $l \notin \sigma_i$. Therefore, l must have held in the initial state and
 239 was never changed by a subsequent event prior to α_j 's occurrence. Note that because ensuring
 240 events are also transition events, it is straightforward to leverage the characterizations of
 241 direct and indirect effects of transition events from Section 3 to learn if events ensured l in
 242 some state σ due to its direct or indirect effects.

243 Outcome Supporting Events

244 In the case where α_j does not set all of the literals of θ , OSEs can be responsible for ensuring
 245 that these remaining literals hold by the time α_j occurs in ρ . Finding OSEs requires first
 246 identifying if any literals in θ were not set as an effect of the transition event α_j . The set of
 247 remaining literals of an outcome θ is given by $R_\theta = \theta \setminus (d_\theta(\alpha_j, \rho) \cup i_\theta(\alpha_j, \rho))$. If $|R_\theta| > 0$,
 248 then a previously occurring event may have supported the outcome θ by ensuring that the
 249 remaining literals held in state σ_{j+1} .

250 ► **Definition 4.** Given a query \mathcal{Q} , a factual path $\rho \in P(\mathcal{Q})$, a transition event α_j of θ , and a
 251 literal $l \in R_\theta$, α_i is an *outcome supporting event (OSE)* via l if α_i ensures $l \in \sigma_{j+1}$.

252 We denote by O^{supp} the set of OSEs and the literals they ensure. Formally, the tuple
 253 $\langle \alpha_i, l \rangle \in O^{supp}$ if α_i is a OSE via l . We denote by O^{init} the set of literals in R_θ that were
 254 not ensured by an event in ρ . Given a literal $l \in R_\theta$, $l \in O^{init}$ if:

$$255 \quad \neg \exists \langle \alpha, l' \rangle \in O^{supp} \text{ s.t. } l' = l$$

256 Intuitively, a literal l is in O^{init} when l has is no outcome supporting event in O^{supp} . In
 257 our example, we already know that we require additional causal information about the set of
 258 remaining outcome literals D , E , and F . Formally, the following literals in the outcome θ_E
 259 have not been explained by C_E^1 :

$$\begin{aligned} 260 \quad R_{\theta_E} &= \theta_E \setminus (d_{\theta_E}(e_3, \rho_E) \cup i_{\theta_E}(e_3, \rho_E)) \\ 261 \quad &= \{A, B, C, D, E, F\} \setminus (\{A, C\} \cup \{B\}) \\ 262 \quad &= \{A, B, C, D, E, F\} \setminus \{A, C, B\} \\ 263 \quad &= \{D, E, F\} \end{aligned}$$

265 Because $|R_{\theta_E}| > 0$, there is more causal information to uncover. As covered in the earlier
 266 discussion on ensuring events, each literal in R_{θ_E} must either be ensured to hold in state
 267 σ_4 by an outcome supporting event or the literal has held consistently from the start of the
 268 scenario. Event e_2 is an outcome supporting event because it ensures that literal D held in
 269 σ_4 . This event meets the three conditions of ensuring $D \in \sigma_4$. First, it is a transition event of
 270 $\{D\}$ because the literal D did not hold in state σ_2 but it did hold in σ_3 after e_2 's occurrence.
 271 It clearly satisfies Conditions 2 because here $i = 2$ and $j = 4$, and so $i < j$. Finally, it
 272 satisfies Condition 3 because event e_i is the most recent transition event of $\{D\}$, and so
 273 $j - i$ is minimal. Similarly, it is straightforward to verify that e_1 is an outcome supporting
 274 event by ensuring that E holds in state σ_4 . The set of outcome supporting events is given
 275 by $O_E^{supp} = \{\langle e_2, D \rangle, \langle e_1, E \rangle\}$. Finally, the set $O_E^{supp} = \{F\}$ because there exists no tuple
 276 $\langle \alpha, F \rangle \in O_E^{supp}$, and so F must have held in the initial state of ρ_E and never changed value.

277 Second Causal Explanation

278 Knowledge of outcome supporting events and remaining outcome literals that held from the
 279 start is represented by the *second causal explanation*. Given the query $\mathcal{Q}_E = \langle AD_E, v_E, \theta_E \rangle$,
 280 the scenario path $\rho_E \in P(\mathcal{Q}_E)$, and the transition event e_3 for θ_E , the *second causal*
 281 *explanation* for θ_E in ρ_E is

$$\begin{aligned} 282 \quad C_E^2 &= \langle O_E^{supp}, O_E^{init} \rangle \\ 283 \quad &= \langle \{\langle e_2, D \rangle, \langle e_1, E \rangle\}, \{F\} \rangle \end{aligned}$$

285 Explanation C_E^2 provides us with information about how the remaining outcome literals
 286 $\{D, E, F\} \in \theta_E$ came to hold in the state σ_4 . Of these remaining literals, D and E were
 287 ensured by events e_2 and e_1 , respectively. The remaining literal F held in the initial state
 288 and was not ensured in σ_4 by any event prior to e_1 .

289 C_E^2 tells us how the remaining outcome literals came to hold in σ_4 , but there is even
 290 more causal information to be revealed in this example. Next, we discuss an approach to
 291 determining if any other events in scenario path ρ_E contributed to e_3 's ability to be a
 292 transition event of θ_E .

293 **Transition Supporting Events**

294 TSEs ensure that the preconditions of α_j are satisfied in state σ_j so that α_j could occur and
 295 cause θ to be satisfied in σ_{j+1} . The approach to identifying TSEs is conveniently similar to
 296 identifying outcome supporting events, and so we will omit the majority of technical details
 297 in favor of working out the example in the interest of space. To determine whether or not
 298 any prior events supported the transition event e_3 , we begin by identifying all preconditions
 299 for e_3 's occurrence and its ability to produce its effects in ρ_E . We obtain α_j 's *preconditions*
 300 in ρ by reasoning over the of laws in AD . In the dissertation work, we introduce notation to
 301 allow reasoning over the components of laws in an action description AD . For example, given
 302 a dynamic causal law λ in AD of form (1), let $e(\lambda) = e$, $c(\lambda) = l_0$, and $p(\lambda) = \{l_1, l_2, \dots, l_n\}$.
 303 We denote by $\mathcal{D}(AD)$ the set of all dynamic causal laws in AD . We use a similar representation
 304 for executability conditions, and we introduce a set of conditions under which preconditions
 305 can be extracted from these laws. In our example, the literals $\neg A$ and $\neg C$ are in $prec(e_3, \rho_E)$
 306 because of laws (8) and (9) in the action description AD_E . By our definition of precondition,
 307 the literals E and F are also in $prec(e_3, \rho_E)$ because of laws (10) and (11) in AD_E . Therefore,
 308 the set of preconditions of e_3 in ρ_E is $prec(e_3, \rho_E) = \{\neg A, \neg C, E, F\}$.

309 Similar to our definition of outcome supporting events, a *transition supporting event* is the
 310 most recent transition event for a precondition of the transition event. It is straightforward
 311 to verify that the set of transition supporting events is given by $T_E^{supp} = \langle e_1, E \rangle$ and the set
 312 of initially set literals is $T_E^{init} = \{\neg A, \neg C, F\}$.

313 **Third Causal Explanation**

314 Knowledge of transition supporting events and precondition literals that held from the start
 315 is represented by the *third causal explanation*. Given the scenario path $\rho_E \in P(Q_E)$, the
 316 transition event e_3 , the set of transition supporting events T_E^{supp} , and the set of uncaused
 317 literals T_E^{init} the *third causal explanation* for θ_E in ρ_E is

$$318 \quad C_E^3 = \langle T_E^{supp}, T_E^{init} \rangle$$

$$319 \quad = \langle \{ \langle e_1, E \rangle \}, \{ \neg A, \neg C, F \} \rangle$$

321 Explanation C_E^3 tells us about the transition event e_3 's preconditions and how they were
 322 met by state σ_3 . The precondition literals of event e_3 were $\neg A$, $\neg C$, E , and F . Of these
 323 precondition literals, E was ensured in σ_3 by the occurrence of event e_1 . The remaining
 324 literals $\neg A$, $\neg C$, and F were not ensured in σ_3 by any scenario event. For relative brevity,
 325 we will not query further for details about the outcome and transition supporting events. It
 326 is easy to see, however, that the framework could tell us that the precondition literal E for
 327 e_3 was made to hold as a *direct effect* of e_1 's occurrence.

328 **Actual Causal Explanation**

329 As the research intends to prove, there exists a space of possible structures for causal
 330 explanation. Recall that when there are remaining outcome literals to explain, there is
 331 a second causal explanation. However, if a transition event has no preconditions in the
 332 scenario path, then there is no third causal explanation. This implies that the structure of
 333 the explanation depends on the information encoded by the corresponding scenario path.
 334 We intend to characterize this space of structures in the dissertation. The framework can
 335 identify all three causal explanations in our example (i.e., C_E^1 , C_E^2 , and C_E^3). To summarize,
 336 the framework has explained that e_3 was a transition event for θ_E through both direct and
 337 indirect effects, e_1 and e_2 were outcome supporting events, and e_1 was a transition supporting
 338 event in the scenario path ρ_E .

339 **4 Overview of Existing Literature**

340 While actual causation has been treated in numerous ways in the Artificial Intelligence
341 literature, the most relevant of which we will cover briefly in this section, existing approaches
342 do not possess the fine-granularity of reasoning and explanation required to meet the reasoning
343 needs of the examples discussed here. Many approaches to reasoning about actual cause have
344 been inspired by the human intuition that cause can be determined by hypothesizing about
345 whether or not a removing X from a scenario would prevent Y from being true [19]. Attempts
346 to mathematically characterize actual causation have largely pursued counterfactual analysis
347 of structural equations [22, 13, 15], neuron diagrams [12], and other logical formalisms
348 [18, 23, 4]. It has been widely documented, however, that the counterfactual criteria alone
349 is problematic and fails to recognize causation in some common cases such as preemption,
350 overdetermination, and contributory cause [21, 10]. More recent approaches such as [14] have
351 addressed some of these shortcomings by modifying the existing definitions of actual cause or
352 by modeling change over time with some improved results. However, there is still no widely
353 agreed upon counterfactual definition of actual cause in spite of a considerably large body of
354 work aiming to find one.

355 The work of [3] departs from the counterfactual approach, using a similar insight to our
356 own that actual causation can be determined by inspecting a specific scenario. Leveraging the
357 Situation Calculus (SC) to formalize knowledge, the approach uses a single step regression
358 approach to identify events deemed relevant to a logical statement becoming true. Although
359 the conceptual approach is similar to our own, the technical approaches differ significantly.
360 For example, [3] identifies a single sequence of causal events without explanation. There
361 are also ramifications due to the choices for the formalization of the domain. Compared to
362 \mathcal{AL} formalizations, SC formalizations incur limitations when it comes to the representations
363 of indirect effects of actions, which play an essential role in our work, and the elaboration
364 tolerance of the formalization. Additionally, SC relies on First-Order Logic, while \mathcal{AL} features
365 an independent and arguably simpler semantics.

366 **5 Open Issues and Expected Achievements**

367 While the core of this framework is fairly well-developed at this stage, there remain some
368 open issues that will be addressed in the dissertation. Evaluation of the framework is
369 a crucial next step, and meaningful progress has been made towards demonstrating the
370 framework's reasoning process when solving examples from causality literature in addition
371 to novel scenarios. We expect to demonstrate that the framework can solve numerous
372 classic examples with finer-grained causal explanations than the current state of the art.
373 Moreover, the dissertation will present a number of empirical studies to compare and evaluate
374 the ability of related approaches to solve the novel example presented in this paper. We
375 expect that related approaches will not be able to explain the causal mechanism of our
376 example in comparable detail. The dissertation will also present a novel set of identified open
377 problems whose investigation can advance the capabilities of the causal reasoning framework.
378 Regarding implementation, the choice of \mathcal{AL} as the underlying formalism has useful practical
379 implications. As demonstrated by a substantial body of literature (see, e.g., [1]), \mathcal{AL} lends
380 itself quite naturally to an automated translation to Answer Set Programming [8, 9], using
381 which, complex reasoning tasks can be specified and executed (see, e.g., [6, 7]). We speculate
382 that a similar approach can also lead to the development of algorithms for our framework,
383 and have begun translating \mathcal{AL} queries, scenario paths, and transition events to ASP.

384 — References

- 385 1 Marcello Balduccini and Michael Gelfond. Diagnostic reasoning with a-prolog. *arXiv*
 386 *preprint cs/0312040*, 2003.
- 387 2 Chitta Baral and Michael Gelfond. Reasoning agents in dynamic domains. In *Logic-based*
 388 *artificial intelligence*, pages 257–279. Springer, 2000.
- 389 3 Vitaliy Batusov and Mikhail Soutchanski. Situation calculus semantics for actual causality.
 390 In *13th International Symposium on Commonsense Reasoning. University College London,*
 391 *UK. Monday, November*, volume 6, 2017.
- 392 4 Sander Beckers and Joost Vennekens. A general framework for defining and extending
 393 actual causation using cp-logic. *International Journal of Approximate Reasoning*, 77:105–
 394 126, 2016.
- 395 5 Charles E Carpenter. Concurrent causation. *University of Pennsylvania Law Review and*
 396 *American Law Register*, 83(8):941–952, 1935.
- 397 6 Thomas Eiter, Wolfgang Faber, Nicola Leone, Gerald Pfeifer, and Axel Polleres. Answer
 398 set planning under action costs. *Journal of Artificial Intelligence Research*, 19:25–71, 2003.
- 399 7 Esra Erdem, Michael Gelfond, and Nicola Leone. Applications of answer set programming.
 400 *AI Magazine*, 37(3), 2016.
- 401 8 Michael Gelfond and Vladimir Lifschitz. The stable model semantics for logic programming.
 402 In *ICLP/SLP*, volume 88, pages 1070–1080, 1988.
- 403 9 Michael Gelfond and Vladimir Lifschitz. Classical negation in logic programs and disjunct-
 404 ive databases. *New generation computing*, 9(3-4):365–385, 1991.
- 405 10 Clark Glymour and David Danks. Actual causation: a stone soup essay. *Synthese*,
 406 175(2):169–192, 2010.
- 407 11 Ned Hall. Two concepts of causation. *Causation and counterfactuals*, pages 225–276, 2004.
- 408 12 Ned Hall. Structural equations and causation. *Philosophical Studies*, 132(1):109–136, 2007.
- 409 13 Joseph Y Halpern. Axiomatizing causal reasoning. *Journal of Artificial Intelligence Re-*
 410 *search*, 12:317–337, 2000.
- 411 14 Joseph Y Halpern. *Actual causality*. MIT Press, 2016.
- 412 15 Joseph Y Halpern and Judea Pearl. Causes and explanations: A structural-model approach.
 413 part i: Causes. *The British journal for the philosophy of science*, 56(4):843–887, 2005.
- 414 16 Steve Hanks and Drew McDermott. Nonmonotonic logic and temporal projection. *Artificial*
 415 *intelligence*, 33(3):379–412, 1987.
- 416 17 Patrick J. Hayes and John McCarthy. Some Philosophical Problems from the Standpoint of
 417 Artificial Intelligence. In B. Meltzer and D. Michie, editors, *Machine Intelligence 4*, pages
 418 463–502. Edinburgh University Press, 1969.
- 419 18 Mark Hopkins and Judea Pearl. Causality and counterfactuals in the situation calculus.
 420 *Journal of Logic and Computation*, 17(5):939–953, 2007.
- 421 19 David Lewis. Causation. *The journal of philosophy*, 70(17):556–567, 1974.
- 422 20 J. McCarthy and P. J. Hayes. Some philosophical problems from the standpoint of artificial
 423 intelligence. *Readings in artificial intelligence*, pages 431–450, 1969.
- 424 21 Peter Menzies. Counterfactual theories of causation. 2001.
- 425 22 Judea Pearl. On the definition of actual cause. 1998.
- 426 23 Joost Vennekens. Actual causation in cp-logic. *Theory and Practice of Logic Programming*,
 427 11(4-5):647–662, 2011.