

# Efficient Discovery of Actual Causality Using Abstraction Refinement

Arshia Rafieioskouei<sup>1b</sup> and Borzoo Bonakdarpour<sup>1b</sup>

**Abstract**—Causality is the relationship where one event contributes to the production of another, with the cause being partly responsible for the effect and the effect partly dependent on the cause. In this article, we propose a novel and effective method to formally reason about the causal effect of events in engineered systems, with application for finding the root-cause of safety violations in embedded and cyber-physical systems. We are motivated by the notion of *actual causality* by Halpern and Pearl, which focuses on the causal effect of particular events rather than type-level causality, which attempts to make general statements about scientific and natural phenomena. Our first contribution is formulating discovery of actual causality in computing systems modeled by transition systems as an satisfiability modulo theory solving problem. Since datasets for causality analysis tend to be large, in order to tackle the scalability problem of automated formal reasoning, our second contribution is a novel technique based on *abstraction refinement* that allows identifying for actual causes within smaller abstract causal models. We demonstrate the effectiveness of our approach (by several orders of magnitude) using three case studies to find the actual cause of violations of safety in 1) a neural network controller for a mountain car; 2) a controller for a Lunar Lander obtained by reinforcement learning; and 3) an MPC controller for an F-16 autopilot simulator.

**Index Terms**—Causality, cyber-physical systems (CPSs), root-cause analysis, safety failures.

## I. INTRODUCTION

**I**N A CAUSAL system, the output of the system is influenced only by the present and past inputs. In other words, in a causal system, the present and future outputs depend solely on past and present inputs, not on future inputs. Causality addresses the logical dependencies between events and reflects the essence of event and action flows in systems. Engineers generally build causal systems, that is, structures, systems, and processes that seek to tie effects to their causes. This also includes approaches to explain the root-cause of failures that violate safety standards, especially in safety-critical systems.

In this context, embedded and cyber-physical systems (CPSs) are no exceptions. In fact, root-cause analysis has been of interest to both academic and industrial circles for decades, aiming not just to find safety violations but also to precisely

explain why they happened. This means proving mathematically that safety would not have been violated in the absence of the identified cause. Formalizing and reasoning about causal explanations is much harder than just finding “bugs” and often aims to identify the earliest flawed decisions by controllers that lead to violations of safety requirements. Finding such causes provides engineers with tremendous insights to design more reliable systems, but it has been a long-standing and very challenging problem for various reasons, from defining a formal definition of causal effect of events to the high-computational complexity of counterfactual reasoning.

There is a wealth of research on causality analysis in the context of embedded and component-based systems from different perspectives [1], [2], [3], [4], [5], [6], [7], [8], [9]. Recently, there has been great interest in using temporal logics to reason about causality and explain bugs [10], [11], [12], [13]. In the CPS domain, using causality to repair AI-enabled controllers has recently gained interest [14], [15]. However, these lines of work either focus on only modeling aspects of causality or do not address the problem of scalability in automated reasoning about causality, which inherently involves a combinatorial blow up to enumerate counterfactuals.

### Objectives

This article is concerned with the following problem. Given are 1) a formal operational description of a computing system  $\mathcal{T}$  (e.g., a transition system of a CPS) and 2) a logical predicate  $\varphi_e$  that describes the effect (e.g., a safety failure) as input. Our goal is to identify a predicate  $\varphi_c$  that describes the cause of  $\varphi_e$  happening (e.g., the earliest bad decision made by a controller). We note that  $\varphi_e$  can be given by the user or can be found by using a verification or testing technique. Hence, the way the effect is identified is irrelevant to the problem studied in this article.

The first natural step is to formalize the definition of causality and in fact, there are several interpretations of the meaning of causality. In this article, we are motivated by the notion of *actual causality* by Halpern and Pearl (HP) [16], which focuses on the causal effect of particular events, rather than type causality, which attempts to make general statements about scientific and natural phenomena (e.g., smoking causes cancer). Actual causality is a formalism to deal with *token-level* causality, which aims to find the causal effect of individual events (in our context, in embedded and CPS), as

Manuscript received 7 August 2024; accepted 9 August 2024. This work was supported by the United States NSF Award CCF under Grant 2320050. This article was presented at the International Conference on Embedded Software (EMSOFT) 2024 and appeared as part of the ESWEEK-TCAD special issue. This article was recommended by Associate Editor S. Dailey. (Corresponding author: Borzoo Bonakdarpour.)

The authors are with Michigan State University, East Lansing, MI 48824 USA (e-mail: rafieios@msu.edu; borzoo@msu.edu).

Digital Object Identifier 10.1109/TCAD.2024.3448299

75 opposed to *type-level* causality, which intends to generalize  
76 the causal effect of types of events.

77 As we aim at analyzing executions of computing systems  
78 (e.g., models or data logs of a CPS), we first formalize causal  
79 models in (possibly infinite-state) *transition systems*, rather  
80 than the classic set of structural equations [16]. We show that  
81 formalizing the three conditions of actual causality yields a  
82 second-order logic formula of the form  $\varphi_{\text{hp}} \triangleq \exists \tau. \exists \tau'. \forall \tau''. \psi$ ,  
83 where  $\tau$ ,  $\tau'$ , and  $\tau''$  range over the set of executions of a  
84 transition system and  $\psi$  stipulates the relation between actual  
85 and counterfactual worlds. More specifically.

- 86 1) The outermost existential quantifier in  $\varphi_{\text{hp}}$  intends to  
87 establishes a relationship between the cause and the  
88 effect in an execution  $\tau$  (known as the *AC1* necessity  
89 condition in the HP framework [16]). That is, the cause  
90 and then the effect *actually* happen in  $\tau$ .
- 91 2) The inner existential quantifier aims at refuting the  
92 causal effect relation in the counterfactual world (known  
93 as the *AC2(a)* condition, stipulating the “but-for” condi-  
94 tion under contingencies). That is, when the cause does  
95 not happen, the effect will not happen.
- 96 3) The universal quantifier requires that if the cause hap-  
97 pens in any execution  $\tau''$  that is to  $\tau$  (as far as  
98 the variables contributing to the cause and effect are  
99 concerned), then the effect should also happen (known  
100 as the *AC2(b)* sufficiency condition).

101 This formula exhibits a quantifier alternation and indeed, the  
102 problem of deciding actual causality in a causal model is  
103 known to be DP-complete [17] in the size of the model,  
104 illustrating the computational complexity of the problem. To  
105 deal with this complexity, we propose an effective method to  
106 formally reason about actual causality using decision proce-  
107 dures to solve satisfiability modulo theory (SMT). Although  
108 there has been tremendous progress in developing efficient  
109 SMT solvers, they may not scale well when dealing with very  
110 large causal models or data logs. To tackle this problem, we  
111 introduce a novel technique based on *abstraction refinement*  
112 that allows identifying causes within smaller abstract causal  
113 models. This abstraction simplifies the model and attempts  
114 to view it from a higher level, while preserving the causal  
115 relations.

116 Although the idea of abstracting causal models in terms of  
117 structural equations has been studied in [18], [19], and [20],  
118 these works develop an exact simulation which may not exist  
119 or do not attempt to establish a relation between actual causes  
120 in the abstract and concrete causal models. Our technique  
121 incorporates two levels of abstraction to reason about actual  
122 causality (i.e., formula  $\varphi_{\text{hp}} \triangleq \exists \tau. \exists \tau'. \forall \tau''. \psi$ ). More specifi-  
123 cally, our approach works as follows (see Figs. 1 and 2). Given  
124 a concrete causal model  $\mathcal{T}$ .

- 125 1) We first compute an under-approximate model  $\tilde{\mathcal{T}}$  of  $\mathcal{T}$ .  
126 This model is used to find witnesses for conditions  
127 *AC1* and *AC2(a)* (i.e., the existential quantifiers). If not  
128 successful, we refine  $\tilde{\mathcal{T}}$  (e.g., by including states that are  
129 in  $\mathcal{T}$  and not in  $\tilde{\mathcal{T}}$ ) and try again.
- 130 2) If the previous step succeeds, we compute an over-  
131 approximate model  $\hat{\mathcal{T}}$  to verify condition *AC2(b)* for the  
132 universal quantifier. If successful, then an actual cause is

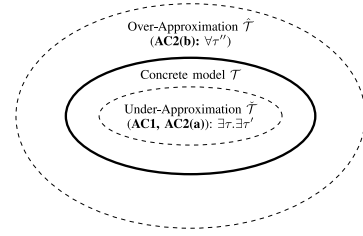


Fig. 1. Over/under-approximations of the concrete model and their relation to HP conditions of the form  $\exists \exists \forall$ .

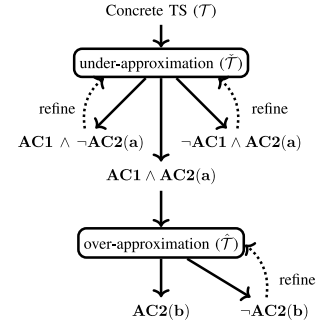


Fig. 2. Overall idea of our algorithm—steps of abstraction-refinement approach.

133 identified and the algorithm terminates. Otherwise, we  
134 refine  $\hat{\mathcal{T}}$  (e.g., by excluding states that are in  $\hat{\mathcal{T}}$  and not  
135 in  $\mathcal{T}$ ) and repeat the second step.<sup>1</sup>

136 We prove the correctness of our approach by showing that our  
137 algorithm is sound (but not necessarily complete).

138 We have implemented<sup>2</sup> our approach using the Python  
139 programming language and utilized the SMT solver Z3 [21]  
140 and data analysis libraries [22], [23] to construct our solver  
141 and abstraction technique. We conduct experiments on three  
142 case studies to find the actual cause of violations of safety  
143 in 1) a neural network controller for a mountain car [24];  
144 2) a controller for a Lunar Lander obtained by reinforcement  
145 learning [24]; and 3) an MPC controller for an F-16 autopilot  
146 simulator [25]. Our experiments demonstrate the effectiveness  
147 of our abstraction-refinement technique by several orders of  
148 magnitude compared to the SMT-based approach for concrete  
149 causal models.

150 In summary, our contributions are the following. We: 1) formu-  
151 late the classic HP framework by transition systems and  
152 introduce an SMT-based decision procedure to identify actual  
153 causes in a computing system; 2) introduce a technique based  
154 on abstraction refinement to deal with scalability of formal  
155 reasoning about actual causality; and 3) conduct three rigorous  
156 experimental evaluations on AI-enabled as well as non-AI  
157 controllers in CPS.

158 Our work is the first step in automating discovery of actual  
159 cause of failures, and our experiments show that we are able  
160 to identify the earliest bad decisions by controllers that lead  
161 to violations of safety requirements.

<sup>1</sup>Alternatively, one can return to the first step and start from scratch.

<sup>2</sup>Source code and all trace logs available at [https://github.com/TART-MSU/Causality\\_abs\\_refinement](https://github.com/TART-MSU/Causality_abs_refinement).

162 *Organization:* The remainder of this article is organized  
 163 as follows. Section II presents the classic HP framework for  
 164 actual causality. In Section III, we introduce our formulation  
 165 of HP for transition systems as well as a translation to an  
 166 SMT-based decision procedure to identify actual causes. Our  
 167 abstraction-refinement technique is introduced in Section IV.  
 168 We present our experimental evaluation in Section V. Related  
 169 work is discussed in Section VI. Finally, we make concluding  
 170 remarks and discuss future work in Section VII. Proofs of  
 171 correctness are available in [26].

## 172 II. PRELIMINARIES—ACTUAL CAUSALITY

173 In this section, we present the notion of *actual causality*  
 174 by HP [16] as the baseline preliminary concept. Since, the  
 175 definition in [16] is not a natural model of computation, in  
 176 Section III, we will adapt the concepts in this section to  
 177 transition systems and second-order logic formulas in order to  
 178 reason about actual causality in computing systems. We will  
 179 consistently use the *Mountain Car* running example to explain  
 180 the definitions and concepts throughout this article.

### 181 A. Causal Models

182 *Definition 1:* A *signature*  $\mathcal{S}$  is a tuple  $(\mathcal{U}, \mathcal{V}, \mathcal{R})$ , where  $\mathcal{U}$   
 183 is set of *exogenous* variables (variables that represent factors  
 184 outside the control of the model),  $\mathcal{V}$  is a set of *endogenous*  
 185 variables (variables whose values are ultimately determined by  
 186 the values of the endogenous and exogenous variables).  $\mathcal{R}$  is  
 187 a function that associates with every variable  $Y \in \mathcal{U} \cup \mathcal{V}$  a  
 188 nonempty set  $\mathcal{R}(Y)$  of possible values for  $Y$ .

189 Following Definition 1, a *state* is a valuation of a vector  
 190 of variables  $\vec{X} = (X_1, \dots, X_n)$  in  $\mathcal{U} \cup \mathcal{V}$ , where each variable  
 191  $X \in \vec{X}$  is assigned a value from  $\mathcal{R}(X)$ .

192 *Definition 2:* A *basic causal model*  $\mathcal{M}$  is a pair  $(\mathcal{S}, \mathcal{F})$ ,  
 193 where  $\mathcal{S}$  is a signature and  $\mathcal{F}_X$  defines a function that asso-  
 194 ciates with each endogenous variable  $X$  a *structural equation*  
 195  $\mathcal{F}_X$  that maps  $\mathcal{R}(\mathcal{U} \cup \mathcal{V} - \{X\})$  to  $\mathcal{R}(X)$ , so  $\mathcal{F}_X$  determines the  
 196 value of  $X$ , given the values of all the other variables in  $\mathcal{U} \cup \mathcal{V}$ .

197 It is important to highlight that exogenous variables cannot  
 198 be linked to a function; thus, assigning values to exogenous  
 199 variables, denoted as  $\vec{u}$ , is referred to as a *context*.

200 *Definition 3:* An *intervention* entails setting the values of  
 201 endogenous variables, denoted as  $\vec{X} \leftarrow \vec{x}$ , and this notation  
 202 signifies that the variables within set  $\vec{X}$  are assigned values  
 203  $\vec{x} = (x_1, \dots, x_n)$ .

204 The structural equations define what happens in the presence  
 205 of interventions. Setting the value of some variables  $\vec{X}$  to  $\vec{x}$   
 206 in a causal model  $\mathcal{M} = (\mathcal{S}, \mathcal{F})$  results in a new causal model,  
 207 denoted  $\mathcal{M}_{\vec{X} \leftarrow \vec{x}}$ , which is identical to  $\mathcal{M}$ , except that  $\mathcal{F}$   
 208 is replaced by  $\mathcal{F}^{\vec{X} \leftarrow \vec{x}}$ : for each variable  $Y \notin \vec{X}$ ,  $\mathcal{F}_Y^{\vec{X} \leftarrow \vec{x}} = \mathcal{F}_Y$   
 209 while for each  $X'$  in  $\vec{X}$ , the equation  $\mathcal{F}_{X'}$  is replaced by  $X' =$   
 210  $x'$ . Thus, we define a *causal model*  $\mathcal{M}$  by a tuple  $(\mathcal{S}, \mathcal{F}, \mathcal{I})$ ,  
 211 where  $(\mathcal{S}, \mathcal{F})$  is a basic causal model (see Definition 2) and  
 212  $\mathcal{I}$  is a set of *allowed interventions*. Following [20], the sets  
 213 of “allowed interventions” ensure that the interventions can  
 214 be appropriately limited to include only those that can be  
 215 abstracted.

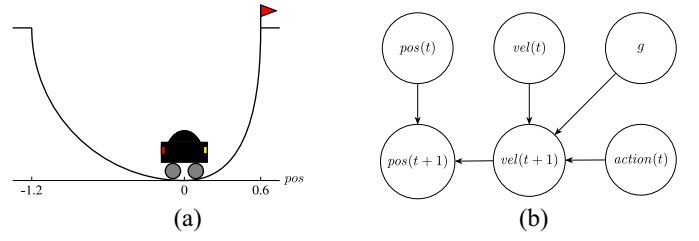


Fig. 3. (a) Schematic of the mountain car example. (b) Graph illustrating the causal model and relationships between the variables at a snapshot in time  $t$ .

216 *Example 1:* Consider a car located in a valley and aiming to  
 217 reach the top of a mountain [see Fig. 3(a)]. At each time step,  
 218 the controller of the car determines whether to apply positive  
 219 or negative acceleration to guide the car toward the mountain  
 220 top. We define signature  $\mathcal{S} = (\mathcal{U}, \mathcal{V}, \mathcal{R})$  for this example as  
 221 follows. Let

$$222 \quad \mathcal{U} = \{\text{pos}(0), \text{vel}(0), g\}$$

223 be the set of exogenous variables, denoting the initial posi-  
 224 tion, initial velocity, and the gravitational force on the car,  
 225 respectively. Let

$$226 \quad \mathcal{V} = \{\text{pos}(t), \text{vel}(t), \text{action}(t)\}$$

227 be the set of endogenous variables, denoting the position,  
 228 velocity, and the controller action, respectively, at each time  
 229 step  $t$ , where  $t \neq 0$ . We also set

$$230 \quad \mathcal{R}(\text{pos}(t)) = [-1.2, 0.6]$$

$$231 \quad \mathcal{R}(\text{vel}(t)) = [-0.07, 0.07]$$

$$232 \quad \mathcal{R}(\text{action}(t)) = \{-1, 0, 1\}$$

233 where  $-1$ ,  $0$ , and  $1$  are assigned as accelerate to the left, do  
 234 not accelerate, and accelerate to the right, respectively, for all  
 235  $t \geq 0$ .

236 Now, we define the causal model  $(\mathcal{S}, \mathcal{F})$  based on the  
 237 system dynamics for each  $t > 0$  by structural equations

$$238 \quad \mathcal{F}_{\text{pos}(t+1)} = \mathcal{F}_{\text{pos}(t)} + \mathcal{F}_{\text{vel}(t)} \quad (1)$$

$$239 \quad \mathcal{F}_{\text{vel}(t+1)} = \mathcal{F}_{\text{vel}(t)} + 0.001\mathcal{F}_{\text{action}(t)} - g \cdot \cos(3\mathcal{F}_{\text{pos}(t)}). \quad (2)$$

240 To illustrate the dependencies of the system, we can use  
 241 a causal graph, as shown in Fig. 3(b). In this model,  
 242  $\mathcal{M}_{\text{action}(t) \leftarrow 1}$  denotes the model obtained by an intervention,  
 243 where the action( $t$ ) is set to 1 at time  $t$  (for some  $t > 0$ ).

### 244 B. Causal Formulas

245 To precisely define actual causality, formal language is  
 246 essential for articulating causal statements with clarity and  
 247 rigor, in particular to formalize causes and effects. We use an  
 248 extension of propositional logic, wherein primitive events take  
 249 the form  $\vec{X} = \vec{x}$ , representing an endogenous variable  $\vec{X}$  and  
 250 a possible value  $\vec{x}$  for  $\vec{X}$ . The combination of primitive events  
 251 is achieved through standard propositional connectives, such  
 252 as  $\{\wedge, \vee, \neg\}$ . Thus, in this article, we are only concerned with  
 253 causal formulas that are state-based (and not temporal).

254 Given a signature  $\mathcal{S} = (\mathcal{U}, \mathcal{V}, \mathcal{R})$ , a *primitive event* is a  
 255 formula of the form  $X = x$ , for  $X \in \mathcal{V}$  and  $x \in \mathcal{R}(X)$ . A *causal*

256 *formula* (over  $\mathcal{S}$ ) is one of the form  $[Y_1 \leftarrow y_1, \dots, Y_k \leftarrow$   
 257  $y_k]\varphi$ , where  $\varphi$  is Boolean combination of primitive events,  
 258  $Y_1, \dots, Y_k$  are distinct variables in  $\mathcal{V}$ , and  $y_i \in \mathcal{R}(Y_i)$ . Such a  
 259 formula is abbreviated as  $[\vec{Y} \leftarrow \vec{y}]\varphi$ . The special case where  
 260  $k = 0$  is abbreviated as  $[\ ]\varphi$  or, more often, just  $\varphi$ . Intuitively,  
 261  $[\vec{Y} \leftarrow \vec{y}]\varphi$  says that  $\varphi$  would hold if  $Y_i$  were set to  $y_i$ , for  
 262  $i = 1, \dots, k$ .

263 A causal formula  $\psi$  is true or false in a causal model,  
 264 given a context. We use a pair  $(\mathcal{M}, \vec{u})$  consisting of a causal  
 265 model  $\mathcal{M}$  and context  $\vec{u}$  as a *causal setting*. Hence, we  
 266 write  $(\mathcal{M}, \vec{u}) \models \psi$  if the causal formula  $\psi$  is true in the  
 267 causal setting  $(\mathcal{M}, \vec{u})$ . We are restricted to recursive models,  
 268 where given a context, no cyclic dependencies exists. In a  
 269 recursive model,  $(\mathcal{M}, \vec{u}) \models X = x$  if the value of  $X$  is  $x$   
 270 once we set the exogenous variables to  $\vec{u}$ . Given a model  
 271  $\mathcal{M}$ , the model that describes the result of this intervention  
 272 is  $\mathcal{M}_{\vec{Y} \leftarrow \vec{y}}$ . Thus,  $(\mathcal{M}, \vec{u}) \models [\vec{Y} \leftarrow \vec{y}]\psi$  iff  $(\mathcal{M}_{\vec{Y} \leftarrow \vec{y}}, \vec{u}) \models$   
 273  $\psi$ . Mathematical formalism serves to express the intuition  
 274 precisely encapsulated within the formula  $[\vec{Y} \leftarrow \vec{y}]\psi$  is true  
 275 in a causal setting  $(\mathcal{M}, \vec{u})$  exactly if the formula  $\psi$  is true  
 276 in the model that results from the intervention, in the same  
 277 context  $\vec{u}$ .

278 *Example 2:* Context  $\vec{u}$  in causal setting  $(\mathcal{M}, \vec{u})$  in our  
 279 example is determined by system inputs: initial velocity, initial  
 280 position, and gravity

$$281 \quad \vec{u} = \left\{ (\text{vel}(0) \leftarrow 0.01), (\text{pos}(0) \leftarrow 0), (g \leftarrow 0.0025) \right\}$$

282 where we defined  $\mathcal{M}$  in Example 1. To conduct causal  
 283 analysis, the car at time  $t = 0$  decides to set  $\text{action}(0) = 1$ ,  
 284 but it fails to reach the goal. We defined causal formula to  
 285 express failure as follows:

$$286 \quad \varphi_{\text{fail}} \triangleq (\text{pos}(n) \neq 0.6)$$

287 where 0.6 is the flag position and  $n$  is the last car state.

### 288 C. Actual Causality

289 *Definition 4:*  $\vec{X} \leftarrow \vec{x}$  is an *actual cause* of  $\varphi$  in causal  
 290 setting  $(\mathcal{M}, \vec{u})$ , if the following three conditions hold.

- 291 1) *AC1:*  $(\mathcal{M}, \vec{u}) \models [\vec{X} \leftarrow \vec{x}]\varphi$ .
- 292 2) *AC2(a):* There is a partition of  $\mathcal{V}$  (set of endogenous  
 293 variables) into two disjoint subsets  $\vec{Z}$  and  $\vec{W}$  (i.e.,  $\vec{Z} \cap \vec{W} =$   
 294  $\emptyset$ ) with  $\vec{X} \subseteq \vec{Z}$  and a setting  $\vec{x}'$  and  $\vec{w}$  of the variables  
 295 in  $\vec{X}$  and  $\vec{W}$ , respectively, such that

$$296 \quad (\mathcal{M}, \vec{u}) \models [\vec{X} \leftarrow \vec{x}', \vec{W} \leftarrow \vec{w}] \neg \varphi.$$

- 297 3) *AC2(b):* For all subsets  $\vec{Z}'$  of  $\vec{Z} - \vec{X}$ , we have

$$298 \quad (\mathcal{M}, \vec{u}) \models \left[ \vec{X} \leftarrow \vec{x}, \vec{W} \leftarrow \vec{w}, \vec{Z}' \leftarrow \vec{z}^* \right] \varphi$$

299 where  $\vec{z}^*$  denotes that variables in  $\vec{Z}'$  are fixed at their  
 300 values in the actual context.

- 301 4) *AC3:*  $\vec{X}$  is minimal; no subset of  $\vec{X}$  satisfies AC1 and  
 302 AC2.

303 Roughly speaking Definition 4 expresses the following.  
 304 AC1 says that  $\vec{X} = \vec{x}$  cannot be considered a cause of  $\varphi$   
 305 unless both  $\vec{X} = \vec{x}$  and  $\varphi$  actually happen. AC2(a) says that the  
 306 but-for condition holds under the contingency  $\vec{W} = \vec{w}$ . Also,

changing the value of some variable in  $\vec{X}$  results in changing 307  
 the value(s) of some variable(s) in  $\vec{Z}$  (perhaps recursively), 308  
 which finally results in the truth value of  $\varphi$  changing. Finally, 309  
 AC2(b) provides a sufficiency condition: if the variables in  $\vec{X}$  310  
 and an arbitrary subset  $\vec{Z} - \vec{X}$  of other variables on the causal 311  
 path are held at their values in the actual context, then  $\varphi$  holds 312  
 even if  $\vec{W}$  is set to  $\vec{w}$  (the setting for  $\vec{W}$  used in AC2(a)). 313  
 The types of events that the HP definition allows as actual 314  
 causes are ones of the form  $X_1 = x_1 \wedge \dots \wedge X_k = x_k$ , that is, 315  
 conjunctions of primitive events; this is often abbreviated as 316  
 $\vec{X}$ . In Section III, Example 3, we will provide an example on 317  
 how actual cause of formula  $\varphi_{\text{fail}}$  can be identified using our 318  
 proposed technique. 319

### 320 III. SMT-BASED DISCOVERY OF ACTUAL CAUSALITY

321 In this section, we transform the components of the HP 321  
 framework presented in Section II into transition systems and 322  
 a second-order formula to express actual causality. Such a 323  
 transition system can model the operational behavior of a 324  
 system (e.g., a controller). Our technique can be agnostic to the 325  
 details of the system and only take a set of execution traces. 326

327 Recall that a causal model  $\mathcal{M}$  is of the form  $(\mathcal{S}, \mathcal{F}, \mathcal{I})$ ,  
 328 where  $\mathcal{S} = (\mathcal{U}, \mathcal{V}, \mathcal{R})$ . Also, recall that a state is a mapping  
 329 from the variables in  $\mathcal{U} \cup \mathcal{V}$  to their respective domain of values.  
 330 We start with representing  $\mathcal{M}$  with a set of traces obtained  
 331 from a transition system that essentially describes how the  
 332 state of all variables in  $\mathcal{U} \cup \mathcal{V}$  evolve over time by structural  
 333 equations  $\mathcal{F}_X$ , for every  $X \in \mathcal{V}$ .

#### 334 A. Transition Systems

335 *Definition 5:* A *transition system*  $\mathcal{T}$  corresponding to a 335  
 causal model  $\mathcal{M}$  is a tuple  $\mathcal{T} = (\Sigma, \Delta, \sigma^0, \Lambda)$ , where: 1)  $\Sigma$  336  
 is the set of all possible states obtained from all possible 337  
 valuations of variables in  $\mathcal{U} \cup \mathcal{V}$ ; 2)  $\Delta$  is a function mapping 338  
 states in  $2^\Sigma$  to a state in  $\Sigma$  (recall that  $\mathcal{F}_X$  is a function); 339  
 3)  $\sigma^0 \in \Sigma$  is the initial state, and 4)  $\Lambda$  is a function mapping 340  
 states in  $2^\Sigma$  to an atomic proposition from a set **AP** (e.g., 341  
 given by causal formulas). 342

343 Following Definition 5, given a causal setting  $(\mathcal{M}, \vec{u})$ , the 343  
 corresponding *causal transition system* is one that is acyclic 344  
 and fixes  $\sigma_0$  according to  $\vec{u}$ . An intervention  $\vec{X} \leftarrow \vec{x}$  is simply 345  
 a set of transitions in  $\Delta$  where in the target state  $\vec{X} = \vec{x}$  346  
 holds, denoted by  $\Delta_{\vec{X} \leftarrow \vec{x}}$ . We denote the set of all possible 347  
 interventions in  $\mathcal{T}$  by  $\mathcal{I}_{\mathcal{T}}$ . 348

349 *Definition 6:* A *path* of a transition systems  $\mathcal{T} =$   
 350  $(\Sigma, \Delta, \sigma^0, \Lambda)$  is a sequence of states of form  $\sigma_0 \sigma_1 \dots$ , where  
 351 for all  $i \geq 0$  1)  $\sigma_0 = \sigma^0$  and 2)  $(\sigma_i, \sigma_{i+1}) \in \Delta$ . The  
 352 *trace* corresponding to a path  $\sigma_0 \sigma_1 \dots$  is the sequence  $\tau =$   
 353  $\Lambda(\sigma_0) \Lambda(\sigma_1) \dots$

354 Let  $\text{Tr}$  denote the set of all traces of a transition system.

355 *Example 3:* Fig. 4 shows three traces  $\tau_0$ ,  $\tau_1$ ,  
 356 and  $\tau_2$  for our mountain car example for context  
 357  $\vec{u} = (\text{pos}(0) = 0.0, \text{vel}(0) = 0.02)$ . In each step, the controller  
 358 makes acceleration decisions. Dotted transitions means the  
 359 next state is not the immediate next time step. The  $n$ th state  
 360 is the last state of the trace. As can be seen, traces  $\tau_0$  and  $\tau_2$   
 361 never reach position 0.6 (i.e., satisfying causal formula  $\varphi_{\text{fail}}$ ,

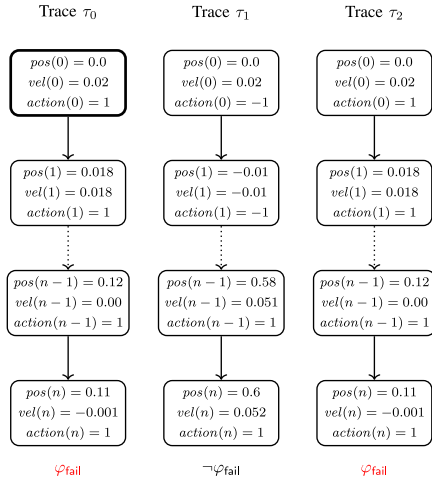


Fig. 4. Three traces for the mountain car example.

362 meaning failing to reach the flag), while trace  $\tau_1$  does (i.e.,  
 363 violating causal formula  $\varphi_{\text{fail}}$ , meaning successfully reaching  
 364 the flag).

365 We introduce three temporal operators to express the occur-  
 366 rence of causes and effects in traces: 1) for a state  $\sigma$  and a  
 367 proposition  $p \in \mathbf{AP}$  iff  $\sigma \models p$  iff  $p \in \Lambda(\sigma)$ ; 2) a trace  
 368  $\tau = \tau_0\tau_1\dots$  satisfies formula  $\Box p$  (read as ‘‘always  $p$ ’’ and  
 369 denoted  $\tau \models \Box p$ ) iff  $\forall i \geq 0. \tau_i \models p$ ; 3) a trace  $\tau_0\tau_1\dots$   
 370 satisfies formula  $\Diamond p$  (read as ‘‘eventually  $p$ ’’ and denoted  $\tau \models$   
 371  $\Diamond p$ ) iff  $\exists i \geq 0. \tau_i \models p$ ; and 4) a trace  $\tau_0\tau_1\dots$  satisfies  
 372 formula  $p \mathcal{U} q$  (read as ‘‘ $p$  until  $q$ ’’ denoted  $\tau \models p \mathcal{U} q$ ) iff  
 373  $\exists i \geq 0. (\tau_i \models q \wedge (\forall j < i. \tau_j \models p))$ .

### 374 B. SMT-Based Formulation of Actual Causality

375 An SMT decision problem generally consists of two com-  
 376 ponents: 1) the SMT instance (i.e., data elements, such as  
 377 variables, domains, functions, sets, etc.) and 2) SMT con-  
 378 straints (i.e., first-order modulo theory involving quantified  
 379 Boolean predicates with arithmetic). In the context of our  
 380 problem, the SMT instance consists of two parts:

- 381 1) A set of elements for expressing a transition system  $\mathcal{T}$   
 382 or 2) a set of traces  $\text{Tr}$  (e.g., from a data log). While the  
 383 latter is simply a set of sequences of states (defined as a  
 384 function from natural numbers to the full set of states),  
 385 the former is specified by Boolean formulas from the  
 386 unrolled transition system, similar to standard bounded  
 387 model checking [27] without loops.
- 388 2) Our SMT model formalize conditions  $AC1$ ,  $AC2(a)$ ,  
 389 and  $AC2(b)$  of Definition 4 for transition systems (see  
 390 Fig. 5). For simplification, we omit  $AC3$  (minimality  
 391 of cause), as it is not the most important constraint  
 392 to reason about causal effect of events in a system.  
 393 Condition  $AC1$  (in Fig. 5) means in the set  $\text{Tr}$ , there  
 394 exists at least one trace  $\tau$ , where effect  $\varphi_e$  appears after  
 395 cause  $\varphi_c$  holds. Condition  $AC2(a)$  requires the existence  
 396 of one trace  $\tau'$ , where neither cause  $\varphi_c$  nor effect  $\varphi_e$   
 397 hold.

398 Additionally, trace  $\tau'$  is not equivalent to trace  $\tau$  (identified  
 399 in  $AC1$ ) as far as variables in  $W$  or  $Z$  are concerned (i.e., the  
 400 counterfactual worlds). The remaining endogenous variables,

the ones in  $\bar{W}$ , are off to the side, so to speak, but may still  
 have an indirect effect on what happens. Condition  $AC2(b)$   
 requires that for all traces  $\tau''$  that are similar to  $\tau$  as far as  
 causal variables in  $Z$  are concerned, if cause  $\varphi_c$  holds, then  
 effect  $\varphi_e$  hold some time in the future.

We clarify that while SMT solvers cannot directly encode  
 temporal operators, one can easily encode them using the  
 above expanded definitions by quantifiers over traces.

#### SMT Decision Problem

Given are 1) a causal transition system  $(\mathcal{T}, \vec{u})$  (or a  
 set of traces  $\text{Tr}$  expressed as a mapping from natural  
 numbers to states); 2) a causal formula  $\varphi_e$ ; 3) an unin-  
 terpreted function representing  $\varphi_c$ ; and 4) constraints  
 $AC1$ ,  $AC2(a)$ , and  $AC2(b)$ . The corresponding SMT  
 instance is satisfiable iff the interpreted  $\varphi_c$  is an  
 actual cause of  $\varphi_e$  in  $\mathcal{T}$ .

Example 4: We aim to identify the cause of the failure,  
 denoted as  $\varphi_{\text{fail}}$ , explained in our running example. For the  
 sake of argument, let  $X = \{\text{action}(0) = 1\}$ . Since both  $\text{pos}$   
 and  $\text{vel}$  are dependent on the value of  $\text{action}$ , they are part of  
 $\bar{Z}$  or the causal path. That is

$$\bar{Z} = \{\text{pos}(t), \text{action}(t), \text{vel}(t) \mid t > 1\}$$

and, hence,  $W = \{\}$  (since  $\bar{W} \cap \bar{Z} = \emptyset$ ). We now analyze the  
 conditions of HP.

- 1) Starting with  $AC1$ , one can instantiate  $\tau$  (in Fig. 5) with  
 concrete trace  $\tau_0$  in Fig. 4, indicating the satisfaction of  
 the first condition.
- 2) Moving to  $AC2(a)$ , which involves counterfactual rea-  
 soning, when we change the actual setting in  $AC1$  to a  
 counterfactual value  $\text{action}(0) = -1$ , the car eventually  
 reaches the goal (i.e.,  $\text{pos} = 0.6$ ). This change allows  
 the car to initiate a leftward movement, acquiring the  
 necessary momentum to reach the flag, so flipping  
 the failure  $\varphi_{\text{fail}}$  to success (i.e.,  $\neg\varphi_{\text{fail}}$ ). Consequently,  
 $AC2(a)$  is satisfied by instantiating  $\tau'$  (in Fig. 5) with  
 concrete trace  $\tau_1$  (in Fig. 4). Also, notice that condition  
 $\tau_1 \not\equiv_Z \tau_0$  is satisfied.
- 3) Considering  $AC2(b)$ , notice that trace  $\tau_2$  is identical to  
 $\tau_0$  as far as the variables in  $\bar{Z}$  are concerned (i.e.,  $\tau_0 \equiv_Z$   
 $\tau_2$ ). Also, since  $\bar{W} = \{\}$ , changing variables in  $\bar{W}$  while  
 preserving the actual context results in an equivalent  
 scenario to  $AC1$ , which is already satisfied. Thus, the  
 only trace that can instantiate  $\tau''$  (in Fig. 5) is  $\tau_2$ , in  
 which  $\varphi_{\text{fail}}$  becomes true. Note that the reason  $\tau_0$  and  $\tau_2$   
 are trace-equivalent is indeed due to the fact that  $\bar{W} = \{\}$ .  
 Hence,  $AC2(b)$  hold.

This means in this set of traces,  $\text{action}(0) = 1$  is the actual  
 cause of failure for the car to reach the flag.

In the ideal world, one has to have all possible traces for  
 combinatorial enumeration to evaluate  $AC2(b)$ . However, this  
 is far from reality and most trace data logs (e.g., by some test-  
 ing mechanisms, fuzzing, mutation testing, some automaton,  
 etc.) include only a subset of possibilities. Our goal in this  
 article is to identify causal effects *within* a given set of traces.

$$\begin{aligned}
\mathbf{AC1} &\triangleq \exists \tau \in \text{Tr.} (\tau \models \neg \varphi_e \mathcal{U} (\varphi_c \wedge \diamond \varphi_e)) && (\varphi_c \text{ causes } \varphi_e \text{ in } \tau) \\
\mathbf{AC2(a)} &\triangleq \exists \tau' \in \text{Tr.} (\tau' \models \square (\neg \varphi_c \wedge \neg \varphi_e)) \wedge (\tau \not\equiv_Z \tau' \vee \tau \not\equiv_W \tau') && (\text{changes in the causal inhibits } \varphi_e) \\
\mathbf{AC2(b)} &\triangleq \forall \tau'' \in \text{Tr.} ((\tau'' \models (\neg \varphi_e \mathcal{U} \varphi_c) \wedge (\tau \equiv_Z \tau'' \wedge \tau \not\equiv_W \tau'')) \rightarrow (\tau'' \models \diamond \varphi_e)) && (\text{in traces similar to } \tau, \varphi_c \text{ causes } \varphi_e)
\end{aligned}$$

Fig. 5. HP conditions adapted for causal transition systems.

448 Finally, as mentioned in the introduction, decision procedure  
449 for verification of actual causality is **DP-complete** [17],  
450 signifying the computation difficulty of automated reasoning  
451 about causality. This means our SMT-based problem is indeed  
452 dealing with a decision problem that is **DP-complete**, setting  
453 the complexity of our SMT-based solution.

#### 454 IV. ABSTRACTION REFINEMENT FOR CAUSAL MODELS

455 In this section, we propose our abstraction-refinement tech-  
456 nique and its application in reasoning about actual causality,  
457 as presented in Section III.

##### 458 A. Overall Idea

459 Generally speaking, the traditional abstraction approach  
460 to handle an existential quantifier is *under-approximation*,  
461 where we start from a subset of behaviors and attempt to  
462 instantiate the quantifier. If successful, then the problem is  
463 solved. Otherwise, we refine the abstraction, by including  
464 addition behaviors and try again. On the contrary, to handle  
465 universal quantifiers, the traditional abstraction approach is  
466 *over-approximation*, where we start from a subset of behaviors  
467 and attempt to verify universality. If successful, then the  
468 problem is solved. Otherwise, we need to ensure that the  
469 counterexample is not *spurious* (due to over-approximation). If  
470 it is, we refine the abstraction by excluding the counterexample  
471 and try again.

472 The overall idea of our technique is as follows (see Fig. 2).  
473 Observe that the logic formula for actual causality is of the  
474 form  $\exists \exists \forall$  (see Fig. 5). Given a transition system  $\mathcal{T}$  and causal  
475 formula  $\varphi_e$  as the effect, we proceed as follows.

- 476 1) *Step 1*: Compute an under-approximation  $\check{\mathcal{T}}$  and an  
477 over-approximation  $\hat{\mathcal{T}}$ . We first attempt to instantiate  
478 the existential quantifiers in *AC1* and *AC2(a)* in  $\check{\mathcal{T}}$ . If  
479 instantiating one of the quantifiers does not succeed, we  
480 refine  $\check{\mathcal{T}}$  and repeat step 1.
- 481 2) *Step 2*: When step 1 succeeds, we compute  $\hat{\mathcal{T}}$  and verify  
482 the universal quantifier in *AC2(b)* for  $\hat{\mathcal{T}}$ . If successful,  
483 the witness to  $\tau$  is a trace where the actual cause happens  
484 and we also obtain a witness to  $\varphi_c$  by the SMT solver.  
485 Otherwise, we can either refine  $\hat{\mathcal{T}}$  and repeat step 2 or  
486 refine  $\check{\mathcal{T}}$  and return to step 1.

487 We show that termination of these steps results in identify-  
488 ing an actual cause  $\varphi_c$  in  $\mathcal{T}$  for  $\varphi_e$ . This algorithm, however,  
489 may never terminate and, thus, our approach is sound but  
490 not complete. We also remark that our heuristic based on  
491 abstraction refinement is sound but not complete (e.g., similar

to the CEGAR [28] technique in model checking) to solve the  
general DP-complete problem. The computation complexity  
of our solution, therefore, does not change.

##### 495 B. Approximating Causal Transition Systems

496 We first fix some notation. For a *concrete* causal transition  
497 system  $\mathcal{T} = (\Sigma, \Delta, \sigma^0, \Lambda)$  (the one given as input for causal  
498 reasoning), let us denote an over-approximate causal transition  
499 system by  $\hat{\mathcal{T}} = (\hat{\Sigma}, \hat{\Delta}, \hat{\sigma}^0, \hat{\Lambda})$  and an under-approximate  
500 causal transition system by  $\check{\mathcal{T}} = (\check{\Sigma}, \check{\Delta}, \check{\sigma}^0, \check{\Lambda})$ . We denote  
501 the domain of endogenous (respectively, exogenous) variables  
502 of  $\mathcal{T}$  by  $\mathcal{R}(\mathcal{V}_{\mathcal{T}})$  (respectively,  $\mathcal{R}(\mathcal{U}_{\mathcal{T}})$ ).

503 Given an over-approximate causal transition system  $\hat{\mathcal{T}}$ , we  
504 construct a sequence  $\hat{\mathcal{T}}_0 \geq \hat{\mathcal{T}}_1 \geq \dots \hat{\mathcal{T}}_k$  of over-approximations,  
505 where (1)  $\hat{\mathcal{T}}_k = \hat{\mathcal{T}}$ , and  $\hat{\mathcal{T}}_{i+1}$  is a refinement of  $\hat{\mathcal{T}}_i$ , for  
506  $0 \leq i < k$ , which we compute using *counterexamples*. A  
507 counterexample is a state of  $\hat{\Sigma}_i$  that is not in  $\Sigma$ . Over-  
508 approximation state mapping is a function which map states  
509 from  $\mathcal{T}$  to  $\hat{\mathcal{T}}$ , i.e.,  $\hat{h}: 2^{\Sigma} \mapsto \hat{\Sigma}$ .

510 *Assumption 1*: In this article, we only allow over-  
511 approximation state mappings  $\hat{h}$  that preserve the equality of  
512 traces as far as variables in  $Z$  are concerned. That is, for two  
513 concrete transitions  $(\sigma_0, \sigma_1)$  and  $(\sigma'_0, \sigma'_1)$ , if 1)  $\sigma_0 \equiv_Z \sigma'_0$  and  
514 2)  $\sigma_1 \not\equiv_Z \sigma'_1$ , then we have 1)  $\sigma_0 \equiv_Z \hat{h}(\sigma'_0)$  and 2)  $\sigma_1 \not\equiv_Z$   
515  $\hat{h}(\sigma'_1)$ . Otherwise, we will not be able to prove the soundness  
516 of Algorithm 1 with regard to causal paths. We will elaborate  
517 more in the requirement in proof of Theorem 2. We will also  
518 explain in Section V, how this assumption is ensured in our  
519 implementation

520 We need an additional function:  $\hat{w}: \mathcal{I}_{\mathcal{T}} \mapsto \mathcal{I}_{\hat{\mathcal{T}}}$  which maps  
521 concrete interventions to over-approximation interventions.

522 *Definition 7*: Given a subset of endogenous variables in  $\Sigma$ ,  
523 called  $\vec{X}$ , and  $\vec{x} \in 2^{\Sigma}$ , let

$$524 \text{Rst}(\Sigma, \vec{x}) = \{\vec{v} \in 2^{\Sigma} : \vec{x} \text{ is the restriction of } \vec{v} \text{ to } \vec{X}\}.$$

525 This definition carries to a transition system  $\mathcal{T} =$   
526  $(\Sigma, \sigma^0, \Delta, \Lambda)$  in a straightforward fashion as follows. The  
527 restriction of a set of values  $\vec{x}$  on  $\Sigma$  is a subset  $\Sigma|_{\vec{x}} \subseteq \Sigma$   
528 restricted to those states, where  $\vec{X} = \vec{x}$ . The set of restricted  
529 transitions is obviously those start and end in states in  $\Sigma|_{\vec{x}}$ .

530 We now explain how we compute the above functions.  
531 Given  $\hat{h}$ , we define:  $\hat{w}(\Delta_{\vec{X} \leftarrow \vec{x}}) = \hat{\Delta}_{\vec{Y} \leftarrow \vec{y}}$  if 1)  $\vec{y} \in 2^{\hat{\Sigma}}$  and  
532 2)  $\hat{h}(\text{Rst}(\Sigma|_{\vec{x}})) = \text{Rst}(\hat{\Sigma}|_{\vec{y}})$ . Hence, for every intervention in  
533  $\Delta_{\vec{X} \leftarrow \vec{x}}$ , there is only one intervention in  $\hat{\Delta}_{\vec{Y} \leftarrow \vec{y}}$ . If such a  $\vec{Y}$   
534 and  $\vec{y}$  do not exist, we take  $\hat{w}(\Delta_{\vec{X} \leftarrow \vec{x}})$  to be *undefined*. Let

---

**Algorithm 1: Finding Actual Cause of  $\varphi_e$  in  $\mathcal{T}$** 


---

**Input:**  $\mathcal{T} = (\Sigma, \Delta, \sigma^0, \Lambda)$ , causal formula  $\varphi_e$ , allowed interventions  $\mathcal{I}_{\mathcal{T}}^h$ ,  $\alpha = [0, 1]$ ,  $\beta \geq 0$   
**Output:** Causal formula  $\varphi_c$

```

1  $\check{\text{Tr}} \leftarrow \check{h}(\text{Tr})$  using  $\alpha$ ;
2 while true do
3    $\{\varphi_c, \check{\tau}, \check{\tau}'\} \leftarrow \text{SMT}(\check{\text{Tr}}, \text{AC1} \wedge \text{AC2(a)});$ 
4   if  $\neg\varphi_c$  then
5      $\hat{\mathcal{T}} \leftarrow \hat{h}(\mathcal{T})$  using  $\beta$  and  $\varphi_c$ ;
6     while true do
7       result  $\leftarrow \text{SMT}(\hat{\mathcal{T}}, \varphi_c, \text{AC2(b)});$ 
8       if result then
9         return  $\varphi_c$ ;
10      else
11         $\hat{\mathcal{T}} \leftarrow \text{Refine}(\hat{\mathcal{T}}, \Sigma - \hat{\Sigma}, \mathcal{I}_{\hat{\mathcal{T}}}^h);$ 
12      end
13    end
14  end
15  Increase  $\alpha$ ;
16   $\check{\text{Tr}} \leftarrow \check{h}(\text{Tr})$  using  $\alpha$ ;
17 end

```

---

$\mathcal{I}_{\hat{\mathcal{T}}}^h$  be the set of interventions for which  $\hat{w}$  is defined, and let  $\mathcal{I}_{\hat{\mathcal{T}}} = \hat{w}(\mathcal{I}_{\hat{\mathcal{T}}}^h)$ .

Based on this definition, it becomes evident that not all interventions in  $\mathcal{I}_{\mathcal{T}}$  will have corresponding mappings in  $\mathcal{I}_{\hat{\mathcal{T}}}$  or  $\mathcal{I}_{\check{\mathcal{T}}}$ . This is due to the fact that  $\check{h}$  and  $\hat{h}$  may aggregate states, resulting in some  $\mathcal{I}_{\mathcal{T}}$  representing only partial interventions on  $\mathcal{I}_{\check{\mathcal{T}}}$  or  $\mathcal{I}_{\hat{\mathcal{T}}}$ . In this context, the introduction of a notion termed allowed intervention becomes crucial. This notion is essential as certain interventions in the abstract model may lack definition or relevance in a well-defined concrete model. Consequently, within this framework of definitions, the translation of interventions is not universal; rather, only essential interventions that can be meaningfully mapped are considered.

We follow a similar but simpler procedure for under-approximations. Given an under-approximate causal transition system  $\check{\mathcal{T}}$ , we construct a sequence  $\check{\mathcal{T}}_0 \leq \check{\mathcal{T}}_1 \leq \dots \leq \check{\mathcal{T}}_k$  of under-approximations, where (1)  $\check{\mathcal{T}}_k = \check{\mathcal{T}}$ , and  $\check{\mathcal{T}}_{i+1}$  is a refinement of  $\check{\mathcal{T}}_i$ , for  $0 \leq i < k$ , which we compute using *counterexamples*. A counterexample is a state of  $\Sigma$  that is not in  $\check{\Sigma}_i$ . In this article, since we begin causal analysis from a trace log  $\Delta$ , we compute an under-approximation by a subset of the input set of traces. That is,  $\check{h}(\text{Tr}) \subseteq \text{Tr}$ .

### C. Detailed Description of Algorithm

The input to Algorithm 1 is a concrete transition systems  $\mathcal{T}$  (more specifically, its trace set) and a causal formula  $\varphi_e$ . Also,  $\alpha$  and  $\beta$  are parameters used in computing  $\check{h}$  and  $\hat{h}$ , respectively.  $\alpha$  indicates the subset size of  $\check{h}$  and  $\beta$  is a threshold to compute Euclidean distance of states for over-approximation. We are restricted to a set of allowed interventions  $\mathcal{I}_{\mathcal{T}}^h$ . Our objective is to identify states of  $\mathcal{T}$ , where causal formula  $\varphi_c$  holds as an actual cause in the trace  $\tau = \hat{\Lambda}(\sigma_0)\hat{\Lambda}(\sigma_1) \dots$ .

Line 1 initializes the under-approximation  $\check{\mathcal{T}}$ , with parameter  $\alpha$  indicating the number of traces to use and map in  $\check{h}$  function. In lines 3–16, the algorithm computes whether

the SMT query returns  $\varphi_c$  as the cause for effect  $\varphi_e$  in the current under-approximation and over-approximation. Specifically, in line 3, the SMT function receives  $\check{\mathcal{T}}$  as the under-approximation and constraints of *AC1* and *AC2(a)* specified in Fig. 5, and it returns the result  $\varphi_c$  as the cause. The SMT solver also returns a witness trace  $\tau \in \check{\text{Tr}}$ . In line 16, if the result of SMT query in line 3 is unsatisfiability, then the algorithm chooses more traces  $\text{Tr}$  by increasing  $\alpha$ . Indeed, lines 15 and 16 establish the refinement for the under-approximate model.

If a cause  $\varphi_c$  is identified by satisfying *AC1* and *AC2(a)*, we use this cause to initialize over-approximation in line 5 (to ensure Assumption 1), where we include all original states as well as potentially unreachable states by creating an abstract representation by function  $\hat{h}$ , such that all states in  $\mathcal{T}$  map to  $\hat{\mathcal{T}}$ , and also similar states are merged into a single abstracted state in  $\hat{\mathcal{T}}$ . The distance threshold for merging states is controlled by the parameter  $\beta$ . If the distance between any pair of states is less than  $\beta$ , those states will be merged. Consequently, a smaller  $\beta$  results in a larger number of abstract states, while a larger  $\beta$  leads to a smaller number of abstract states.

In lines 7–11, we focus on verifying *AC2(b)* using the over-approximation. In line 7, the SMT query takes  $\hat{\mathcal{T}}$  as the current over-approximate model and  $\varphi_c$  as output from line 3. It then examines whether all traces for which  $\varphi_c$  and  $\varphi_e$  hold can be modified by changing states such that  $\varphi_e$  still holds. If the SMT solver returns SAT, then  $\varphi_c$  is returned as the actual cause, where  $\varphi_c$  is a Boolean expression on the atomic propositions related to states in a specific trace. If the result is not SAT, in line 11, we use counterexample(s) in  $\Sigma - \hat{\Sigma}$ , allowed interventions identified by  $\hat{w}(\mathcal{I}_{\mathcal{T}})$ . These counterexamples are then eliminated by *Refine*, and the resulting model is assigned to the new  $\hat{\mathcal{T}}$ . We emphasize that in the refinement step for over-approximation (line 11), it is crucial to consider restricted interventions, denoted as  $\mathcal{I}_{\hat{\mathcal{T}}}^h$ . This consideration is necessary because, in a concrete model, certain interventions may not be directly mapped to their counterparts in the over-approximated model. Consequently, the refinement process must incorporate  $\mathcal{I}_{\mathcal{T}}^h$  as an essential input, utilizing it effectively during the mapping process to ensure consistency of model translation.

*Theorem 1:* Let  $\mathcal{T}$  be a concrete causal transition system and  $\varphi_c$  and  $\varphi_e$  be two causal formulas. If  $\varphi_c$  is an actual cause of  $\varphi_e$  identified by Algorithm 1 (for  $\hat{\mathcal{T}}$  and  $\check{\text{Tr}}$ ), then  $\varphi_c$  is an actual cause of  $\varphi_e$  in  $\mathcal{T}$ .

### D. Correctness

In this section, we formally prove the soundness of Algorithm 1.

*Theorem 2:* Let  $\mathcal{T}$  be a concrete causal transition system and  $\varphi_c$  and  $\varphi_e$  be two causal formulas. If  $\varphi_c$  is an actual cause of  $\varphi_e$  identified by Algorithm 1 (for  $\hat{\mathcal{T}}$  and  $\check{\text{Tr}}$ ), then  $\varphi_c$  is an actual cause of  $\varphi_e$  in  $\mathcal{T}$ .

## V. EXPERIMENTAL EVALUATION

This section first provides an overview of the implementation details of the algorithm proposed in Section IV-C. We also evaluate our technique on three case studies: 1) mountain

626 car; 2) Lunar Lander environments from OpenAI Gym [24]—  
 627 commonly used evaluation benchmarks for learning-enabled  
 628 CPS; and 3) an F-16 autopilot simulator [25] that uses an MPC  
 629 controller.

### 630 A. Implementation

631 To identify the actual cause of failures in our studies, we  
 632 need to generate traces consisting of those that do not violate  
 633 safety and those that do violate safety. We need successful  
 634 traces to find counterfactuals for failure scenarios, where the  
 635 same conditions lead to success through different decisions.

636 In our experiments, we use 47 networks for the mountain car  
 637 experiment [15], [29] and generate over 570 neural networks,  
 638 trained with Deep Reinforcement Learning [30], for the Lunar  
 639 Lander case study. Consequently, the success rates of traces  
 640 in satisfying  $\varphi_e$  used in our experiments were 17% and 11%  
 641 for the case studies in Sections V-C and V-D, respectively. In  
 642 the case study in Section V-E, the success rates were 21%  
 643 and 33% for the first and second scenarios, respectively. This  
 644 is because the non-AI MPC controller typically makes better  
 645 decisions than the AI controller.

646 We have implemented Algorithm 1 using the Python pro-  
 647 gramming language. Algorithm 1 is implemented through two  
 648 approaches. First, the Z3 SMT solver [21], and second (for  
 649 nonsymbolic cases), by employing a search method to find  
 650 traces in datasets that meet the HP conditions. For instance,  
 651 if we find a trace that leads to failure, we take this sample  
 652 and search for other traces with the same features, except  
 653 for the decision that caused the failure in the original trace.  
 654 To accomplish this, we utilize built-in data science search  
 655 algorithms in [22] and [23]. In fact, in our case studies, we are  
 656 dealing with large-sized data rather than symbolic properties.  
 657 Therefore, in Section V-F, we will demonstrate that searching  
 658 through the dataset is more efficient compared to Z3. While  
 659 Z3 is primarily employed for its robust capabilities in theorem  
 660 proving and constraint solving, it is not as effective for finding  
 661 traces in a large set of already generated traces that meet  
 662 certain conditions.

663 For execution of Algorithm 1, specific strategies are adopted  
 664 in refinement of the under- and especially over-approximate  
 665 (function Refine in line 11) models in cases of unsatisfiability.  
 666 In the under-approximation model, a parameter  $\alpha$  is utilized  
 667 to incorporate additional traces. This parameter can be pro-  
 668 gressively increased to obtain more traces, thereby refining  
 669 the under-approximation. In the over-approximation model, a  
 670 parameter  $\beta$  is used within the mapping function to dictate the  
 671 threshold for the distance between states. When the distance  
 672 between a group of states is less than this threshold, they are  
 673 merged into a single state to simplify the model. Moreover, in  
 674 refining the over-approximate model, the algorithm checks for  
 675 the existence of counterexample states that violate the over-  
 676 approximation. If such states are identified, they are removed  
 677 from the model to ensure its accuracy [28].

678 Assumption 1 for both our case studies is implemented as  
 679 follows that in the over-approximation function, it is crucial  
 680 not to merge states that transition to different outcomes. For  
 681 instance, in the mountain car example, if there are two traces

that differ only in their actions but have the same position and  
 velocity, and the under-approximation model identifies that  
 action might be a possible cause of failure, these states cannot  
 be merged in the over-approximation model. This is because  
 merging them would obscure the distinction between a trace  
 leading to failure and another leading to success.

### B. Experimental Settings

688 All of our experiments were conducted on a single core of the  
 689 Apple M2 Pro CPU, which features a 10-core architecture and  
 690 operates @3.7 GHz. Given a set of collected traces, we applied  
 691 our techniques in four different modes to identify the cause of  
 692 failure (safety violations): 1) *Only\_Z3* is the implementation,  
 693 where we only use the SMT solver Z3 to discover actual  
 694 causality (the technique proposed in Section III); 2) *Abs\_Z3* is  
 695 the implementation, where Algorithm 1 uses Z3; 3) *Only\_DA*  
 696 is the implementation, where we only use the search algorithms  
 697 in [22] and [23] in lieu of an SMT solver; and 4) *Abs\_DA* is the  
 698 implementation, where Algorithm 1 uses the search algorithms  
 699 in [22] and [23] in lieu of an SMT solver.  
 700

### C. Case Study 1: Mountain Car

701  
 702 Out first case study is the continuation of our running  
 703 example. In Fig. 3(a), the car is initially positioned in the  
 704 valley between two mountains with the objective being to  
 705 navigate it to the peak of the right mountain before a set dead-  
 706 line. The system incorporates three variables in accordance  
 707 with (1) and (2), specifying the domain for each variable as  
 708  $\text{pos}(t) \in [-1.2, 0.6]$ ,  $\text{vel}(t) \in [-0.07, 0.07]$ , and  $\text{action}(t) \in$   
 709  $[-1, 1]$ . Here, action represents a learning-based function  
 710  $f$ , implemented using various pretrained neural networks of  
 711 different dimensions

$$\text{action}(t) = f(\text{pos}(t), \text{vel}(t)).$$

712  
 713 The car's mission is to achieve  $\text{pos}(t) = 0.6$  before the time  
 714 limit of  $t = 100$  episodes. Our study explores various initial  
 715 settings for  $\text{pos}(0)$ ,  $\text{vel}(0)$ , and the function  $f$  to find the cause  
 716 of the vehicle's failure to reach its target. In our study, we  
 717 began by collecting data by assigning different initial values to  
 718 the variables  $\text{pos}(0)$  and  $\text{vel}(0)$ , which were treated as external  
 719 (exogenous) variables. We also utilized various combinations  
 720 of pretrained neural networks as the decision-making mech-  
 721 anism for acceleration. The action controller in the mountain  
 722 car scenario employs a neural network characterized by a  
 723 rectangular architecture with varied dimensions. The sigmoid  
 724 function serves as the activation mechanism for both the input  
 725 and hidden layers, whereas the Tanh function is utilized for  
 726 the output layer. This approach represents a modification of  
 727 the methodology detailed in [15] and [29].

728 By executing multiple initial valuations with distinct neural  
 729 networks, we generated a substantial set of traces, each  
 730 indicating whether the car reached its destination within 100  
 731 episodes.

### D. Case Study 2: Lunar Lander

732  
 733 In this case study, the space lander is initially positioned  
 734 at a certain altitude from the ground, aiming to land on the



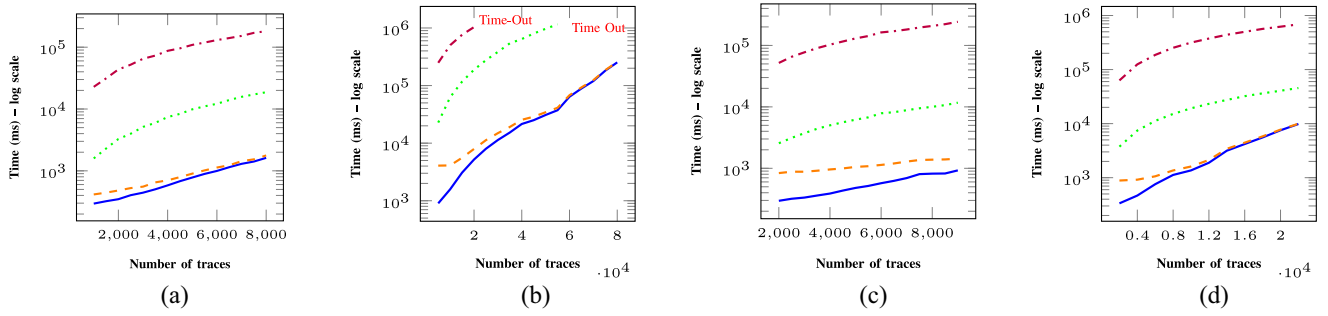


Fig. 6. Comparison of four modes of our implementation for various case studies. The legend is as follows: (a) — represents *Abs\_DA*, (b) - - - represents *Abs\_Z3*, (c) ····· represents *Only\_DA*, and (d) - · - · represents *Only\_Z3*.

735 designated landing pad. The landing pad is always located  
 736 at  $(0 \pm \epsilon, 0 \pm \epsilon)$ . In the Lunar Lander system, there are  
 737 eight variables (e.g.,  $x$  and  $y$  coordinates, velocity, angular  
 738 velocity, angle, etc.), which are intrinsic to the system, and an  
 739 additional seven variables that are configured to represent the  
 740 environment (e.g., wind, gravity, turbulence power, etc.). For a  
 741 comprehensive overview of this case study, refer to [24]. The  
 742 exogenous variables we consider are four different values for  
 743 wind:  $\{0, 5, 10, 15\}$ , three values for gravity:  $\{-8, -10, -12\}$ ,  
 744 and three values for turbulence power:  $\{0.8, 1.5, 2\}$ . Moreover,  
 745 in our experiment, we focus on a subset of the endogenous  
 746 variables, specifically  $\text{pos}_x(t)$ ,  $\text{pos}_y(t)$ ,  $\text{vel}_x(t)$ ,  $\text{vel}_y(t)$ , and  
 747  $\text{action}(t)$ . These variables correspond to the horizontal and  
 748 vertical positions, horizontal and vertical velocities, and the  
 749 action controlling the engines of the lander, respectively.

750 In our model, action denotes a learning-based function  
 751  $f$ , which is implemented using various pretrained neural  
 752 networks with different dimensions

$$753 \quad \text{action}(t) = f(\text{pos}_x(t), \text{pos}_y(t), \text{vel}_x(t), \text{vel}_y(t)).$$

754 In the Lunar Lander environment, there are four discrete  
 755 actions available for controlling the lander, denoted as  
 756  $\text{action} = \{0, 1, 2, 3\}$ .

- 757 1) 0: Do nothing.
- 758 2) 1: Fire the left orientation engine.
- 759 3) 2: Fire the main engine.
- 760 4) 3: Fire the right orientation engine.

761 The action controller designed for the Lunar Lander exper-  
 762 iment is based on deep reinforcement learning principles,  
 763 as explored in [30]. The neural networks employed in this  
 764 experiment are rectangular in shape and utilize the rectified  
 765 linear unit (ReLU) activation function to introduce nonlin-  
 766 earity and enhance the learning capability of the model. We  
 767 performed multiple simulations with varying initial values  
 768 for  $\text{pos}_x(0)$ ,  $\text{pos}_y(0)$ ,  $\text{vel}_x(0)$ ,  $\text{vel}_y(0)$ , wind, gravity, and  
 769 turbulence each paired with different neural networks. This  
 770 procedure produced a substantial set of traces, with each  
 771 trace indicating whether the lander successfully landed on the  
 772 landing pad within the time frame of  $t < 500$  episodes or not.

### 773 E. Case Study 3: F-16 Autopilot MPC Controller [25]

774 This benchmark models both the inner-loop and outer-loop  
 775 controllers of the F-16 fighter jet. We explore two scenarios.  
 776 The first scenario involves reaching a specified altitude set

point while maintaining a certain speed. The second scenario  
 tests whether the automated collision avoidance system can  
 recover the aircraft from a critical moment.

777  
 778  
 779  
 780  
 781  
 782  
 783  
 784  
 785  
 786  
 787  
 788  
 789  
 790  
 791  
 792  
 793  
 794  
 795  
 796  
 797  
 798  
 799  
 800  
 801  
 802  
 803  
 804  
 805  
 806  
 807  
 808  
 809  
 810  
 811  
 812  
 813  
 814  
 815  
 816

1) *First Scenario*: In this scenario, the aircraft's goal is to reach a certain altitude while maintaining a specified speed within a timeline of  $t$ . There are 16 state variables (e.g., altitude, airspeed, pitch, yaw, roll, power-lag, angle of attack (AoA) noted as  $\alpha$ , etc.). Our exogenous variables are the initial settings for altitude(0),  $\alpha(0)$ , airspeed(0), pitch(0), and the power lag that the engine suffers (power-lag). Our endogenous variables are altitude( $t$ ),  $\alpha(t)$ , airspeed( $t$ ), pitch( $t$ ), power-lag( $t$ ), and the actions of the autopilot system for  $t > 0$ , which include changing the throttle  $\delta_t(t)$  and adjusting the angle of the elevators  $\delta_e(t)$  to control the pitch (nose up or down). In this experiment, we investigate the actions ( $\delta_t(t)$  and/or  $\delta_e(t)$ ) that determine whether the plane succeeds or fails in reaching the desired checkpoint, achieving the desired speed, or violating aircraft limits, such as upward acceleration, AoA, or minimum airspeed, that could lead to stalling.

2) *Second Scenario*: Here, we place the aircraft in a critical position near the ground to evaluate its collision avoidance system. This scenario involves using a larger set of variables, thereby increasing the dimensionality of our problem compared to the previous scenario. These critical moments involve high degrees of pitch, roll, and yaw, as well as low airspeed near the ground, which may lead to failures, such as ground collision and violations of the aircraft's aerodynamic limits. Our exogenous variables are the initial settings for altitude(0), airspeed(0), pitch(0),  $\alpha(0)$ , yaw(0), roll(0), and power-lag, while the endogenous variables are altitude( $t$ ), airspeed( $t$ ), pitch( $t$ ), yaw( $t$ ), roll( $t$ ), for  $t > 0$  and the actions of the autopilot system. These actions include adjusting the degree of the rudder  $\delta_r(t)$  to change the yaw of the plane, changing the degree of the aileron  $\delta_a(t)$  to modify the roll of the plane, and controlling the throttle  $\delta_t(t)$  and elevator  $\delta_e(t)$ . As in the previous scenario, we are examining the autopilot decisions that influence whether the aircraft can successfully recover from a potential collision or avoid violating aerodynamic constraints. Additionally, we aim to identify the actual cause of the failures.

### F. Performance Analysis

Fig. 6(f)–(h) illustrate the results of our experiments for the mountain car, Lunar Lander, and both F-16 simulation

TABLE I  
EXPERIMENT ON 1000 TRACES

Case Study	Algorithm	$\alpha$	Refinement steps	Time (ms)
Mountain Car	Abs_DA	0.01	28	1024
		0.05	7	475
		0.1	2	102
	Abs_Z3	0.01	22	9793
		0.05	6	4757
		0.1	2	1306
Lunar Lander	Abs_DA	0.01	24	494
		0.05	7	396
		0.1	3	150
	Abs_Z3	0.01	19	2809
		0.05	4	794
		0.1	3	239

820 scenario, respectively. Indeed all graphs show a similar profile  
821 in terms of the behavior of the four modes of experiments  
822 mentioned in Section V-A.

823 As shown in the graphs, the abstraction algorithms (*Abs\_DA*  
824 and *Abs\_Z3*) demonstrate significantly better performance by  
825 orders of magnitude than the conventional solvers (*Only\_DA*  
826 and *Only\_Z3*), with the latter exhibiting exponential growth  
827 in runtime with an increasing number of traces. This  
828 demonstrates the effectiveness of our abstraction-refinement  
829 technique: it identifies the actual causes of failures while  
830 running much faster than techniques on concrete traces. As  
831 shown in Fig. 6(f), our technique processes up to 80 000 traces  
832 in under 250 s, whereas *Only\_Z3* times out with threshold  
833 1200 s at 20 000 traces and *Only\_DA* at 55 000 traces.

834 Notably, *Abs\_DA* outperforms *Abs\_Z3*, and *Only\_DA* shows  
835 better performance than *Only\_Z3*. This observation can be  
836 attributed to the fundamental differences between SMT  
837 solvers, which focus on logical consistency, and the search-  
838 ing methods developed in data analysis libraries, which are  
839 tailored for efficient searching in large datasets.

840 In Table I, we present a comparison between different  
841 valuations of the parameter  $\alpha$ , which represents the subset  
842 size of  $\tilde{h}$ . We conducted an experiment to find an optimal  
843 value for  $\alpha$ . Our findings indicate that a very small  $\alpha$  may  
844 require numerous refinements, as it needs to add more traces  
845 to identify the cause, which is inefficient. On the other hand,  
846 large values of  $\alpha$  needs fewer refinements, but the under-  
847 approximation function has to process a larger amount of  
848 data, which increases the processing time. Therefore, there is  
849 a tradeoff between the number of refinements and the total  
850 time spent on them. We note that for row that have equal  $\alpha$ ,  
851 we shuffle the trace set, which impact computing the under-  
852 approximation.

### 853 G. Causality Analysis

854 This section demonstrates an important aspect of this  
855 research in investigating the actual cause of safety failures in  
856 CPS to *explain* the underlying reason. Our case studies involve  
857 simulations that specifically focus on the intersection of AI-  
858 enabled decision-making (mountain car and Lunar Lander),  
859 environmental dynamics feedback, and the correctness of a  
860 non-AI controller within an F-16 aircraft simulation.

861 1) *Mountain Car*: In Example 3 (see Fig. 4), we prove  
862 that making a poor decision to accelerate to the right (i.e.,  
863  $\text{action}(0) = 1$ ) leads to failing in reaching the mountain top

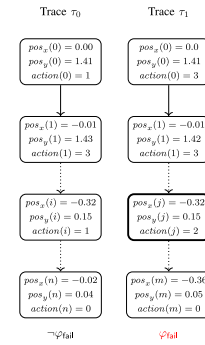


Fig. 7. Simulated traces in Lunar Lander and causal effect of decision by the main engine.

(i.e., formula  $\varphi_{\text{fail}}$ ). Instead, in the counterfactual scenario  
864 we observe that it is necessary to accelerate to the left to  
865 gain momentum in order to climb the mountain. This not  
866 only shows the earliest bad decision by the controller but  
867 also identifies the “but-for” scenario, meaning what would  
868 have happened if a different action was taken. Additionally,  
869 counterfactual reasoning demonstrates how to fix the bad  
870 decision made by the neural network.  
871

872 2) *Lunar Lander*: We observe that when there is a strong  
873 wind from left to right, some controllers tend to overuse the  
874 right engine, resulting in  $\text{action} = 3$  during the initial steps.  
875 This causes the lander to drift to the left. However, we observe  
876 that even in this situation where the lander is positioned to the  
877 left of the landing pad, the controller can use its left engine,  
878  $\text{action} = 1$ , to move the lander to the right and land safely.  
879 However, some controllers use their main engine,  $\text{action} = 2$ ,  
880 resulting in the lander not reaching the landing pad. This  
881 results in  $\text{pos}_x(t) < 0 - \epsilon$ , constituting a failure.

882 To illustrate this further, Fig. 7 shows two traces starting  
883 from the same point but taking different actions in the  
884 first step. Dotted transitions means the next state is not the  
885 immediate next time step. The final state of the traces  $\tau_0$  and  $\tau_1$   
886 is the  $n$ th and  $m$ th state, respectively. In trace  $\tau_0$ ,  $\text{action}(0) = 1$ ,  
887 while in trace  $\tau_1$ ,  $\text{action}(0) = 3$ . However, in both traces, the  
888 controllers overuse the right engine in the initial steps (both  
889 controllers in  $\tau_0$  and  $\tau_1$  use the right engine  $\text{action}(1) = 3$ ),  
890 causing the lander to drift far to the left, resulting in  $\text{pos}_x(i) =$   
891  $-0.32$  in  $\tau_0$  and  $\text{pos}_x(j) = -0.32$  in  $\tau_1$ . At this state, where in  
892 both scenarios the lander has the same position and setting, the  
893 controller in  $\tau_1$  decides to use the main engine  $\text{action}(j) = 2$ ,  
894 while the controller in  $\tau_0$  opts to use the left engine  $\text{action}(i) =$   
895  $1$  to move the lander to the right. These decisions under similar  
896 conditions lead to the failure of  $\tau_1$  (i.e.,  $\text{pos}_x(m) = -0.36 <$   
897  $0 - \epsilon$ ) and the success of  $\tau_0$  (i.e.,  $0 - \epsilon < \text{pos}_x(n) = -0.02 <$   
898  $0 + \epsilon$ ). This finding indicates that the failure in  $\tau_1$  using  
899 decision  $\text{action}(j) = 2$ , while the counterfactual scenario in  $\tau_0$   
900 succeeds with a different decision  $\text{action}(i) = 1$ , highlighting  
901 that  $\text{action}(j) = 2$  in  $\tau_1$  is the actual cause of the failure.

902 3) *F-16 Autopilot Simulation*: Here, we identify the cause  
903 of failures and analyze counterfactual scenarios (alternative  
904 actions) under the same conditions that could lead to success.  
905 When the aircraft needs to gain altitude at *low speed*, some  
906 traces show the controller lowering the nose to gain speed and

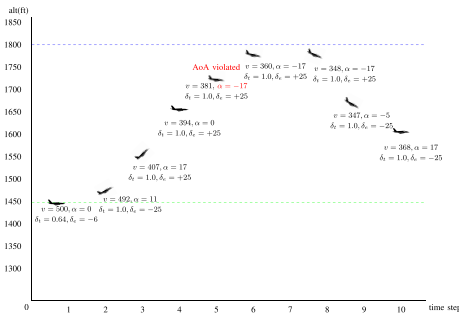


Fig. 8. F-16 scenario leads to failure due to a violation of the AoA limit.

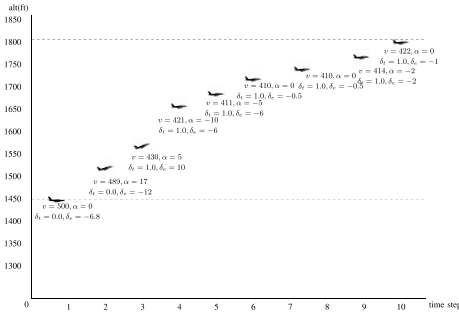


Fig. 9. F-16 counterfactual scenario leads to success.

the aircraft’s maximum negative AoA limit is  $-15$ ,  $\alpha(5) = -17$  violates this limit, and the controller fails to achieve its objective. On the contrary, in the counterfactual scenario (see Fig. 9), the controller starting with less aggressive throttle and elevator adjustments, such as  $\delta_t(1) = 0$  and  $\delta_e(1) = -6.8$ , resulting in a slight pitch. This strategy continues similarly with  $\delta_t(2) = 0$  and  $\delta_e(2) = -12$ , avoiding harsh climbs to reach the destination. By examining this scenario, we find that in the first time step, Fig. 8 makes the decisions  $\delta_t(1) = 0.64$  and  $\delta_e(1) = -6.8$ , while Fig. 9 makes  $\delta_t(1) = 0.0$  and  $\delta_e(1) = -6.8$  under the same conditions (same altitude, speed, etc.). This counterfactual example shows that an alternative decision by the controller leads to success, providing sufficient evidence that the initial decision is the actual cause of failure.

VI. RELATED WORK

There is a wealth of research on causality analysis in the context of embedded and component-based systems from different perspectives. In [2], [3], [4], [6], [7], [8], and [9], a new structure of formal causal analysis is proposed that can serve as a substitute for the HP causal model. This approach is distinct from our work, which utilizes a framework of causal analysis to identify the cause of a specific effect. Recently, there has been great interest in using temporal logics to reason about causality and explaining bugs [10], [11], [12], [13]. However, these lines of work either focus on only modeling aspects of causality or do not address the problem of scalability in automated reasoning about causality, which inherently involves a combinatorial blow up for counterfactual reasoning. In the CPS domain, using causality to repair AI-enabled controllers has recently gained interest [15]. This work explored the construction of HP models on AI-enabled controllers, the search for the cause of failure using a search algorithm, and the verification of these causes using HP constraints. In contrast, our work focuses on identifying the cause of failure efficiently in traces using HP constraints and proposes an efficient method for doing so. In [31], causal analysis is performed on system models and system execution traces. In contrast, our algorithm is designed to efficiently identify the cause of any potential failure. Additionally, our work is focused on systems, such as CPS, that interact with their environment.

Although the idea of abstracting causal models in terms of structural equations has been studied in [18], [19], and [20], these works do not attempt to establish a relation between *actual* causes in the abstract and concrete causal models. In the studies [18], [19], [20], the concept of abstraction in causal models was introduced, along with the preliminaries required to construct an abstraction function that maps low-level variables to high-level variables. The work in [18] presents a more general form of abstraction, while [19], [20] focus on the concept of intervention in causal models and how to build an abstraction that preserves them. The distinction between our work and these studies lies in our objective; we are not aiming to construct causal models, but rather, we are

avoid stalling before attempting to climb. This approach results in a loss of altitude and insufficient time to reach the desired altitude within the specified time frame, leading to failure. However, in counterfactual scenarios, the controller opts to gain speed by using more throttle and then gradually raises the nose using the elevators, eventually reaching the desired altitude. This demonstrates that the decision to lower the nose is the actual cause of the failure to reach the desired altitude within the specified time frame.

In another scenario, when transitioning from a lower to a higher altitude, some traces show controllers using excessive elevator and throttle, which places the aircraft in a danger zone and violates the AoA limits, leading to catastrophic failure. However, in alternative counterfactual scenarios with the same starting conditions, the controller gradually uses the throttle and adjusts the elevator more cautiously. This approach allows the aircraft to reach the desired altitude without violating its aerodynamic limits.

To illustrate the latter scenario in detail, Figs. 8 and 9 show two flight real paths starting from the same altitude,  $altitude(1) = 1450$ , and the same speed,  $airspeed(1) = 500$ , with the goal of reaching an altitude of  $altitude(n) = 1800$ . This process should occur within a specified time frame while not violating aircraft limits. In Fig. 8, the controller starts by using throttle  $\delta_t(1) = 0.64$  and setting the elevator to a negative position,  $\delta_e(1) = -6$ , to achieve a positive pitch angle. This decision continues in subsequent steps in a more extreme manner, with  $\delta_t(2) = 1.0$  (full throttle) and  $\delta_e(2) = -25$ , resulting in nearly a 45-degree pitch. Next, to counteract this situation, the controller attempts to use  $\delta_e(3) = 25$  and  $\delta_e(4) = 25$  to stabilize the aircraft’s sharp nose-up attitude, leading to a negative AoA,  $\alpha(5) = -17$ . Since

utilizing abstraction to identify the cause of an effect in a more efficient manner.

In [10] and [32], the concept of explaining counterexamples returned from the model checker is proposed, with one focusing on specifications in LTL format and the other in HyperLTL format. However, in our work, we aim to efficiently identify the cause of failure in an embedded system.

## VII. CONCLUSION

We concentrated on designing an efficient technique to reason about actual causality. We proposed an SMT-based formulation to determine whether for an input transition system or a set of traces and a state formula (the effect), there exists an actual cause. Since identifying an actual cause involves counterfactual reasoning and, hence, a combinatorial blow up, we also introduced an efficient heuristic based on abstraction refinement. We evaluated our techniques on three case studies from the CPS domain: AI-enabled controllers for a 1) mountain car; 2) Lunar Lander [24]; and 3) an MPC controller for an F-16 autopilot simulator [25].

One natural extension is to consider *probabilistic* actual causality, where either occurrence of events in the system are associated with probabilities, or, data points follow some distribution. Another important direction is causal models where the system is partially observable.

## REFERENCES

[1] Y. Zhou and E. A. Lee, "Causality interfaces for actor networks," *ACM Trans. Embedd. Comput. Syst.*, vol. 7, no. 3, pp. 1–35, 2008.

[2] M. Broy, "Time, causality, and realizability: Engineering interactive, distributed software systems," *J. Syst. Softw.*, vol. 210, Apr. 2024, Art. no. 111940.

[3] G. Goessler and L. Astefanoaei, "Blaming in component-based real-time systems," in *Proc. Int. Conf. Embedd. Softw.(EMSOFT)*, 2014, pp. 1–10.

[4] G. Gössler and J.-B. Stefani, "Causality analysis and fault ascription in component-based systems," *Theor. Comput. Sci.*, vol. 837, pp. 158–180, Oct. 2020.

[5] S. Cherrared, S. Imadali, E. Fabre, and G. Göbller, "SAKURA a model based root cause analysis framework for vIMS," in *Proc. 17th Annu. Int. Conf. Mobile Syst., Appl., Services(MobiSys)*, 2019, pp. 594–595.

[6] G. Göbller, O. Sokolsky, and J. Stefani, "Counterfactual causality from first principles?" in *Proc. 2nd Int. Workshop Causal Reason. Embedd. Safety-Critical Syst. Technol. (CREST)*, 2017, pp. 47–53.

[7] S. Wang, Y. Geoffroy, G. Göbller, O. Sokolsky, and I. Lee, "A hybrid approach to causality analysis," in *Proc. 6th Int. Conf. Runtime Verif. (RV)*, 2015, pp. 250–265.

[8] G. Göbller and D. L. Métayer, "A general trace-based framework of logical causality," in *Proc. 10th Int. Symp. Formal Aspects Compon. Softw. (FACS)*, 2013, pp. 157–173.

[9] G. Göbller, D. L. Métayer, and J. Raclet, "Causality analysis in contract violation," in *Proc. 1st Int. Conf. Runtime Verif. (RV)*, 2010, pp. 270–284.

[10] N. Coenen et al., "Explaining hyperproperty violations," in *Proc. 34th Int. Conf. Comput. Aided Verif. (CAV)*, 2022, pp. 407–429.

[11] B. Finkbeiner and A. Kupriyanov, "Causality-based model checking," in *Proc. 2nd Int. Workshop Causal Reason. Embedd. Safety-Critical Syst. Technol. (CREST)*, 2017, pp. 31–38.

[12] N. Coenen, B. Finkbeiner, H. Frenkel, C. Hahn, N. Metzger, and J. Siber, "Temporal causality in reactive systems," in *Proc. 20th Int. Symp. Autom. Technol. Verif. Anal. (ATVA)*, 2022, pp. 208–224.

[13] B. R. B. Finkbeiner, H. Frenkel, and J. Siber, "Checking and sketching causes on temporal sequences," in *Proc. 21st Int. Symp. Autom. Technol. Verif. Anal. (ATVA)*, 2023, pp. 314–327.

[14] P. Lu, M. Cleaveland, O. Sokolsky, I. Lee, and I. Ruchkin, "Repairing learning-enabled controllers while preserving what works," in *Proc. 15th Int. Conf. Cyber-Phys. Syst. (ICCPSS)*, 2024, pp. 1–11.

[15] P. Lu, I. Ruchkin, M. Cleaveland, O. Sokolsky, and I. Lee, "Causal repair of learning-enabled cyber-physical systems," in *Proc. IEEE Int. Conf. Assured Auton. (ICAA)*, 2023, pp. 1–10.

[16] J. Y. Halpern, *Actual Causality*. Cambridge, MA, USA: MIT Press, 2016.

[17] G. Aleksandrowicz, H. Chockler, J. Y. Halpern, and A. Ivrii, "The computational complexity of structure-based causality," *J. Artif. Intell. Res.*, vol. 58, pp. 431–451, Feb. 2017.

[18] P. K. Rubenstein et al., "Causal consistency of structural equation models," in *Proc. 33rd Conf. Uncertain. Artif. Intell. (UAI)*, 2017, pp. 1–15.

[19] S. Beckers, F. Eberhardt, and J. Y. Halpern, "Approximate causal abstractions," in *Proc. 35th Conf. Uncertain. Artif. Intell. (UAI)*, 2019, pp. 606–615.

[20] S. Beckers and J. Y. Halpern, "Abstracting causal models," in *Proc. 33rd AAAI Conf. Artif. Intell. (AAAI)*, 2019, pp. 2678–2685.

[21] L. M. de Moura and N. Bjørner, "Z3: An efficient SMT solver," in *Proc. 14th Int. Conf. Tools Algorithms Constr. Anal. Syst. (TACAS)*, 2008, pp. 337–340.

[22] C. R. Harris et al., "Array programming with NumPy," *Nature*, vol. 585, no. 7825, pp. 357–362, Sep. 2020.

[23] W. McKinney, "Data structures for statistical computing in python," in *Proc. 9th Python Sci. Conf.*, 2010, pp. 56–61.

[24] G. Brockman et al., "OpenAI Gym," 2016, *arXiv:1606.01540*.

[25] P. Heidlauf, A. Collins, M. Bolender, and S. Bak, "Verification challenges in F-16 ground collision avoidance and other automated maneuvers," in *Proc. 5th Int. Workshop Appl. Verif. Continuous Hybrid Syst. (ARCH)*, 2018, pp. 208–217.

[26] A. Rafieioskouei and B. Bonakdarpour, "Efficient discovery of actual causality using abstraction-refinement," 2024, *arXiv:2407.16629*.

[27] E. M. Clarke, A. Biere, R. Raimi, and Y. Zhu, "Bounded model checking using satisfiability solving," *Formal Methods Syst. Design*, vol. 19, no. 1, pp. 7–34, 2001.

[28] E. M. Clarke, O. Grumberg, S. Jha, Y. Lu, and H. Veith, "Counterexample-guided abstraction refinement," in *Proc. 12th Int. Conf. Comput. Aided Verif. (CAV)*, 2000, pp. 154–169.

[29] R. Ivanov, J. Weimer, R. Alur, G. J. Pappas, and I. Lee, "Verisig: Verifying safety properties of hybrid systems with neural network controllers," 2018, *arXiv:1811.01828*.

[30] V. Mnih et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.

[31] F. Leitner-Fischer and S. Leue, "Towards causality checking for complex system models," in *Proc. 8th Dagstuhl-Workshop Modellbasierte Entwicklung eingebetteter Systeme (MBEES) Model-Based Develop. Embedd.Syst.*, 2012, pp. 71–80.

[32] I. Beer, S. Ben-David, H. Chockler, A. Orni, and R. J. Treffer, "Explaining counterexamples using causality," in *Proc. 21st Int. Conf. Comput. Aided Verif. (CAV)*, 2009, pp. 94–108.