

# Temporal Causality in Reactive Systems

Norine Coenen<sup>1</sup>, Bernd Finkbeiner<sup>1</sup>, Hadar Frenkel<sup>1</sup>,  
Christopher Hahn<sup>2</sup>, Niklas Metzger<sup>1</sup>, and Julian Siber<sup>1</sup>

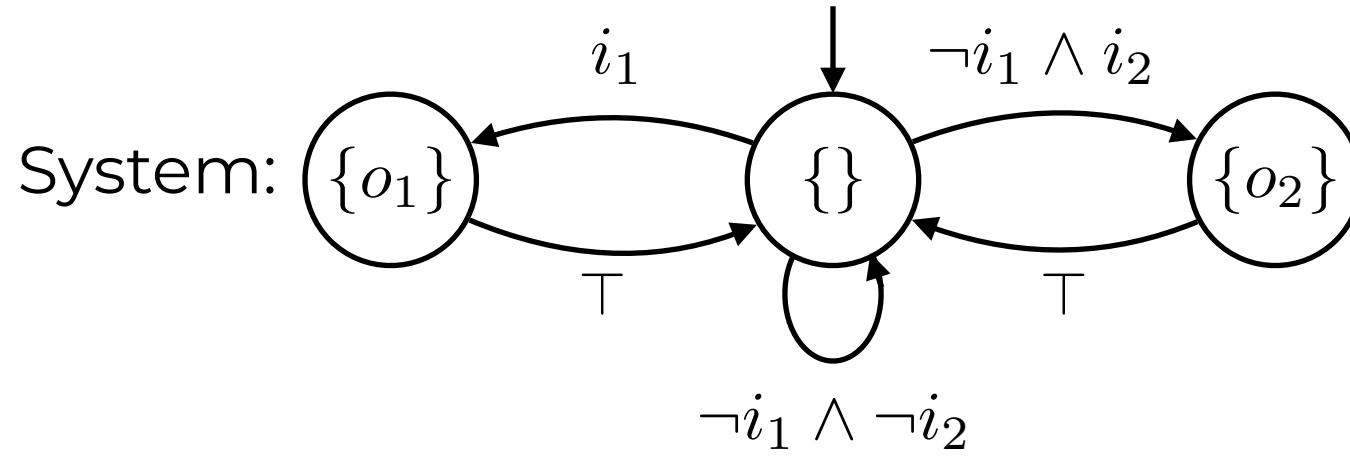
<sup>1</sup>*CISPA Helmholtz Center for Information Security*

<sup>2</sup>*Stanford University*

ATVA 2023 | 27 October 2023 | In Proceedings of ATVA 2022



# Temporal Properties as Causes



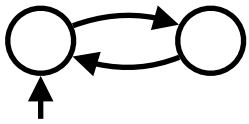
Trace:  $(\{i_1, i_2\} \{i_1, i_2, o_1\})^\omega$

Does  $\diamond i_1$  cause  $\square \neg o_2$ ?



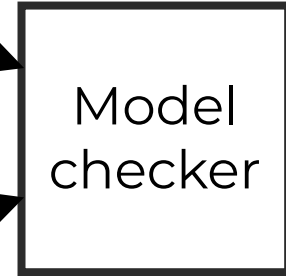
# Causes as Explanations for Model Checking

System model



$\varphi$

Specification



Counterexample  $\{a, b\}\{a\}(\{a, b\})^\omega$

The *causes* for  $\neg\varphi$  can explain the counterexample.

One solution: Highlighting

$\{a, \mathbf{b}\}\{a\}(\{a, \mathbf{b}\})^\omega$

$\diamond b ?$

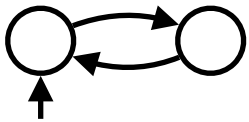
$\square \diamond b ?$

$b \wedge \bigcirc \bigcirc b ?$



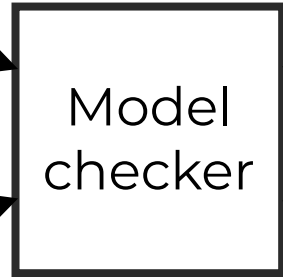
# Causes as Explanations for Model Checking

System model



$\varphi$

Specification



Counterexample  $\{a, b\}\{a\}(\{a, b\})^\omega$

The *causes* for  $\neg\varphi$  can explain the counterexample.

One solution: Highlighting

$\{a, \mathbf{b}\}\{a\}(\{a, \mathbf{b}\})^\omega$

**Our Solution:** Property causes

$\diamond b$

$\square \diamond b$

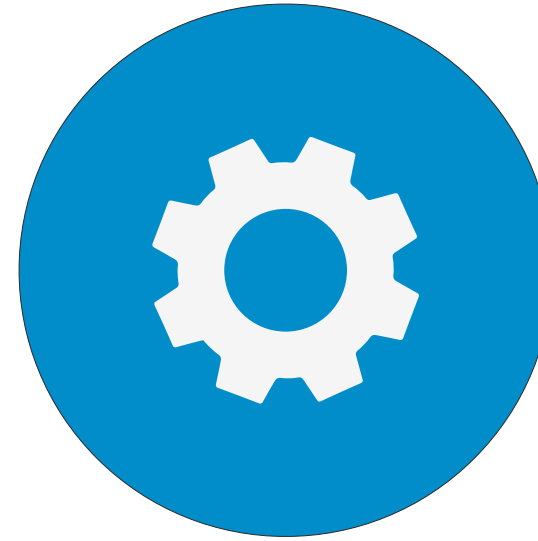
$b \wedge \bigcirc \bigcirc b$



# Outline



Define causality for  
trace properties



Algorithm for  
checking causality





# Actual Causality<sup>1,2</sup>



For finite sets of events  $(a, n) \in AP \times \mathbb{N}$  (proposition and time-point).

**SAT:** Cause and Effect have happened.

**CF:** If the cause had not happened (but everything else stayed the same), the effect would not have happened either.

**MIN:** There is no subset that satisfies the above.

---

<sup>1</sup>*Causes and Explanations: A Structural-Model Approach*. Halpern and Pearl (2005).

<sup>2</sup>*A Modification of the Halpern-Pearl Definition of Causality*. Halpern (2015).



# $C$ is a Cause for $E$ iff...



**SAT:**  $\pi$  satisfies  $C$  and  $E$  .

**CF:** If the cause had not happened (but everything else stayed the same), the effect would not have happened either.

**MIN:** There is no smaller cause candidate that satisfies the above.



# Distance Metrics



An adaption of similarity relations<sup>3,4</sup>.

$$\pi = \left( \{i_1, i_2\} \{i_1, i_2, o_1\} \right)^\omega$$

A distance metric  $<_{\pi}^C$  orders traces w.r.t. their similarity to  $\pi$ .

$$\pi_1 <_{\pi}^C \pi_2 \quad \text{iff} \quad zip(\pi, \pi_1, \pi_2) \models$$
$$\square \bigwedge_{i \in I} \left( (i_{\pi} \not\leftrightarrow i_{\pi_1}) \rightarrow (i_{\pi} \not\leftrightarrow i_{\pi_2}) \right) \wedge \diamond \bigvee_{i \in I} (i_{\pi_1} \not\leftrightarrow i_{\pi_2})$$

$\Rightarrow$  Causality is a *hyperproperty*.

---

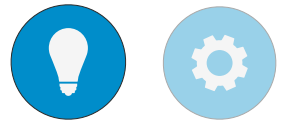
<sup>3</sup>A *Theory of Counterfactuals*. Stalnaker (1968).

<sup>4</sup>*Counterfactuals*. Lewis (1973).





# Counterfactual Input Sequences



$$\pi = (\{i_1, i_2\} \{i_1, i_2, o_1\})^\omega$$

The counterfactual input sequences are the closest sequences that negate  $C$ :

$$C_1 = i_1 \vee \bigcirc i_1 \quad \blacktriangleright \quad \sigma_{\neg C_1} = \{i_2\} \{i_2\} \{i_1, i_2\}^\omega$$

and not, e.g.:  $\{i_2\} \{i_2\} \{i_2\}^\omega$

$$C_2 = i_1 \wedge \bigcirc i_1 \quad \blacktriangleright \quad \sigma_{\neg C_2}^1 = \{i_2\} \{i_1, i_2\} \{i_1, i_2\}^\omega$$
$$\sigma_{\neg C_2}^2 = \{i_1, i_2\} \{i_2\} \{i_1, i_2\}^\omega$$



# Limit Assumption



$$\pi = (\{i_1, i_2\} \{i_1, i_2, o_1\})^\omega$$

The naive distance metric could be vacuously satisfied:

$$C_3 = \square \diamond i_1 \quad \Rightarrow \quad \begin{aligned} \sigma_{\neg C_3}^1 &= \{i_2\}^\omega && \text{is closer than} \\ \sigma_{\neg C_3}^2 &= \{i_1, i_2\} \{i_2\}^\omega && \text{is closer than} \\ \sigma_{\neg C_3}^k &= \{i_1, i_2\}^k \{i_2\}^\omega \\ &\dots \end{aligned}$$

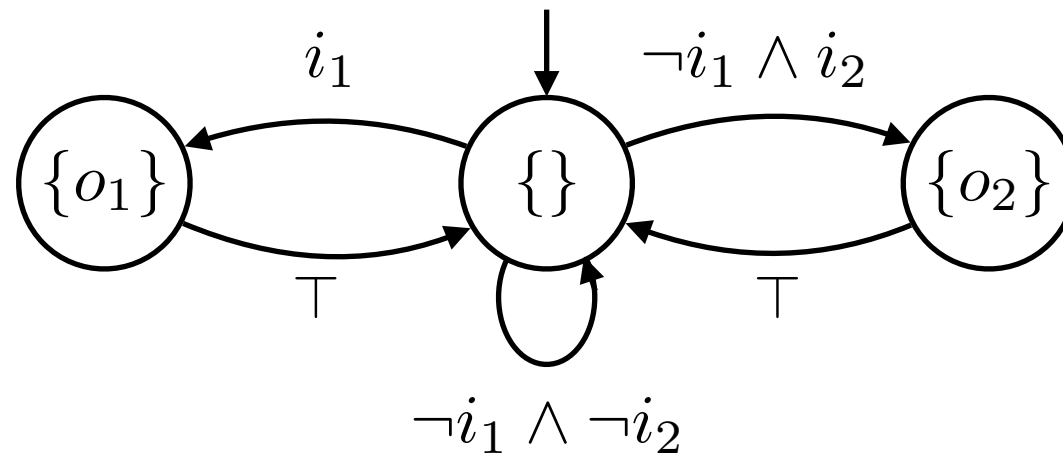
We propose an extension that satisfies the *limit assumption*.

(See our recent work<sup>5</sup> on how to accommodate more general metrics).

<sup>5</sup>Counterfactuals Modulo Temporal Logics. Finkbeiner and Siber (2023).



# Contingencies on Traces



Counterfactuals alone are often imprecise.

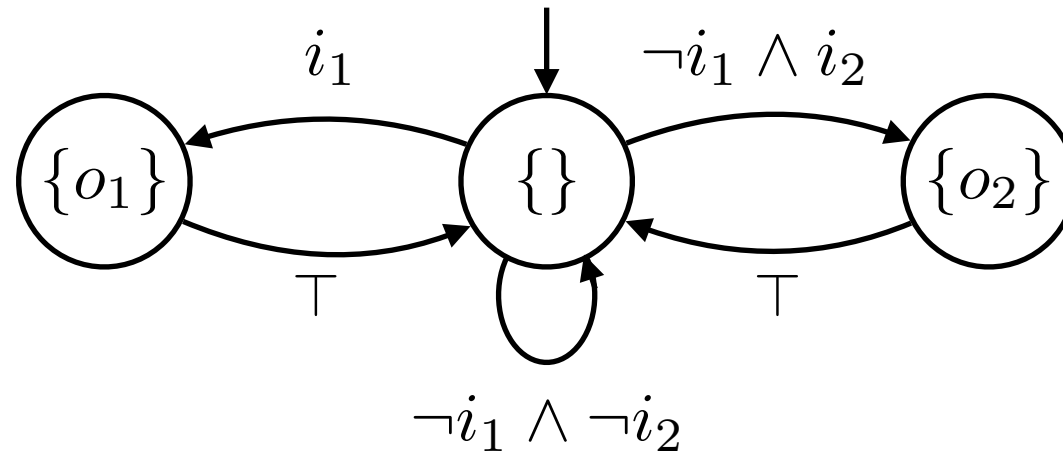
Consider:  $E = \text{O}(o_1 \vee o_2)$  and  $\pi = \{i_1, i_2\}\{o_1\}\{\}$  <sup>$\omega$</sup> .

$$\pi_{\neg i_1} = \{\cancel{i_1}, i_2\}\{o_2\}\{\}$$
 <sup>$\omega$</sup>

$\implies C = i_1$  alone does not negate the effect.



# Contingencies on Traces



Counterfactuals alone are often imprecise.

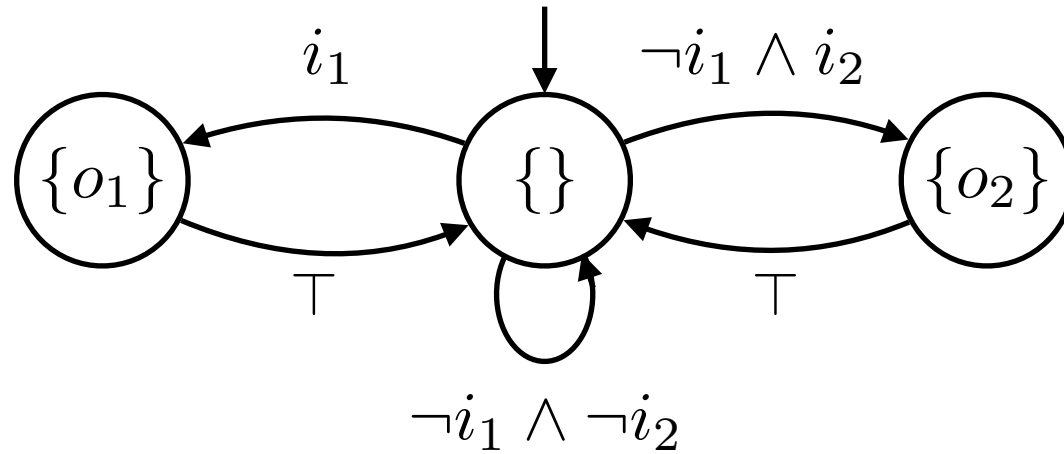
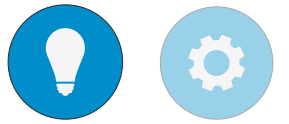
Consider:  $E = \mathcal{O}(o_1 \vee o_2)$  and  $\pi = \{i_1, i_2\}\{o_1\}\{\}$  <sup>$\omega$</sup> .

$$\pi_{\neg(i_1 \vee i_2)} = \{\cancel{i_1}, \cancel{i_2}\}\{\}\{\}$$
 <sup>$\omega$</sup>

$\implies C = i_1 \vee i_2$  works but is too imprecise.



# Contingencies on Traces



Counterfactuals alone are often imprecise.

Consider:  $E = \bigcirc(o_1 \vee o_2)$  and  $\pi = \{i_1, i_2\}\{o_1\}\{\}\omega$ .

$$\pi_{\neg i_1}^{cont.} = \{\cancel{i_1}, i_2\}\{\cancel{o_2}\}\{\}\omega$$

Contingency resets value.

$\implies C = i_1$  with the contingency  $\bigcirc \neg o_2$  works.



# $C$ is a Cause for $E$ iff...



**SAT:**  $\pi$  satisfies  $C$  and  $E$  .

**CF:** For every counterfactual input sequence  $\sigma$  , there exists a contingency trace  $\pi'$  such that  $\sigma =_{inputs} \pi'$  and  $\pi'$  does not satisfy  $E$  .

**MIN:** There is no smaller cause candidate that satisfies the above.



# Minimality



**SAT** and **CF** define a lot of potential causes.

$\pi_{\neg(i_1 \vee i_2)} = \{\cancel{i_1}, \cancel{i_2}\}\{\}\{\}\omega \implies C = i_1 \vee i_2$  works but is too imprecise.

$\pi_{\neg i_1}^{cont.} = \{\cancel{i_1}, i_2\}\{\cancel{o_2}\}\{\}\omega \implies C = i_1$  with the *contingency*  $\circ \neg o_2$  works.

Solution: prefer semantically minimal properties as causes, i.e., check:

$$i_1 \rightarrow (i_1 \vee i_2)$$



# $C$ is a Cause for $E$ iff...



**SAT:**  $\pi$  satisfies  $C$  and  $E$  .

**CF:** For every counterfactual input sequence  $\sigma$  , there exists a contingency trace  $\pi'$  such that  $\sigma =_{inputs} \pi'$  and  $\pi'$  does not satisfy  $E$  .

**MIN:** There does not exist a  $C'$  such that  $C \rightarrow C'$  and  $C'$  satisfies **SAT** and **CF**.

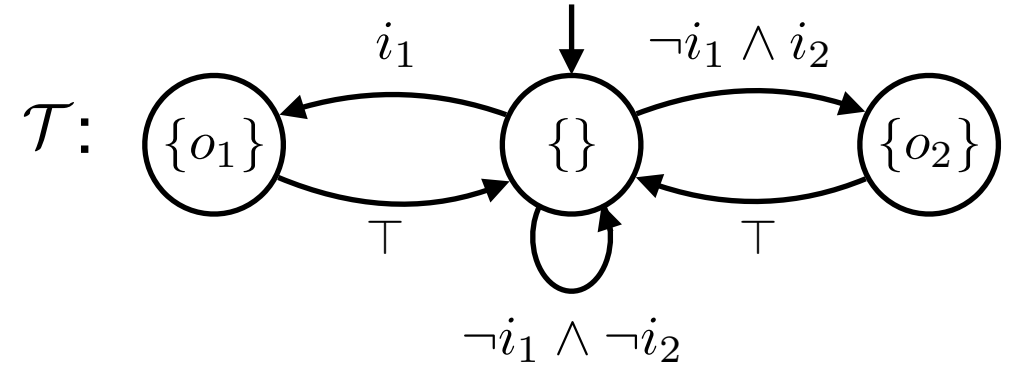




# Temporal Properties as Causes



$$\pi = (\{i_1, i_2\} \{i_1, i_2, o_1\})^\omega$$



Does  $\diamond i_1$  cause  $\square \neg o_2$  on  $\pi$  in  $\mathcal{T}$ ?

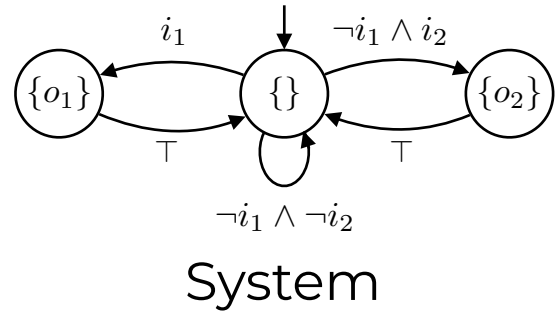
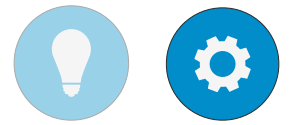
**SAT:** ✓

**CF:** ✓  $\pi' = (\{i_2\} \{i_2, o_2\})^\omega$  is the only counterfactual trace.

**MIN:** ✗  $\exists q. q \wedge \square(\bigcirc q \leftrightarrow \neg q) \wedge \square(q \rightarrow i_1)$  is more minimal (and the cause).



# Checking Temporal Causality



Counterfactual automaton (for contingencies)

Counterexample  $\pi = (\{i_1, i_2\} \{i_1, i_2, o_1\})^\omega$

$$\forall \pi. \exists \pi'. \varphi(C, E, \pi)$$

Cause  $C$

HyperQPTL encoding

Effect  $E$

+

Distance metric  $< \frac{C}{\pi}$

HyperQPTL model checking

✓ Is a cause

✗ Not a cause



# Conclusion

