Causality and Explainability for Black Box software Hana Chockler causaLens and Department of Informatics King's College, London causa**Lens** 

# Motivation: Modern computerized systems are huge and difficult to understand



# Motivation: Modern computerized systems are huge and difficult to understand





# Actual Causality

# A theoretical concept from AI Extends causal counterfactual reasoning

Turns out to be very useful!



# <u>Intractable</u> - but there are efficient approximation algorithms and sufficient partial solutions



# Formal Verification Is the system correct?



# Formal Verification Is the system correct?





## Counterexamples in hardware

#### A huge timing diagram that is very difficult to understand

• 🔀 Vivado 2017.2.1								
<u>F</u> ile <u>E</u> dit <u>T</u> ools	<u>W</u> indow La <u>v</u> o	ut <u>V</u> iew <u>R</u> un	Help Q- Qu	ick Access				
E = 1 × 1 ≠ 1 ≤ 1 / × 1 + 1 + 1 1000 ns ∨ Ξ    C								
SIMULATION - Simulation	n – sim							? ×
Scope × Sources		_ 0 8	Objects	?	_ 🗆 🖒 X	Untitled 1*		? 🗆 🖒 X
Q   ¥   €	Q ≚ ≑ 🌣		Q 🌣			Q 💾 🤤	Q   👷   🔸	•   •
Name ~ 🏮 top	Design Unit top(RTL)	Block Type VHDL En	Name	Value 0	Data ^ Logic			
v i top i sub_inst	top(RTL) sub(RTL)	VHDL En VHDL En	通 B 通 C 通 Q し 、 comb	1 1 0 1	Logic Logic Logic	Logic   Name   Value   ns   200 ns   400 ns     Logic   1   0<	0 ns 1200 ns 1200 ns 1,200 ns	



next(!START Until (STATUS\_VALID and READY))

works and is really useful!



#### Reasoning about black-boxes





We can reason about various properties of the system without opening the black box

## Explanations for Deep Neural Network's decisions



## Subtle misclassification - uncovered by explanations



# Reinforcement learning - causal simplification of policies





#### Original policy

Simplified policy





# Proposed regulatory framework





Mortgage approval



# Bibliography

- Chockler and Halpern. "Responsibility and Blame: A Structural-Model Approach". J. Artif. Intell. Res. 22: 93-115 (2004)
- Beer, Ben-David, Chockler, Orni, Trefler. "Explaining Counterexamples Using Causality". FMSD (2012)
- Aleksandrowicz, Chockler, Halpern, Ivrii. "The Computational Complexity of Structure-Based Causality". AAAI'14: 974-980.
- Alrajeh, Chockler, Halpern. "Combining Experts' Causal Judgments". Artif. Intell. (2020).
- Sun, Chockler, Huang, Daniel Kroening. "Explaining Image Classifiers Using Statistical Fault Localization". ECCV'20: 391-406.
- Chockler, Kroening, Sun. "Explanations for Occluded Images". ICCV'21: 1234-1243.
- Pouget, Chockler, Sun, Kroening. "Ranking Policy Decisions". NeurIPS'21.
- McNamee and Chockler: Causal policy ranking. OSC ICLR'2022.
- Chockler and Halpern. "On Testing for Discrimination Using Causal Models". AAAI'22.

