

Why do things go wrong (or right)? Applications of causal reasoning to verification

Hana Chockler

causaLens

and

Department of Informatics
King's College, London

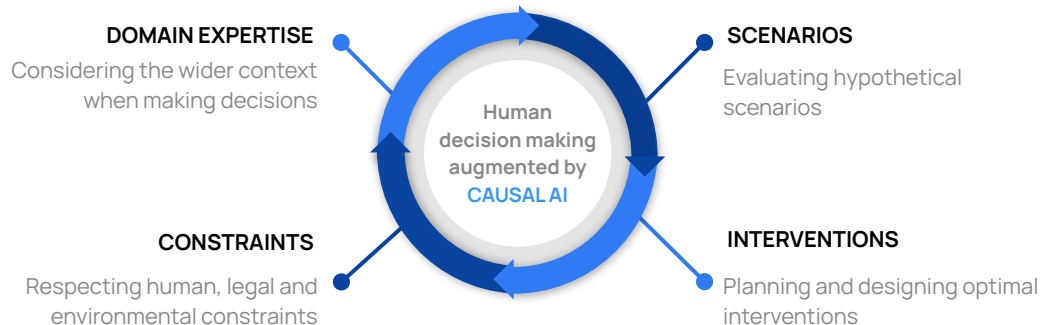


Humans trust Causal AI with complex decisions

Correlation ML systems learn
to perform simple predictions

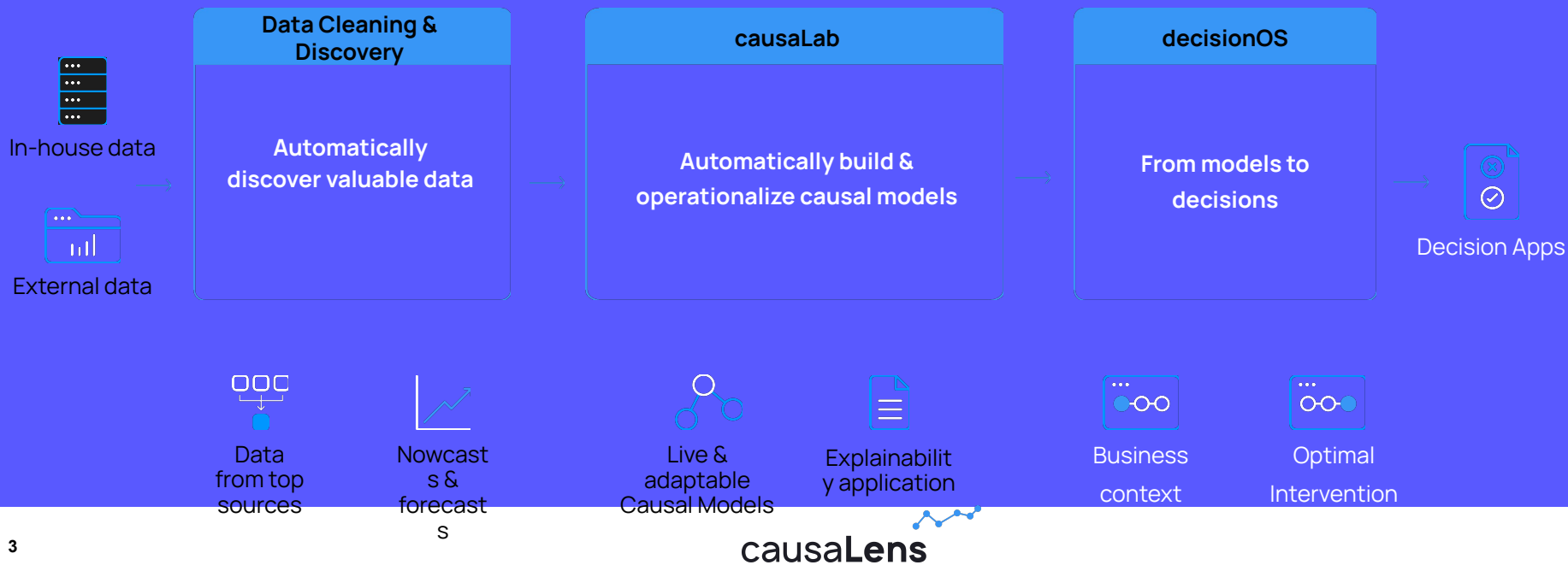
But predictions are a very small
element of decision making.

Causal AI is the only technology that can augment human decision making



World's First Full-Stack Causal AI Platform

We launched the World's First Causal AI Enterprise Platform, which automates everything from Raw Data to Improved Business Decisions.



Motivation: Modern computerized systems are huge and difficult to understand



Motivation: Modern computerized systems are huge and difficult to understand

Can we understand and fix errors?

What does the system do?

black
box

Can we be sure it is correct?

Motivation

Modern computerized systems are huge and difficult or even impossible to understand

Can we understand and fix errors?

What does the system do?

black box

Can we be sure it is correct?

black box

Deep Neural Networks

Actual Causality

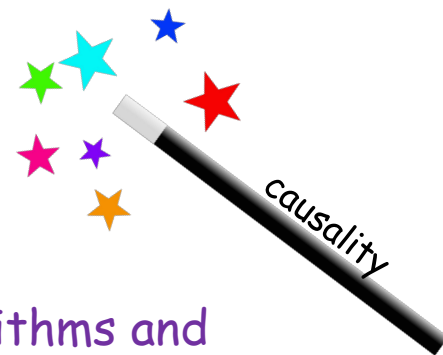
A theoretical concept from AI
Extends causal counterfactual reasoning

+

Quantification of causality,
allowing to rank causes by importance

©Chockler & Halpern, 2004

Turns out to be very useful!



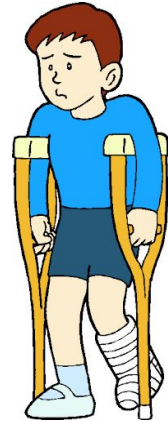
Intractable - but there are efficient approximation algorithms and sufficient partial solutions

A-priori (type) and a-posteriori (actual) causality



There is a terrible pollution, so
my next patient is likely to
suffer from breathing problems

Turns out he broke his
leg



Background

:

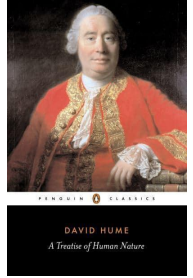


Causality

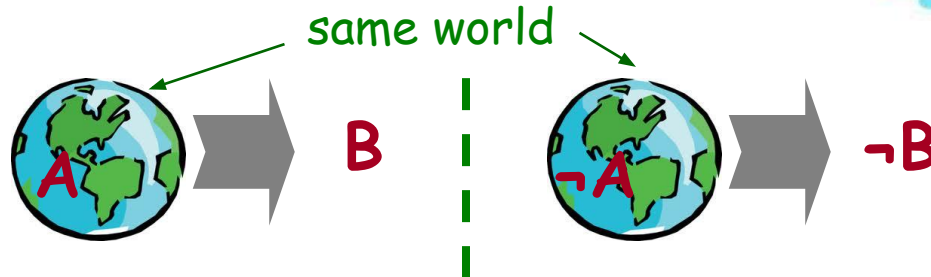
When do we say that **A** is a cause of **B**?

Common approach: **counterfactual causality**.

A is a **cause** of **B** if, had **A** not happened, then **B** would not have happened.



Rain is a cause of me being drenched.



Causality

When do we say that **A** is a cause of **B**?

Common approach: **counterfactual causality**.

We need to capture more complex causal connections!

over
determination



~~Rain is a cause of me
being drenched.~~



Causality

When do we say that **A** is a cause of **B**?

Common approach: counterfactual causality.

We need to capture more complex causal connections!

preemption



Car is a cause of me being drenched, but not the rain

Actual causality

Extends the counterfactual reasoning
by having expressive causal models
allowing overdetermination, preemption,
and complex causal structures

Overdetermination: A is a cause of B if there exists some contingency C
(change in the current world)
in which B counterfactually depends on A .

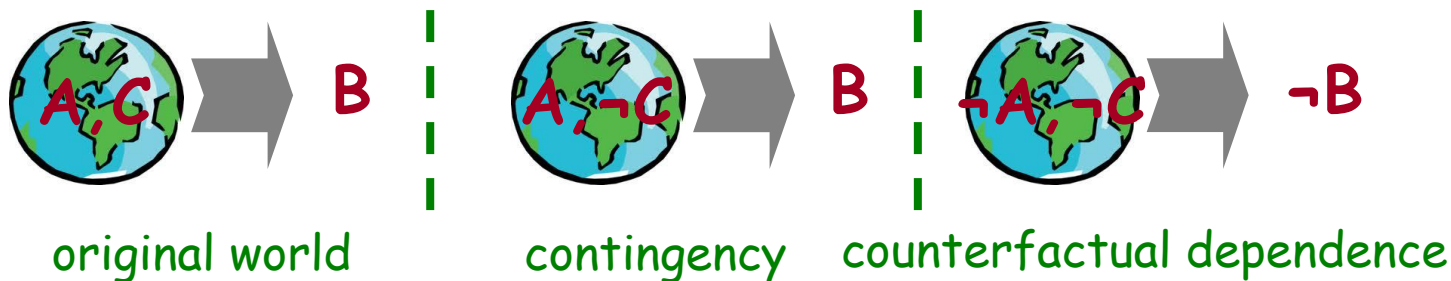


Illustration of overdetermination in actual causality



Rain is an actual cause of me being drenched.



Contingency = the car



Rain is
a counterfactual cause

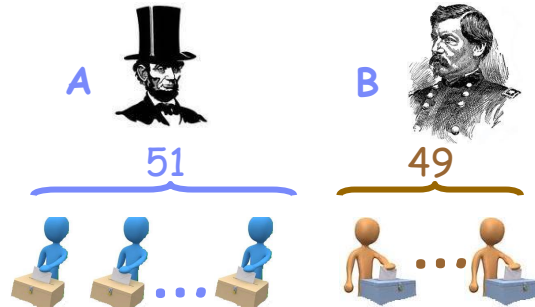


Responsibility: a quantitative measure of causality

Voting example

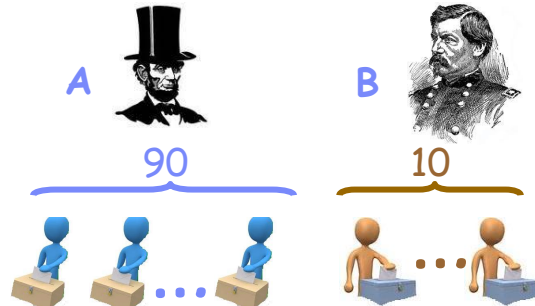


California:



Each **blue** voter is a cause of Lincoln's win

New York:

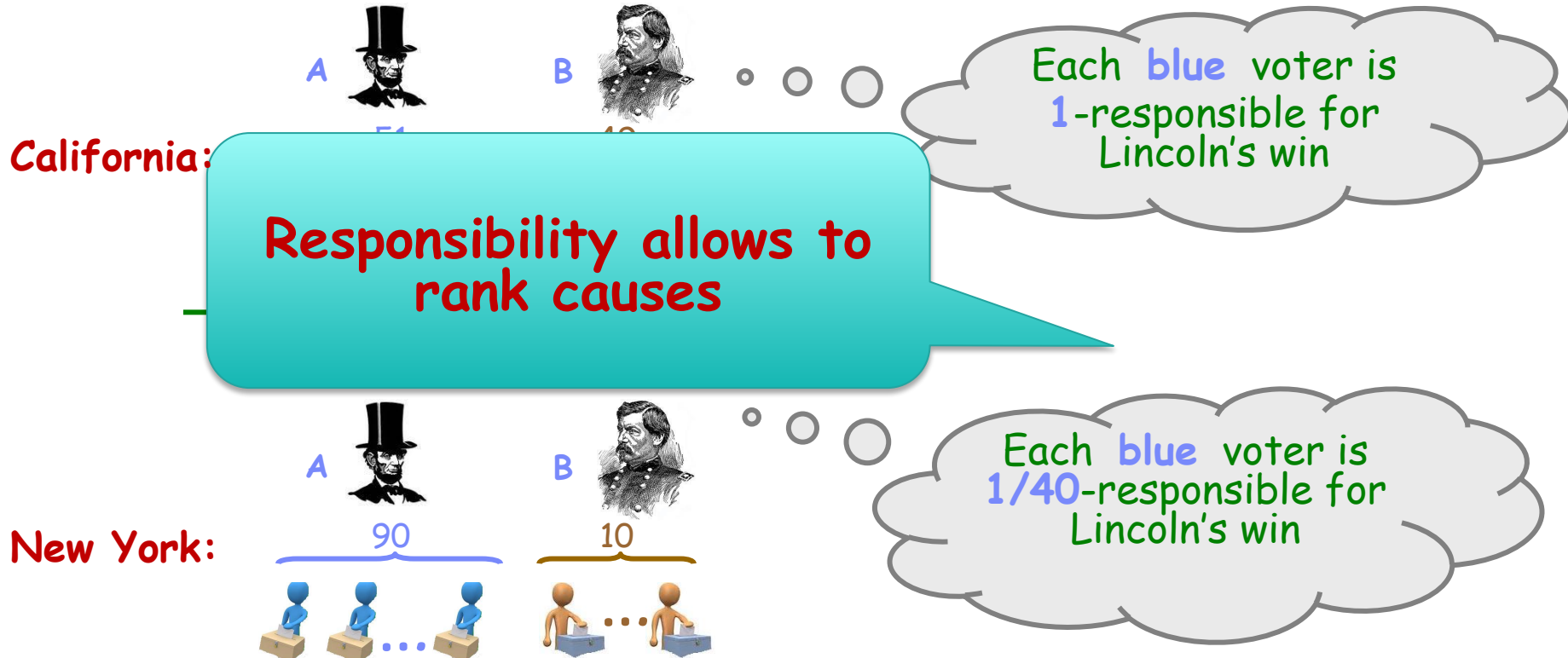


We need to distinguish between the cases!

Each **blue** voter is a cause of Lincoln's win

Responsibility: a quantitative measure of causality

Voting example



Complexity of Computing Causality and Responsibility

Causality:

- ◆ Σ_2 - **complete** for singleton causes
- ◆ D_2 - **complete** in general case

D_2 is the difference class of Σ_2 and Π_2

Responsibility:

- ◆ $\text{FP } \Sigma_2 [\log(n)]$ complete.

INTRACTABLE

Complexity of Computing Causality and Responsibility

Causality:

- ◆ Σ_2 - complete for singleton causes
- ◆ D_2 - complete in general case

Responsibility:

- ◆ $FP^{\Sigma_2[\log(n)]}$ - complete.

INTRACTABLE

The good news:

- ◆ There are linear-time approximation algorithms
 - o Accurate on most problems
- ◆ We usually care only about highest-ranked causes
 - o Polynomial to compute the exact set



Motivation

Modern computerized systems are huge and difficult or even impossible to understand

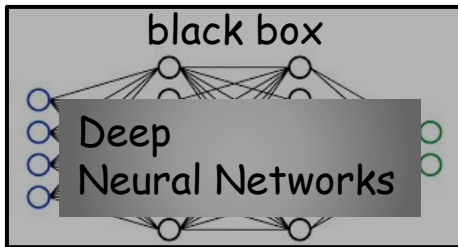
Can we understand and fix errors?

Can we be sure it is correct?

verification

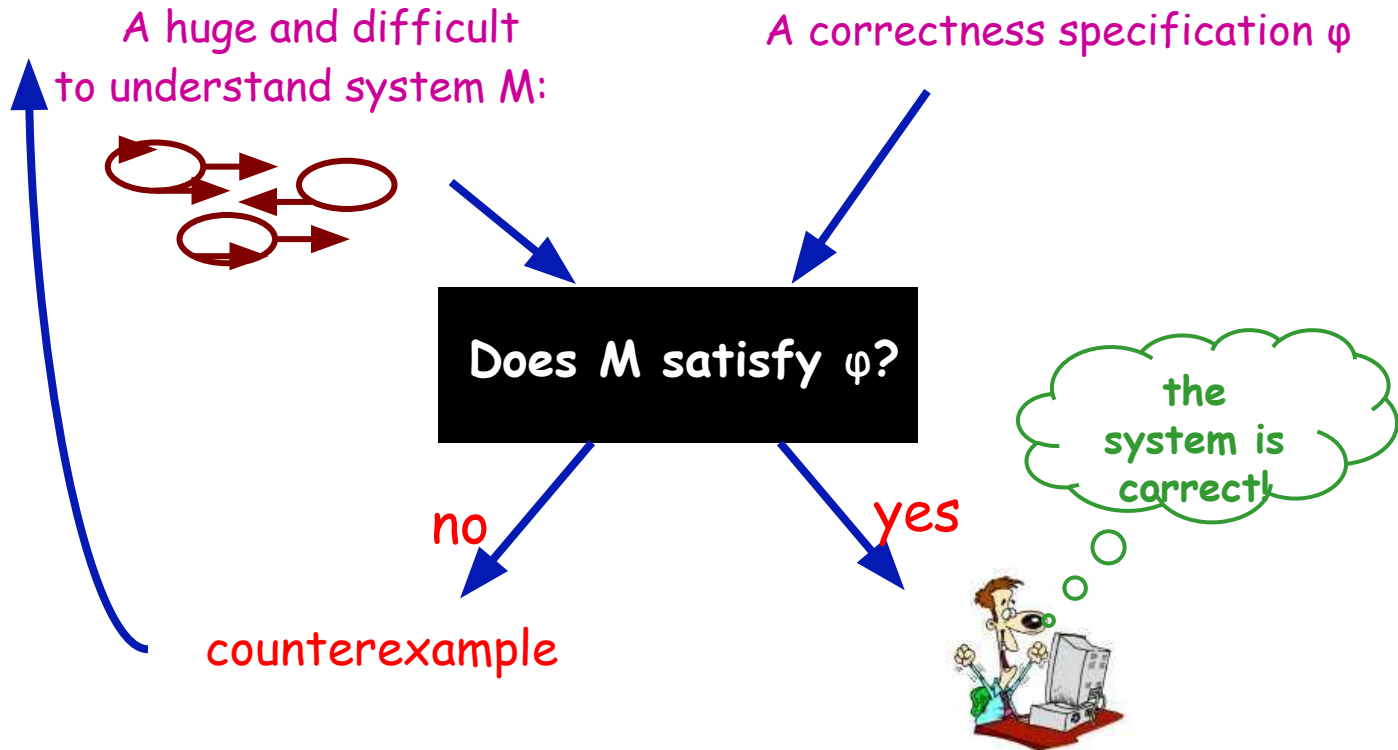
black box

What does the system do?



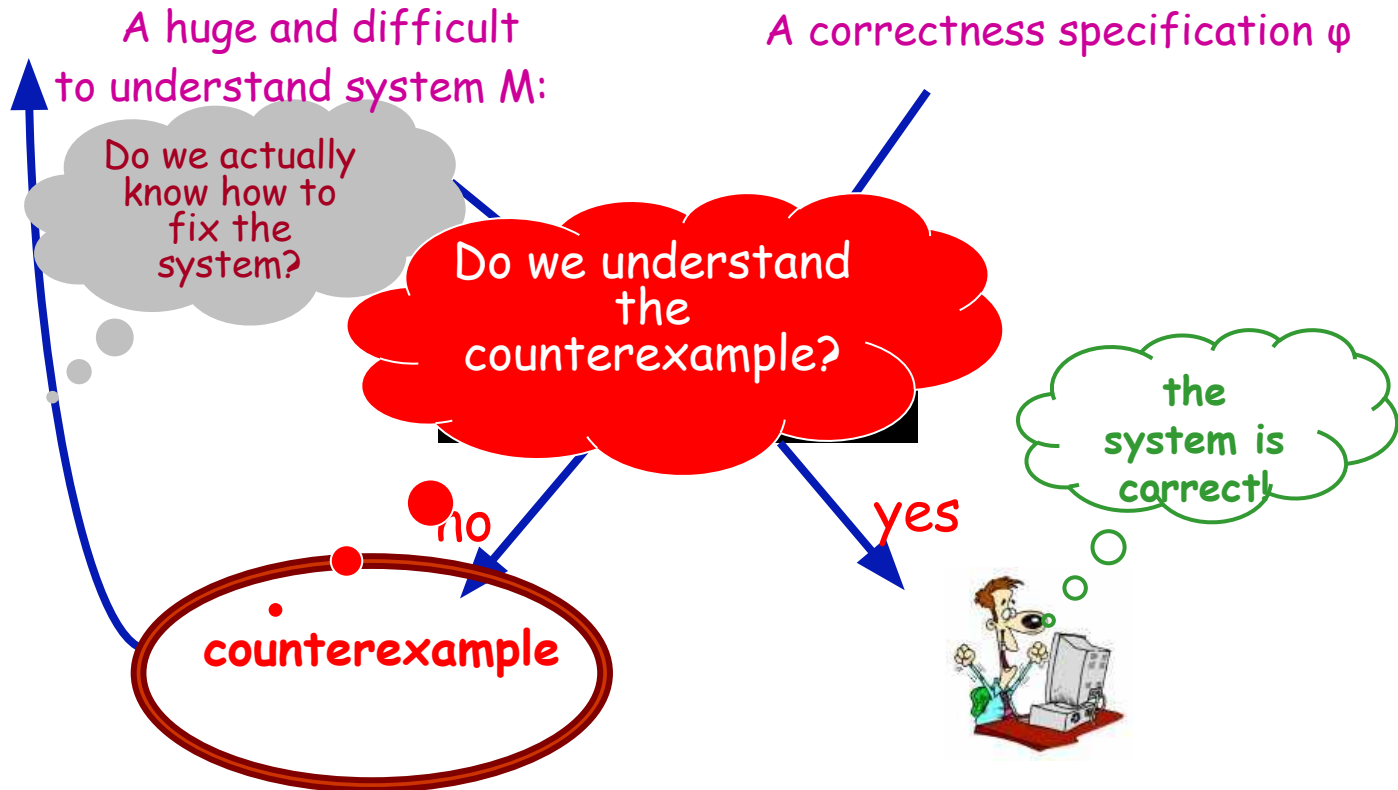
Formal Verification

?Is the system correct

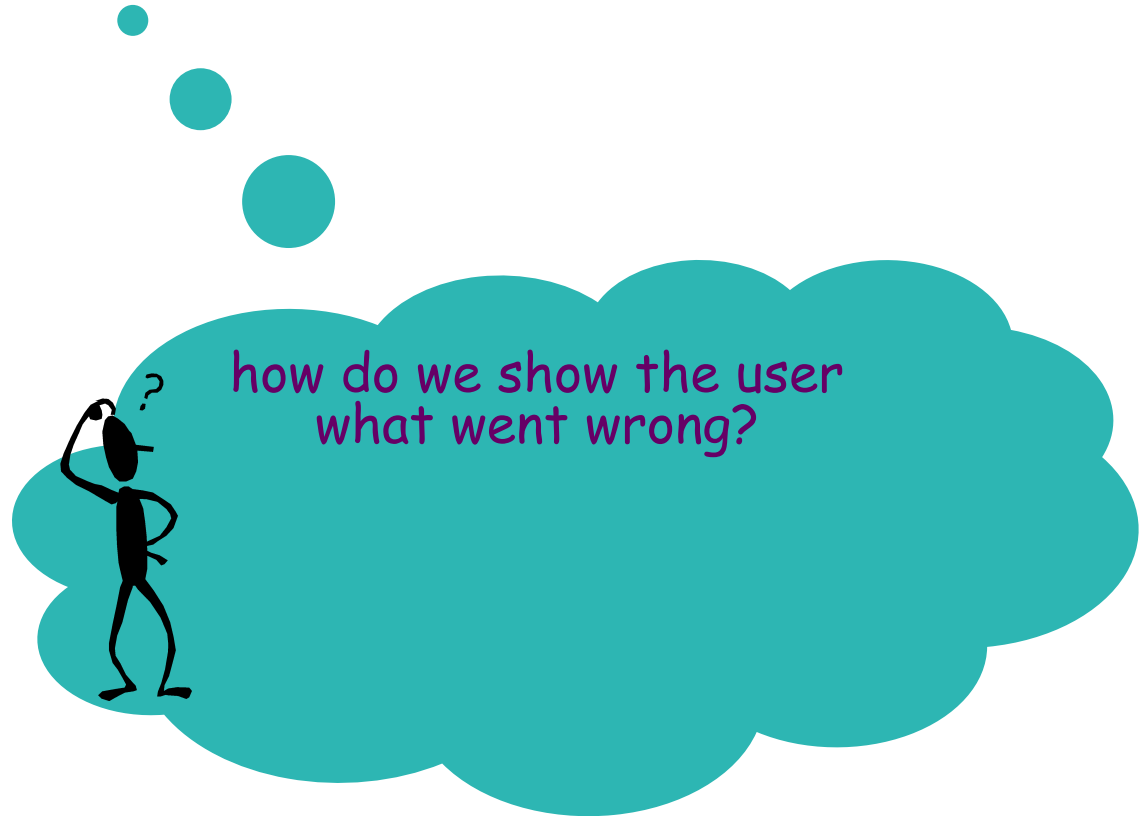


Formal Verification

?Is the system correct

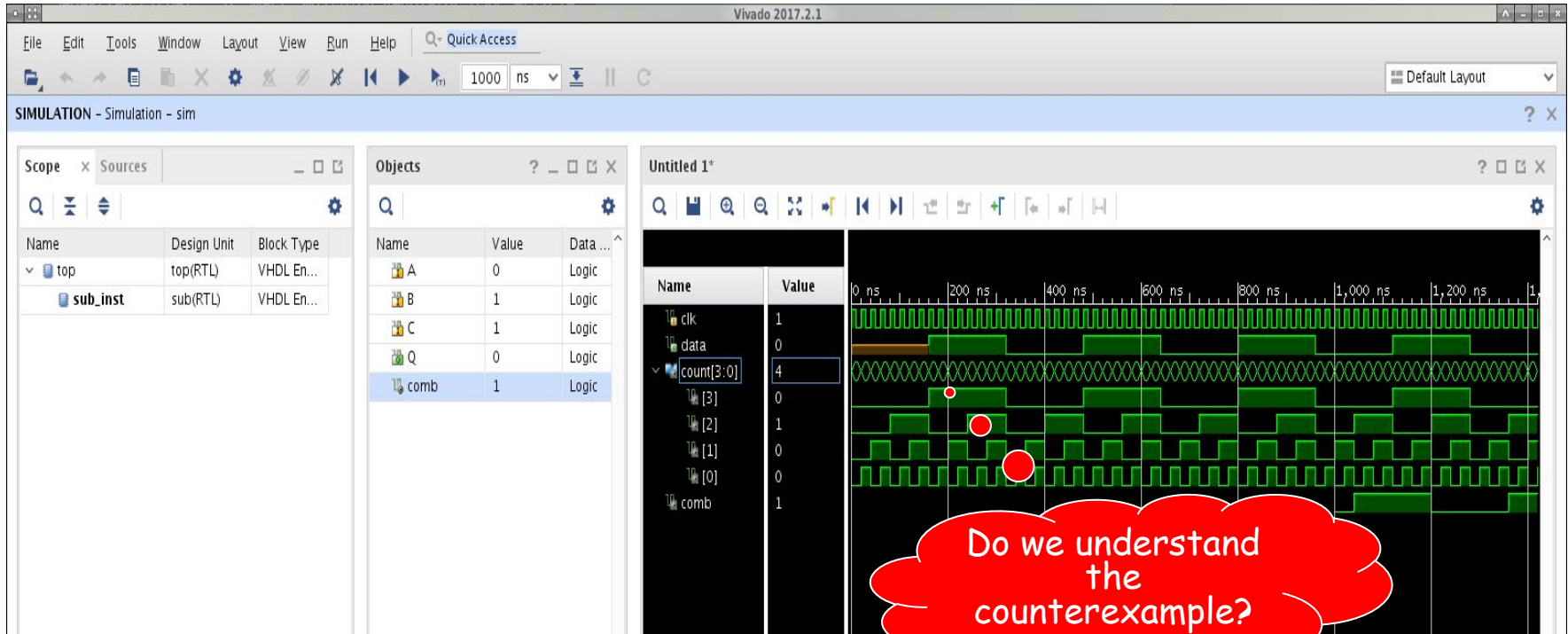


Counterexamples



Counterexamples in hardware

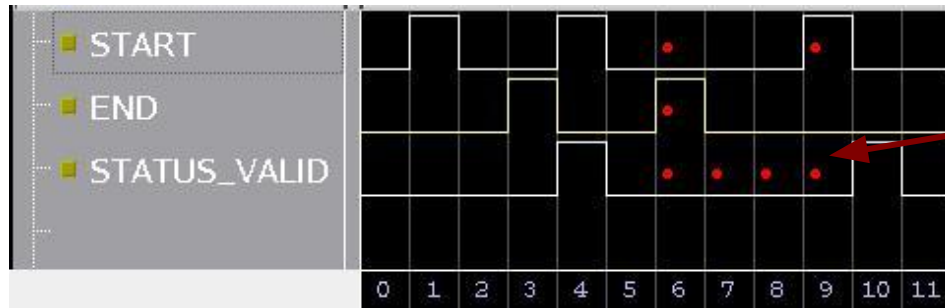
A huge timing diagram that is very difficult to understand



Explaining counterexamples using causality (Red Dots) part of IBM tool



A timing diagram of a buggy hardware execution

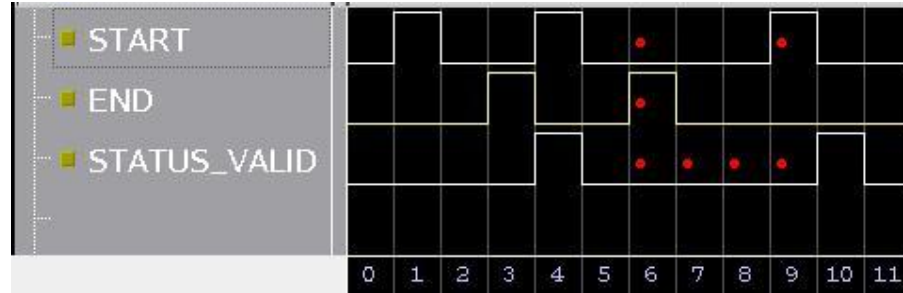


causes
marked as
red dots

$\varphi = \text{always } ((\neg \text{START} \text{ and } \neg \text{STATUS_VALID} \text{ and } \text{END}) \rightarrow \text{next}(\neg \text{START} \text{ Until } (\text{STATUS_VALID} \text{ and } \text{READY})))$

works and is really
useful!

Explaining counterexamples using causality (Red Dots) part of **IBM** tool



Following this work...

Many applications
of causality and
responsibility to
software
engineering

CREST workshop

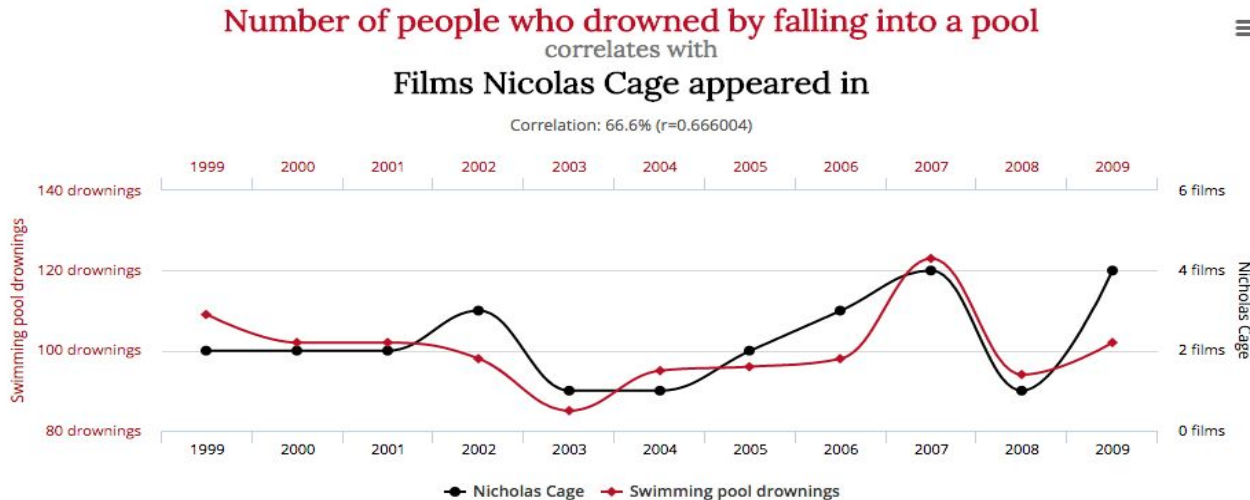


Ongoing work:
causal debugging
for software

Explanation of faults in software testing - SOA

◆ Statistical Analysis for Fault Localisation

- o Looks for correlation - elements that appear more in failing traces than in passing ones are suspicious
- o Elements are ordered by their degree of suspiciousness



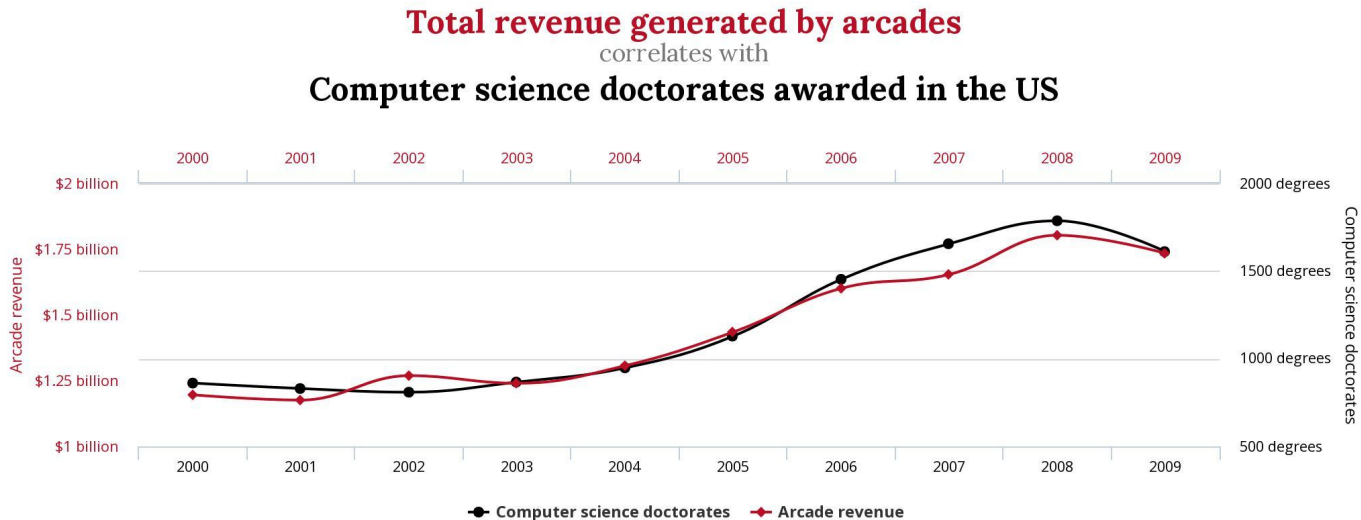
Data sources: Centers for Disease Control & Prevention and Internet Movie Database

<http://www.tylervigen.com/spurious-correlations>

Explanation of faults in software testing - SOA

◆ Statistical Analysis for Fault Localisation

- o Looks for correlation - elements that appear more in failing traces than in passing ones are suspicious
- o Elements are ordered by their degree of suspiciousness

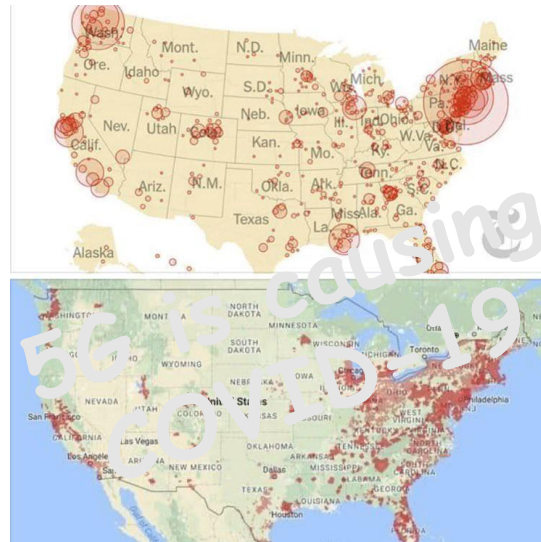


<http://www.tylervigen.com/spurious-correlations>

Explanation of faults in software testing - SOA

◆ Statistical Analysis for Fault Localisation

- o Looks for correlation - elements that appear more in failing traces than in passing ones are suspicious
- o Elements are ordered by their degree of suspiciousness

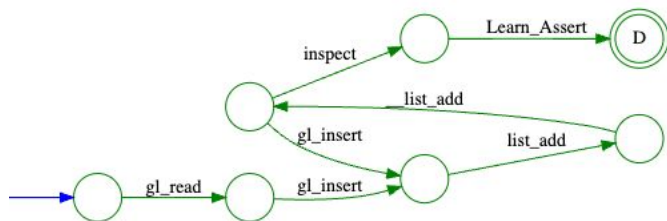


Explanation of faults in software testing - SOA

◆ Statistical Analysis for Fault Localisation

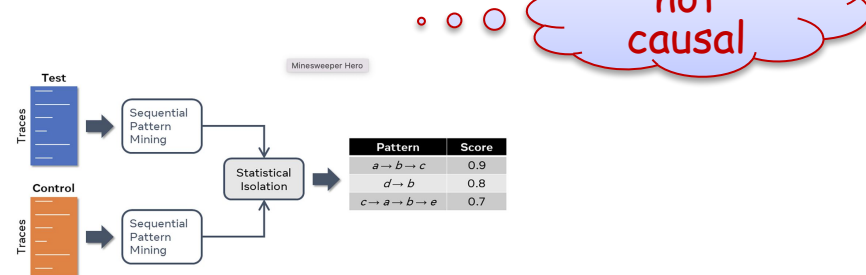
- o Looks for correlation - elements that appear more in failing traces than in passing ones are suspicious
- o Elements are ordered by their degree of suspiciousness

Learning the language of software errors



Recent work from  Meta

Minesweeper automates root cause analysis as a first-line defense against bugs



Motivation

Modern computerized systems are huge and difficult or even impossible to understand

Can we understand and fix errors?

What does the system do?

black box

Can we be sure it is correct?

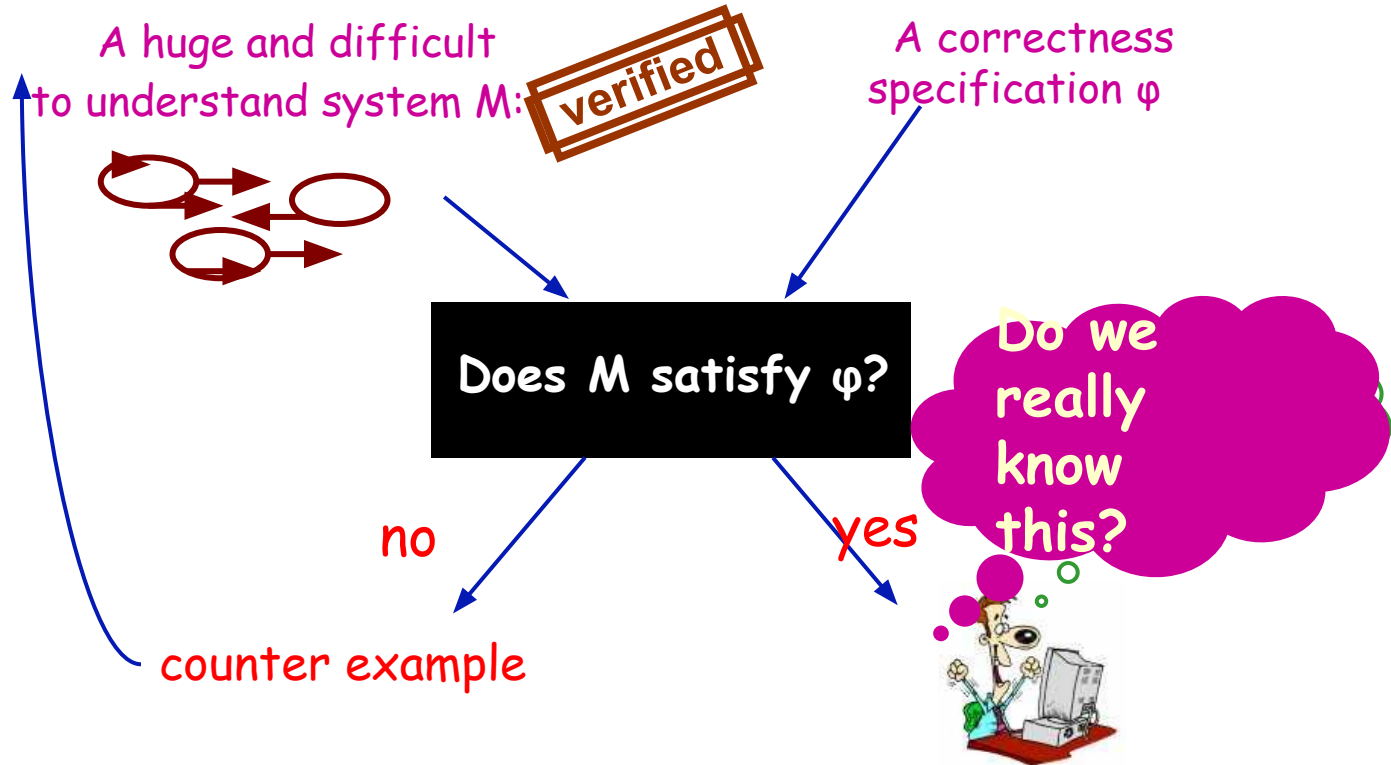
black box

Deep Neural Networks

verification

Formal Verification (Model Checking)

?Is the system correct

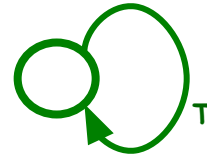


Vacuity - the main idea

Vacuous satisfaction of ϕ in M means that some part of ϕ is irrelevant in M



system
 M :



Printer that
doesn't print

ϕ = always (req \rightarrow eventually grant)

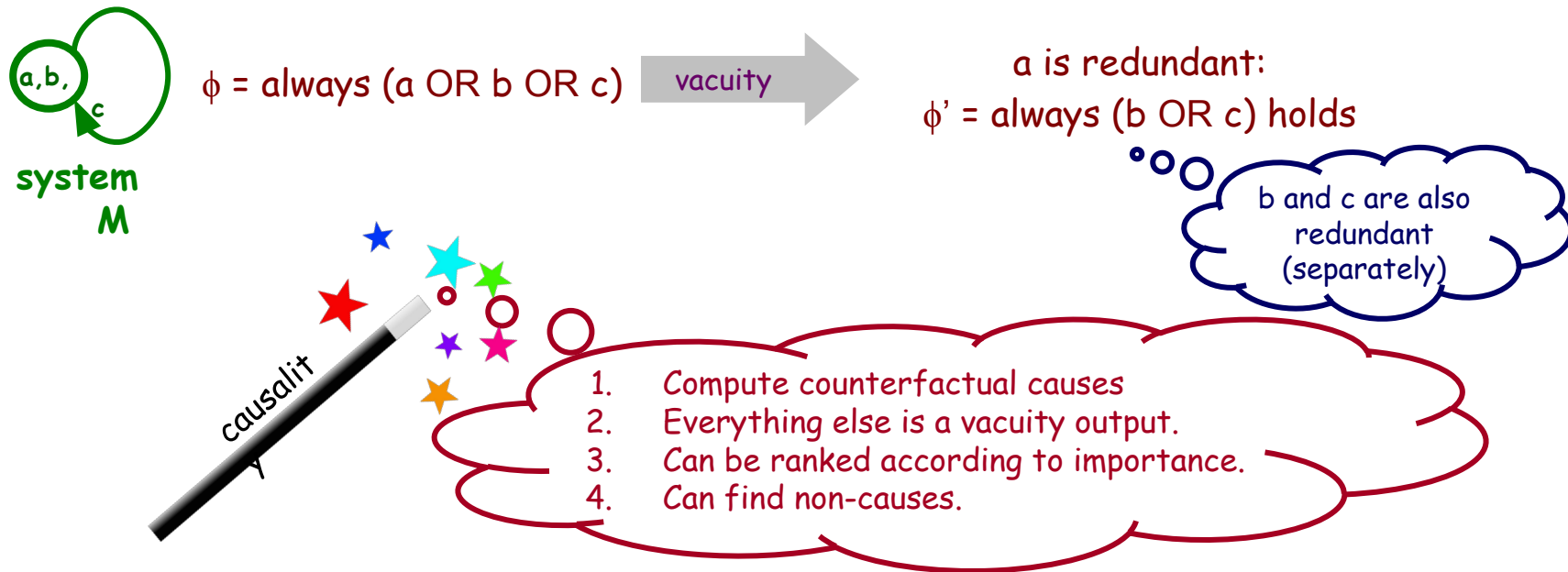


Vacuous pass can
point to problems

What is the output of vacuity check?

Vacuous satisfaction of ϕ in M means
that some part of ϕ is irrelevant in M

Standard vacuity checks output (some) redundant parts of ϕ



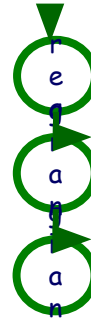
Coverage - the main idea

Low coverage of M by ϕ means that some part of M is irrelevant for the satisfaction of ϕ



system
 M :

Printer that prints
everything twice



ϕ = always (req \rightarrow eventually grant)



Low coverage can
point to problems

What is the output of coverage check?

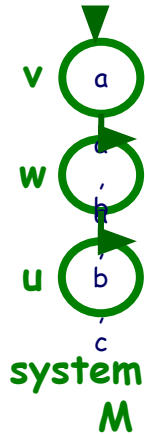
Low coverage of M by φ means that some parts of M are irrelevant for the satisfaction of φ

There is no standard coverage check... but if there was one...

What is the output of coverage check?

Low coverage of M by ϕ means that some parts

Impractical for huge and difficult to understand systems
(so there is a good reason it is not done)
... but can be a good idea for small critical systems



There

$\phi = \text{all}$

c is completely irrelevant; a, b are causes

causality

1. Compute counterfactual causes
2. Everything else is not covered.
3. Can be ranked according to importance.
4. Can find non-causes.

Why is verification a good application for causality?



♦ Interventions are always possible

- o An intervention amounts to a change in the value of a variable
- o Unlike other domains, where changes can be impossible (like healthcare)

$x=y$  $x=1$



Now let's replace
you with the same
person but
without peanut
allergy

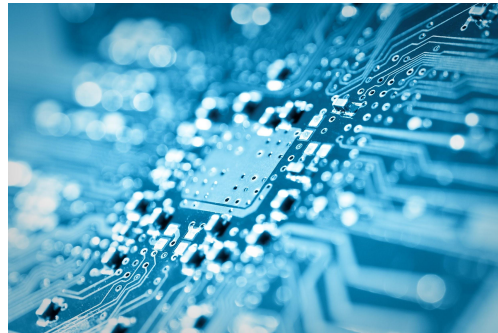


-
- Diagram illustrating a 4x12 grid world environment. The grid is divided into three sections: START (yellow square), END (yellow square), and STATUS_VALID (yellow square). Red dots indicate specific states or goals within the grid.

Why is verification a good application for causality?



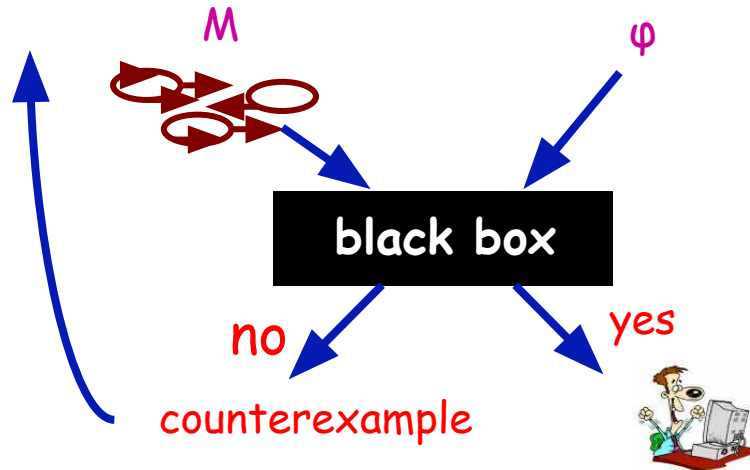
- ◆ Interventions are always possible
- ◆ It is usually clear what are the variables and easy to calculate the equations
- ◆ The systems are deterministic and all variables are known
 - o No noise, no hidden confounders
 - o Not quite true for concurrent systems, but still better than in other domains



Why is verification a good application for causality?



- ◆ Interventions are always possible
- ◆ It is usually clear what are the variables and easy to calculate the equations
- ◆ The systems are deterministic and all variables are known
- ◆ The approach is agnostic to the model-checking algorithm



Motivation

Modern computerized systems are huge and difficult or even impossible to understand

Can we understand and fix errors?

black box

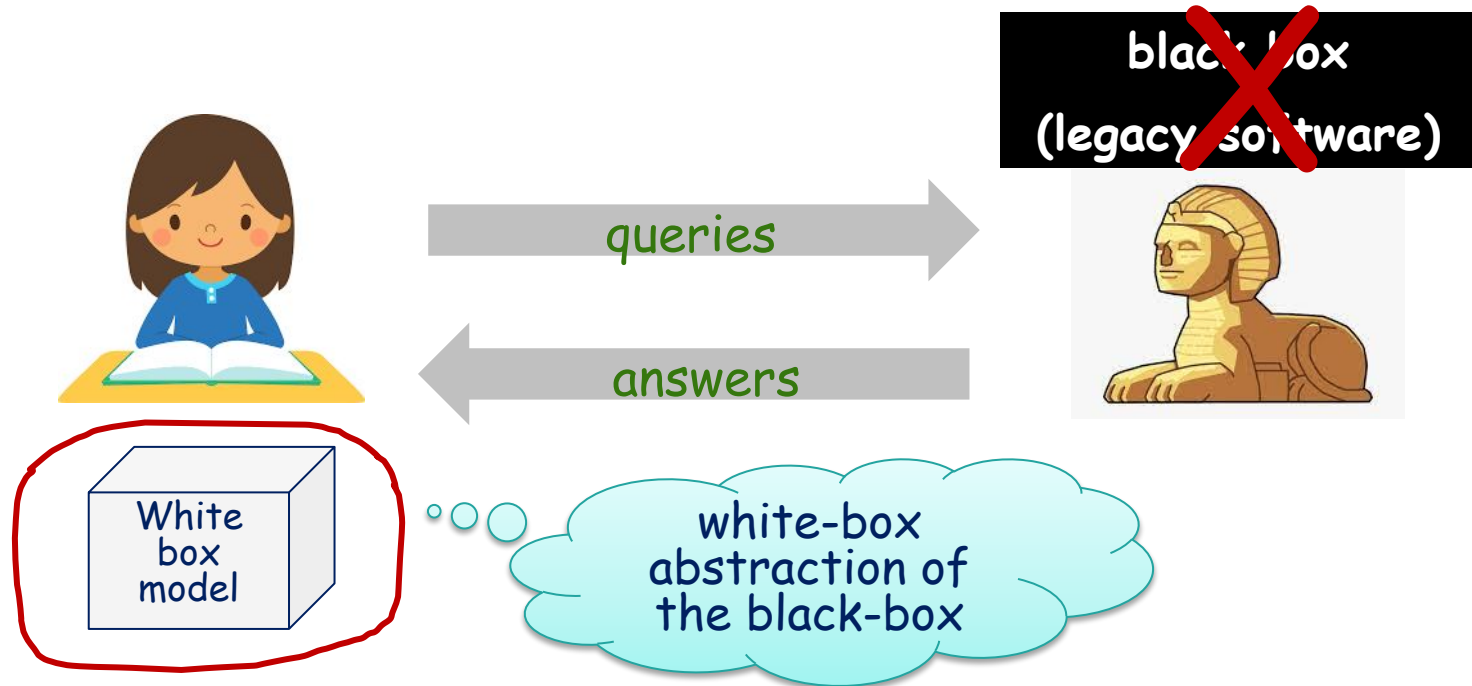
What does the system do?

Can we be sure it is correct?

black box

Deep Neural Networks

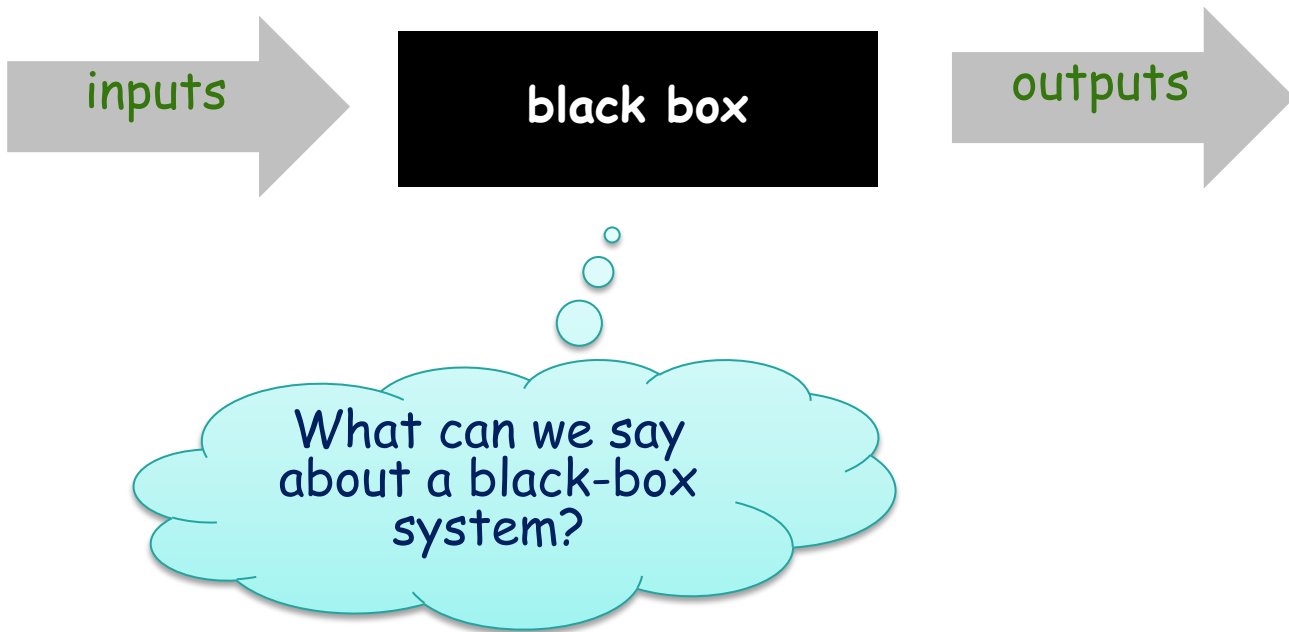
Model learning



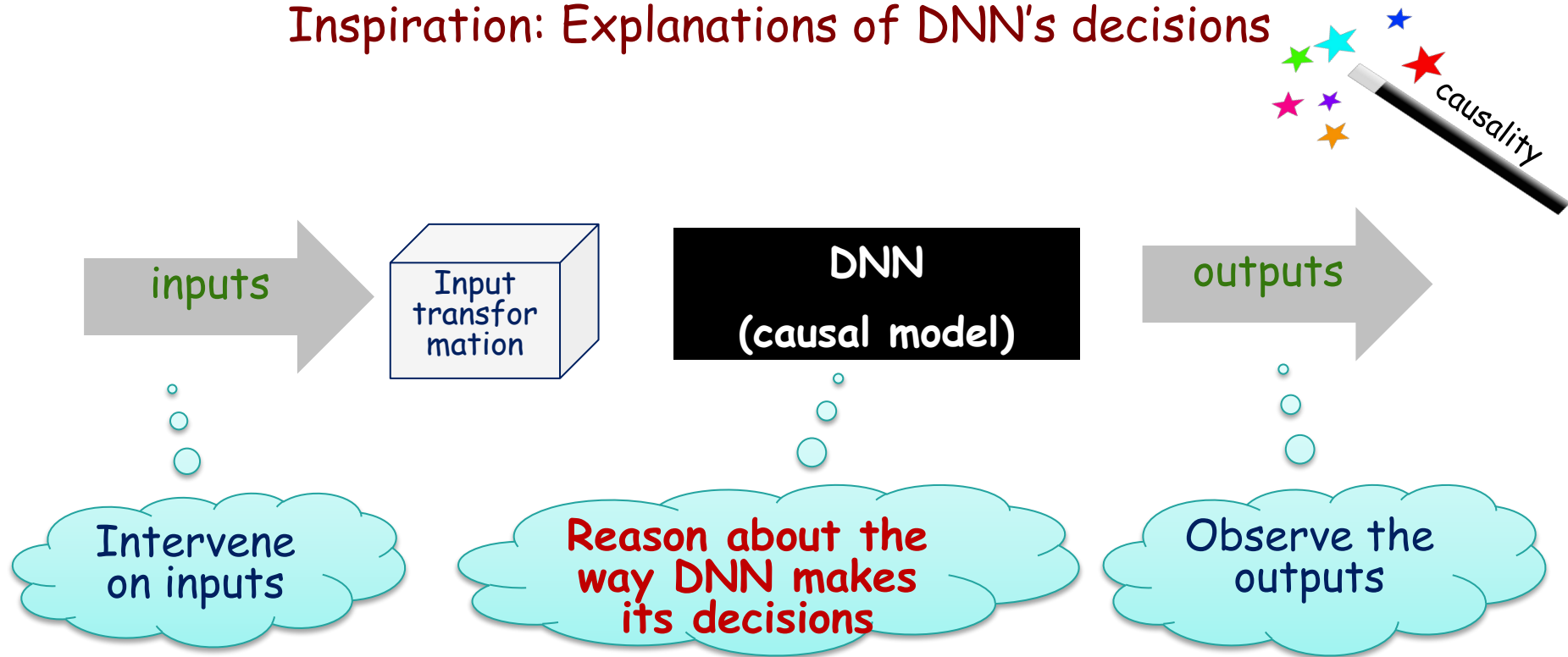
Can be viewed as a causal model

Reasoning about black-boxes

?Do we need to construct a white box at all



Inspiration: Explanations of DNN's decisions



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

Explanations for Deep Neural Network's decisions



**DNN for
classifying animals**



red panda

Causal explanation:
minimal, sufficient,
non-trivial subset of
the pixels of the image

**Because
of this part:**



Subtle misclassifications - uncovered by explanations



DNN for
classifying images



cowboy hat

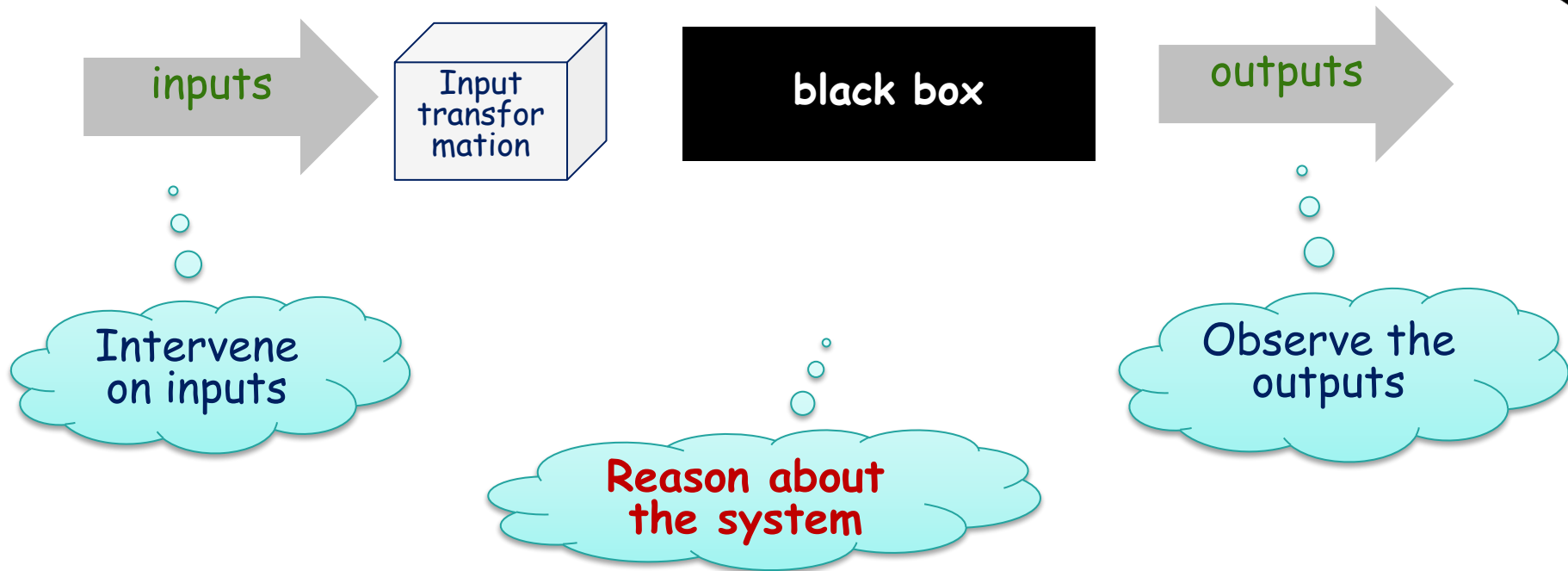
seems
ok

Explanation
uncovered
misclassification!

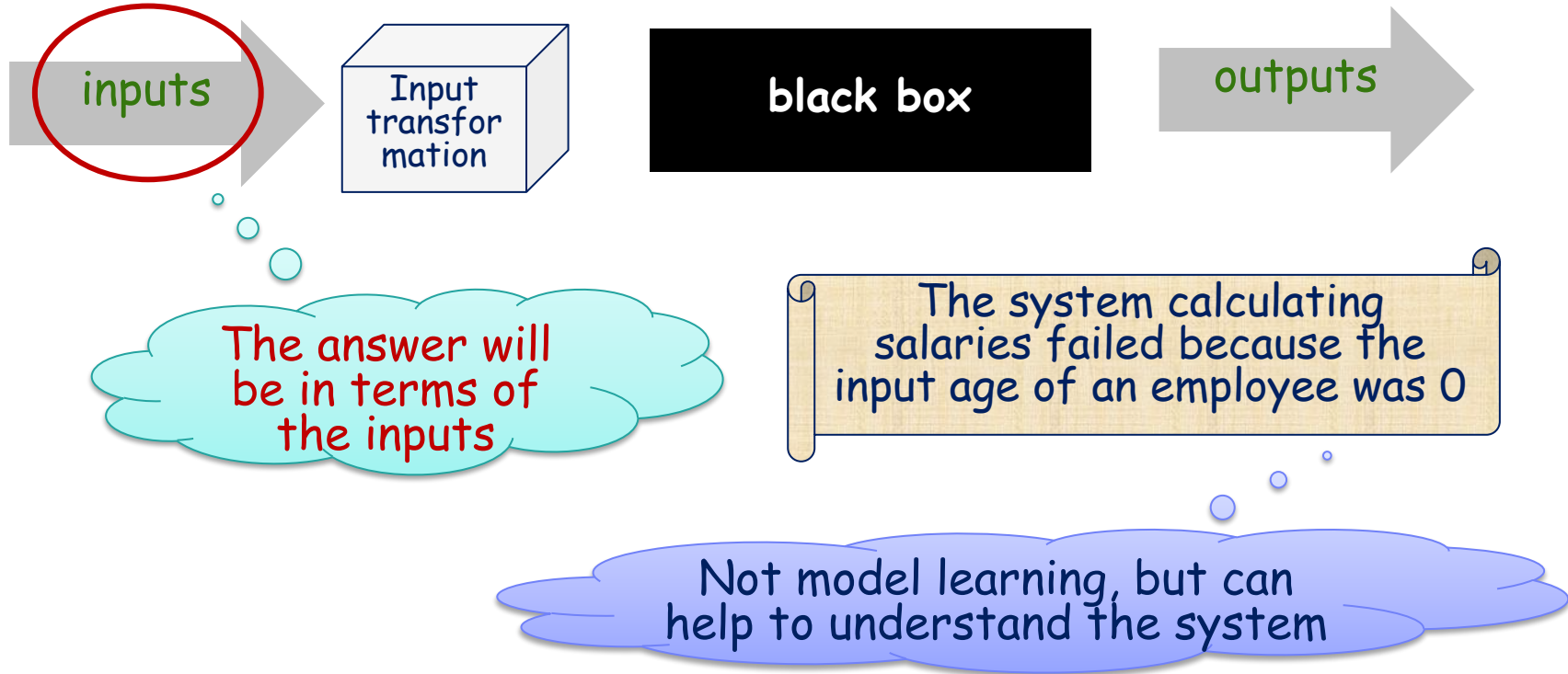
Because
of this part:



Can we use a similar approach to answer the question
? "What does the system do"

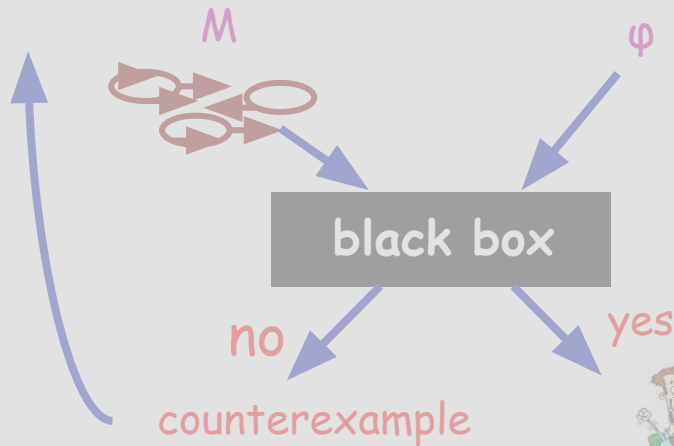


Can we use a similar approach to answer the question
? "What does the system do"



Bibliography

- ◆ Chockler and Halpern. "Responsibility and Blame: A Structural-Model Approach". J. Artif. Intell. Res. 22: 93-115 (2004)
- ◆ Chockler, Halpern, Kupferman. What causes a system to satisfy a specification? ACM Trans. Comput. Log. 9(3): 20:1-20:26 (2008)
- ◆ Beer, Ben-David, Chockler, Orni, Trefler. "Explaining Counterexamples Using Causality". FMSD (2012)
- ◆ Chockler, Gurfinkel, Strichman. Beyond vacuity: towards the strongest passing formula. FMSD (2013)
- ◆ 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).
- ◆ Chockler, Kesseli, Kroening, Strichman. "Learning the Language of Software Errors". J. Artif. Intell. Res. 67: 881-903 (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.



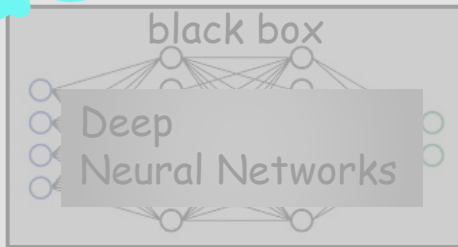
Questions?

Suspecting positive answers

Explaining counterexamples



The system calculating salaries failed because the input age of an employee was 0



Reasoning about black boxes

DNN for
classifying images