# Towards True Reproducibility of Findings in Cybersecurity Research

**Emma Tosch, Northeastern University, 6 December 2022**

# Disclaimer
## Design with the spirit of LASER in mind

How to capture

the essence of a keynote

without making it feel like

a speech or lecture?

# Disclaimer

**Design with the spirit of LASER in mind**

# Disclaimer
## Design with the spirit of LASER in mind

# Disclaimer
## Design with the spirit of LASER in mind

Work in progress

Speculation

Incomplete Ideas

Foster Discussion

Selfishly…

# Outline

Oral History of Artifact Evaluation (student perspective)

Evaluators produce replicates

Language design for reproducibility

# Formal languages unlock great power

# Backstory

- September 2022: different workshop…

  - Full disclosure…
    completely forgot!

- 2013/14 — lab mates participated in one of the earliest AECs for SIGPLAN

- 2014 — submitted artifact (OOPSLA 2014)

- 2014 — began AEC review (POPL 2015)

**HOWTO for AEC Submitters**
(http://bit.ly/HOWTO-AEC)
(Last updated May 2022)

**Dan Barowy** (dbarowy@cs.umass.edu) - *now at Williams College*
**Charlie Curtsinger** (charlie@cs.umass.edu) - *now at Grinnell College*
**Emma Tosch** (etosch@cs.umass.edu) - *now University of Vermont*
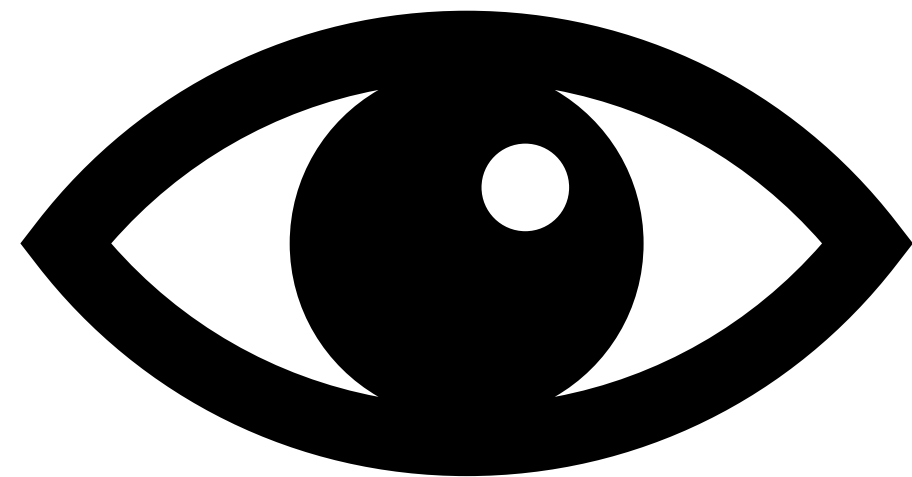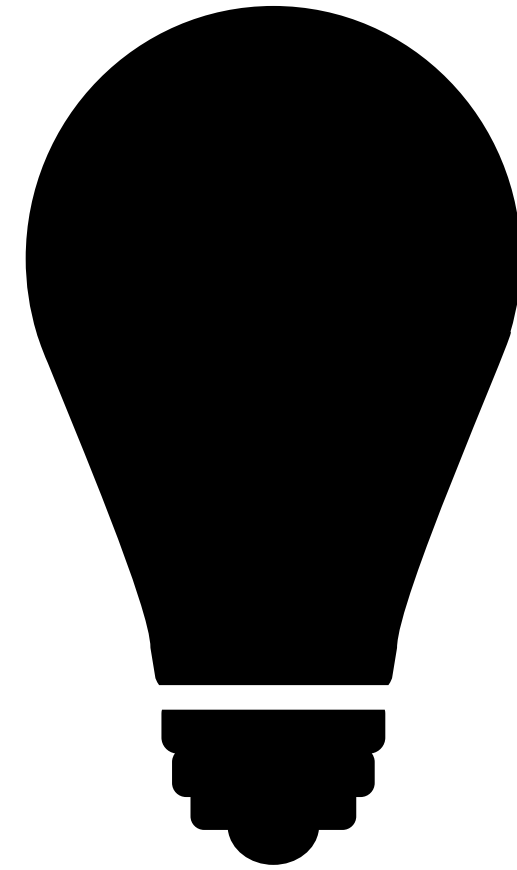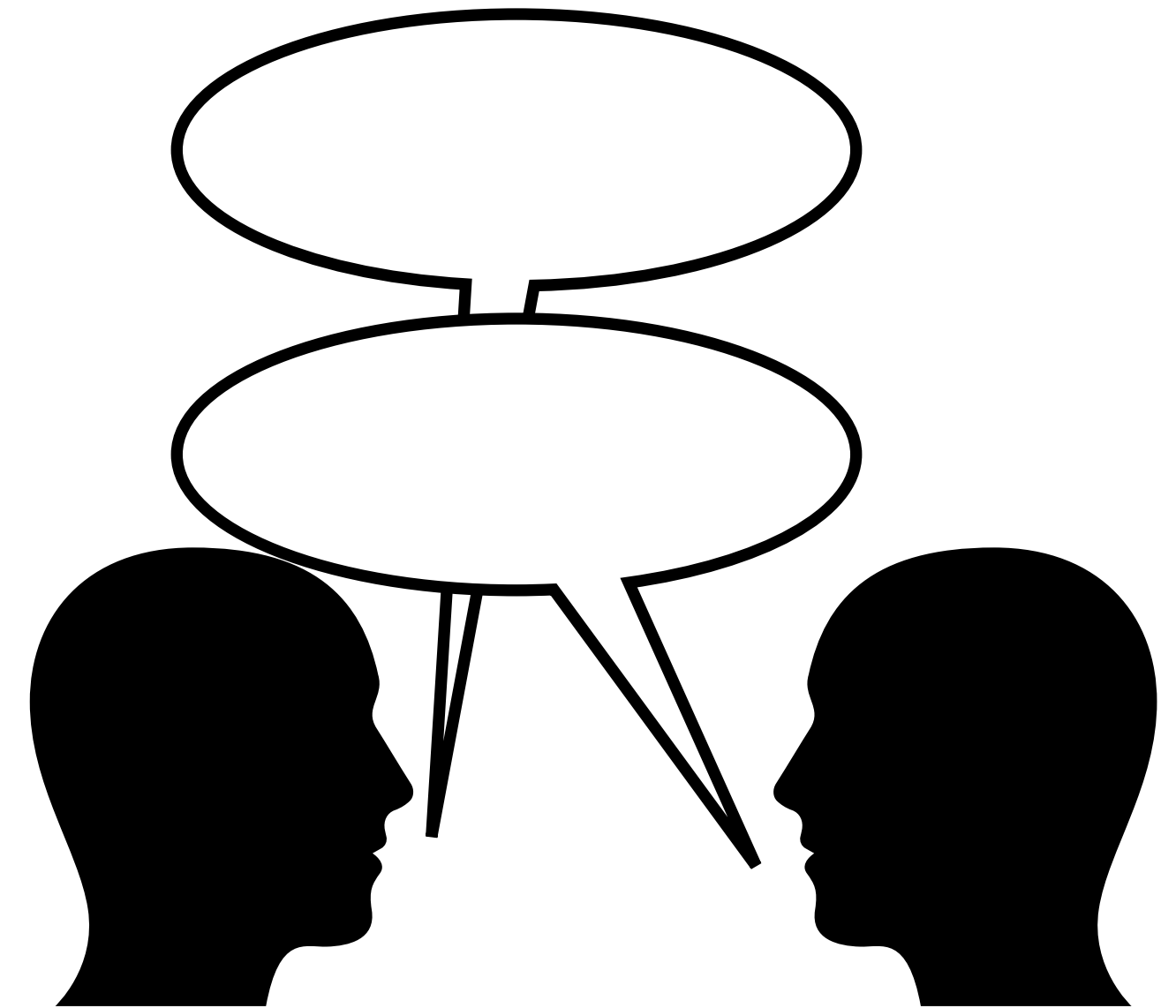**John Vilk** (jvilk@cs.umass.edu) - *now at Stripe*
of the PLASMA group (http://plasma.cs.umass.edu) *at University of Massachusetts Amherst*
with encouragement and support from **Emery Berger** (emery@cs.umass.edu)

After serving on several Artifact Evaluation Committees and winning two Distinguished Artifact Awards, we put together this HOWTO document to help you submit an artifact that will pass the AEC process with flying colors.

**How to Build a Good Software Artifact**

1. Provide documentation with your artifact. We recommend that you prepare a Getting Started Guide. It should explain:
   a. how to download your artifact
   b. how to install your artifact
   c. how to run your artifact
   d. how to compare your artifact's outputs to outputs described in your paper.
2. Explicitly enumerate your claims in both your paper and in your artifact's documentation.
3. Provide a VM if possible, and when appropriate. VMs aid reproducibility because they help control for nuisance factors that are not central to an author's claims, significantly facilitating the review process. Nonetheless, reviewers may need to accept performance tradeoffs for VMs (e.g., because of the absence of special hardware). These tradeoffs are acceptable as long as authors explain to reviewers how and why they should adjust their expectations.
4. Provide step-by-step instructions, but make it easy for reviewers to supply their own inputs to your artifact. When reviewers can "play" with your artifact, it gives them confidence that your ideas were implemented robustly.
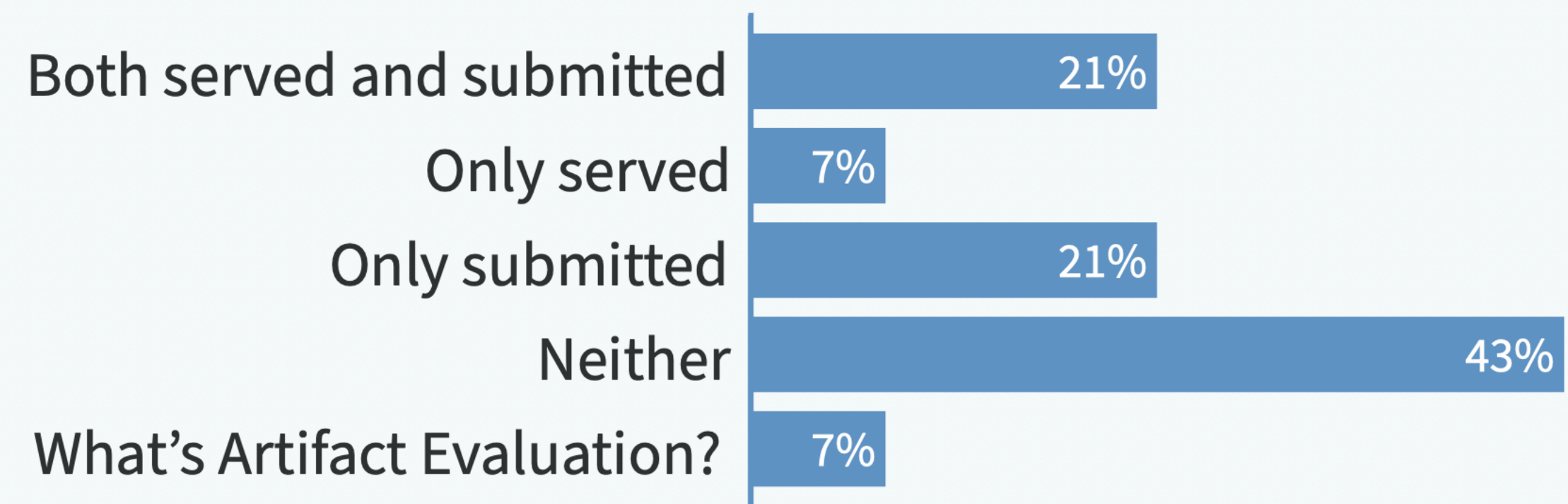
**Source Code**

1. If you are not bound by a nondisclosure agreement, make every effort to supply reviewers with source code. Good reviewers may read and modify your source code to learn the true capabilities of your artifact.
2. Document your code. You should sufficiently explain what is going on so that people who want to build on your work can do so.
3. If you discuss a new algorithm or unique implementation approach in your paper, have a reference to its implementation in the source code.



Bla, bla, bla. Proper name, place name, backstory stuff.

Oral History of Artifact Evaluation (student perspective)

# Have you ever served on an AEC committee or submitted an artifact to an AEC?

Both served and submitted **21%**

Only served **7%**

Only submitted **21%**

Neither **43%**

What's Artifact Evaluation? **7%**

# The idea of submitting to artifact evaluation makes me feel...

supporting **walkies** contributing

**confused** **excited** **nervous**

bored science-y

**rigorous** surprised same

helpful time hopeful

# The idea of serving on an artifact evaluation committee makes me feel..

stressed
overwhelmed
upside
responsible deadweight
busy time clueless need
excited arrogant unprepared
overburdened supportive

# The rise and fall of expectations
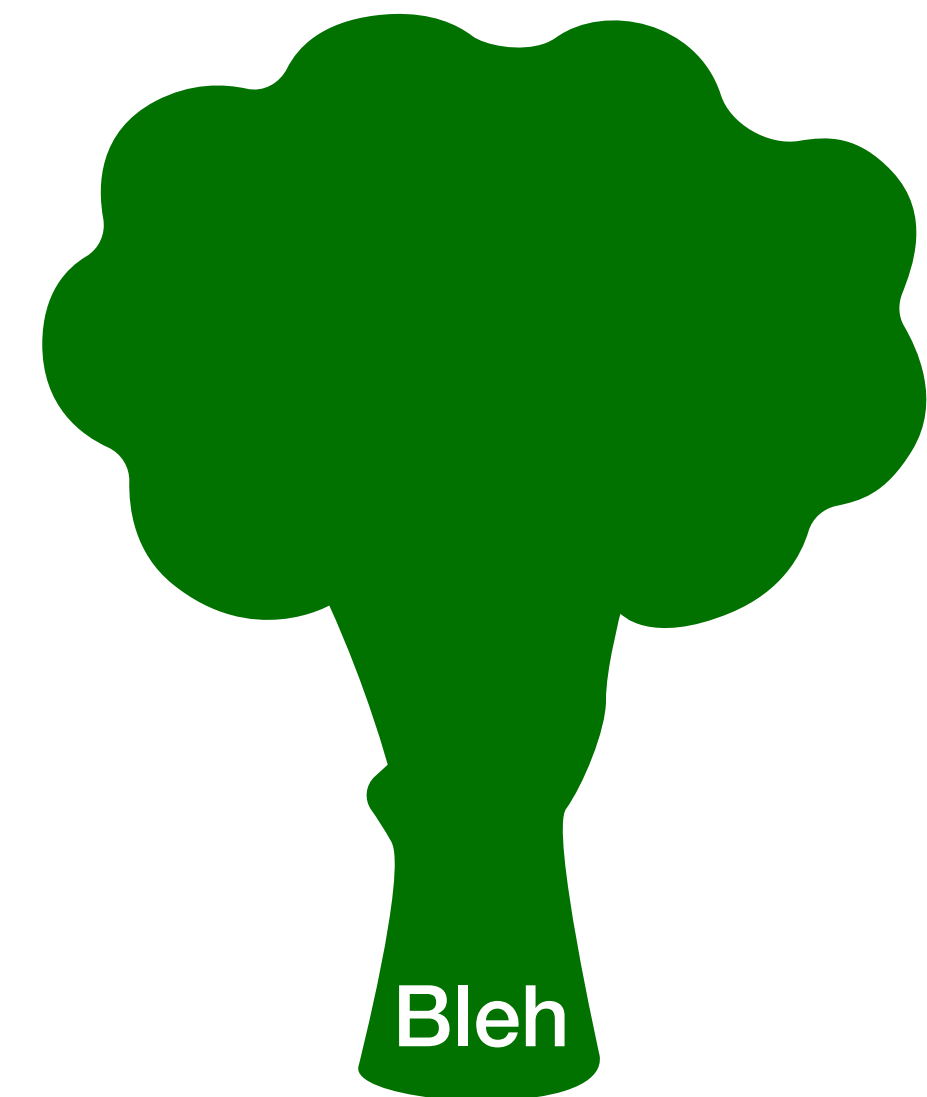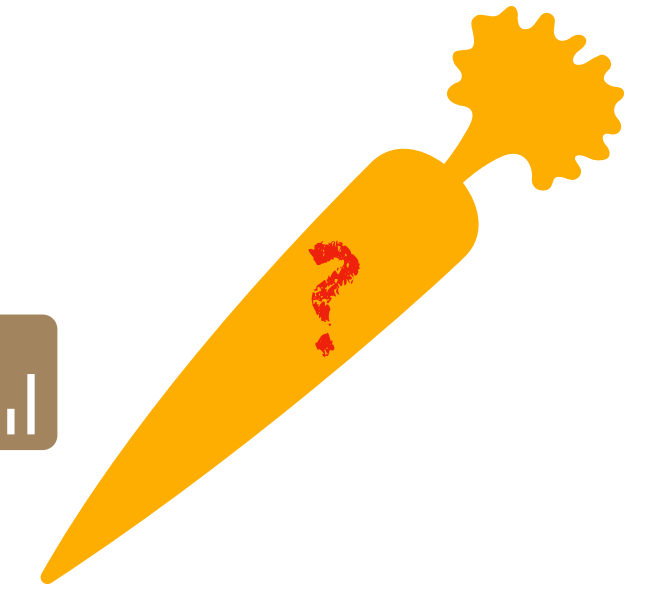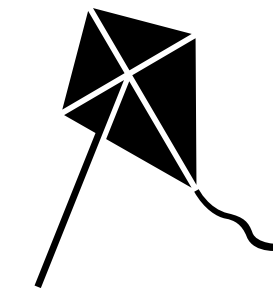
# Student perspective
## Process early on

1. Read abstract, note expectations *set by abstract*

2. Read paper, revise expectations, *in light of the paper*

3. Write out expected software components and datasets*

4. Sketch a plan for something novel to do with the software

5. Early days: no separate guide

# Student perspective
## Reality

- Retrospective: assumed goal was reusability

  - *Then*: one badge. *Now*: Five

- Arguments in favor (at the time)

  - Promote best practices

  - Disincentivize "runs on my machine"

  - Temper reader's expectations (inflated abstracts)
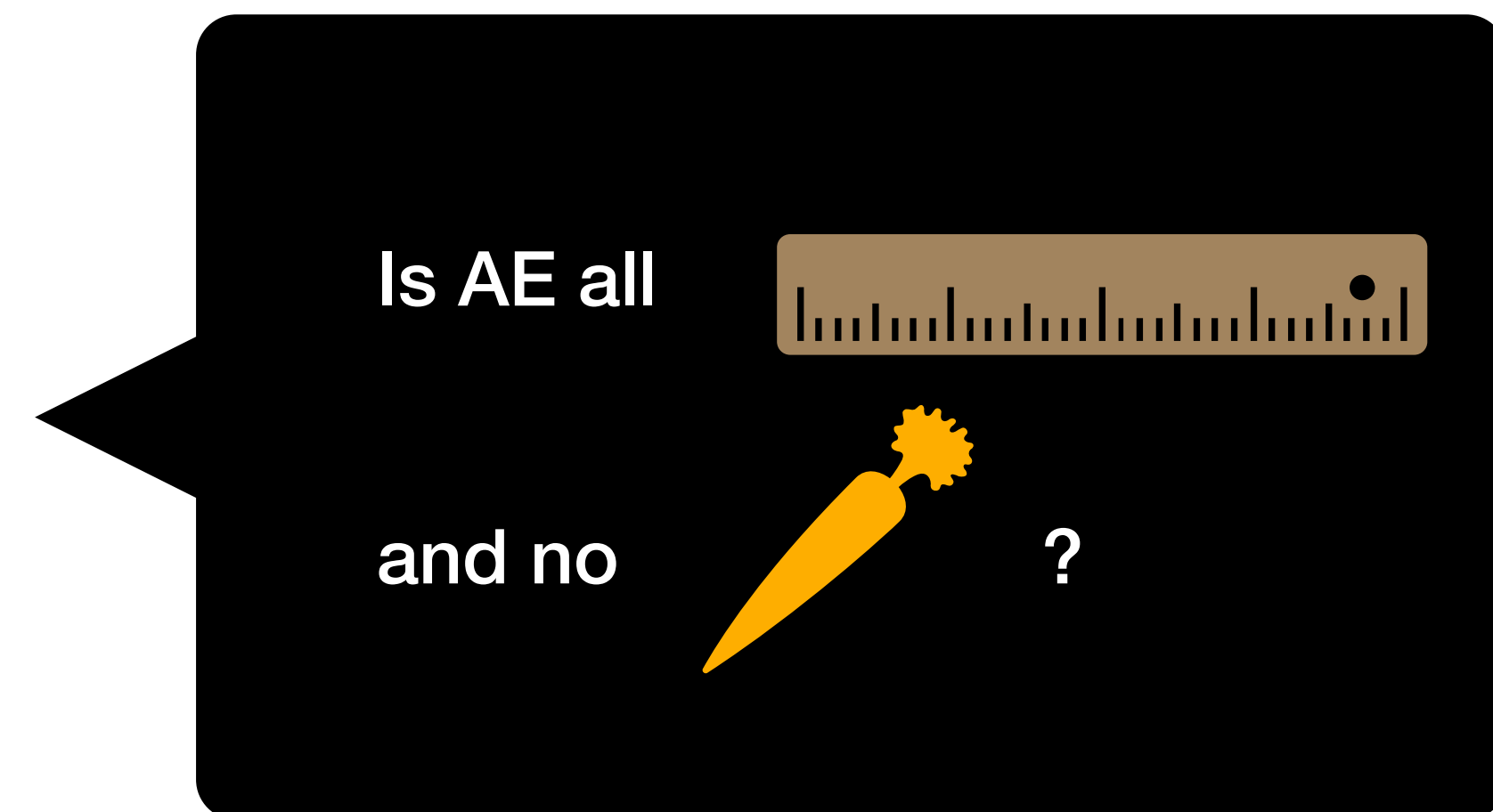
Bleh

# Only ever submitted once…

Oral History of Artifact Evaluation (student perspective)

# Why do we cite papers in the first place?

To please Reviewer #2.

Oral History of Artifact Evaluation (student perspective)
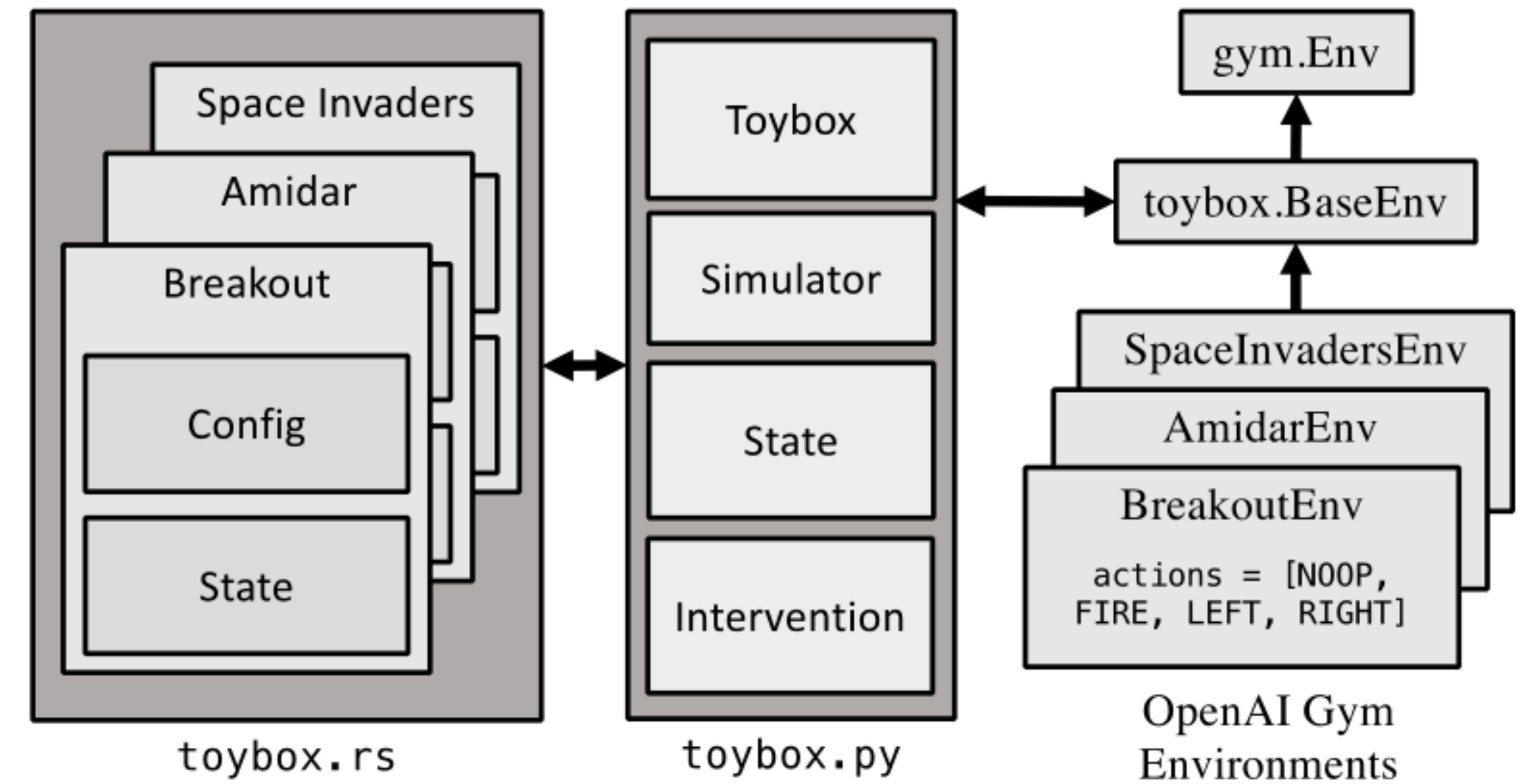
# Why do we cite papers in the first place?

- Findings
  - Don't *want* to have to start from scratch
- Contributions
  - New Software
  - New Datasets
  - New Methods
  - New Research Areas

Is AE all

and no                ?

# Software will be cited if it works*

## …regardless of AE results

- Incentive: Public artifact

  - Don't need artifact eval

  - Do we even want users?

    - Parable of SurveyMan

- Incentive: Good citizenship

  - Stand on the shoulder of giants!

  - Have you ever used someone else's artifact? (Not repo)



toybox.rs     toybox.py     OpenAI Gym Environments

**Toybox**



The Machine Learning Toybox for testing of Atari Reinforcement Learning Agents.

View My GitHub Profile

**toybox.rs**

Welcome to toybox.rs! This is the main organization and point of entry for using the Toybox platform for testing and experimentating with autonomous agents.

- Main repository with tests/experimentation support provided by a customized openai/baselines: toybox-rs/Toybox
- Core repository with implementations of the games: toybox-rs/toybox-rs. Releases available on PyPI: `pypi package 0.5.0`

**What is Toybox?**

Toybox is a set of *highly intervenable* environments for testing autonomous agents. While our efforts have focused on the efficient testing of deep RL agents, this work can be used in a variety of contexts that involve white-box testing of black-box agents.

If you use this code, or otherwise are inspired by our white-box testing approach, please cite our NeurIPS workshop paper:

Oral History of Artifact Evaluation (student perspective)

# Not able to convince collaborators to submit

# What about student evaluators?

# Student perspective
## (Students: feel free to share your thoughts)

- I liked serving on AECs

  - I learned new technologies

  - Reading others' code makes your code better

  - Scalable training in methods

- Other incentives:

  - Be on a PC (now students officially on PCs)

  - Early on: part of something important

- Problem: evaluation is a lot of work

# How to find more appealing carrots?

**What do student stakeholders want out of the process?**

# Artifact Evaluators as contributors

**Proposition 1**

# Submit… something else

**Proposition 2**

# Outline

Oral History of Artifact Evaluation (student perspective)

Evaluators produce replicates

Language design for reproducibility

# Where's the carrot?

**Answer: In empirical evaluation.**

# Focus efforts on replication
## Eval as improvement to the science

The value of <u>replicates</u>…

DISCLAIMER: Work in Progress

Evaluators produce replicates

# Focus efforts on replication

## Eval as improvement to the science

The value of <u>replicates</u>…

Specify as a causal graphical model

**Two values; assign equal weight**

Belief ⊙

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

Evaluators produce replicates

# Focus efforts on replication
## Eval as improvement to the science

The value of <u>replicates</u>…

Specify as a causal graphical model

Belief     **O**

**Y**

**Performance**

**Over what???**

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com
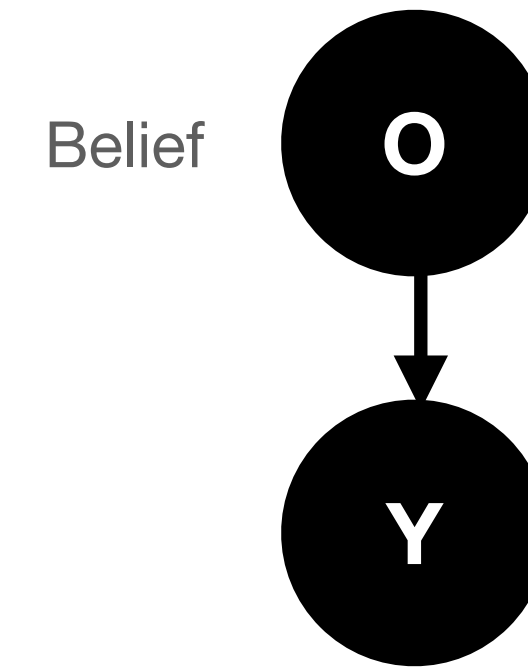
Evaluators produce replicates

# Focus efforts on replication

## Eval as improvement to the science

The value of <u>replicates</u>…

Specify as a causal graphical model

Belief

**O**

**Y**

**Performance**

**Over what???**

**Population of all possible programs on <u>my</u> machine**

**Population of all possible programs <u>on all suitable</u> machines**

**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

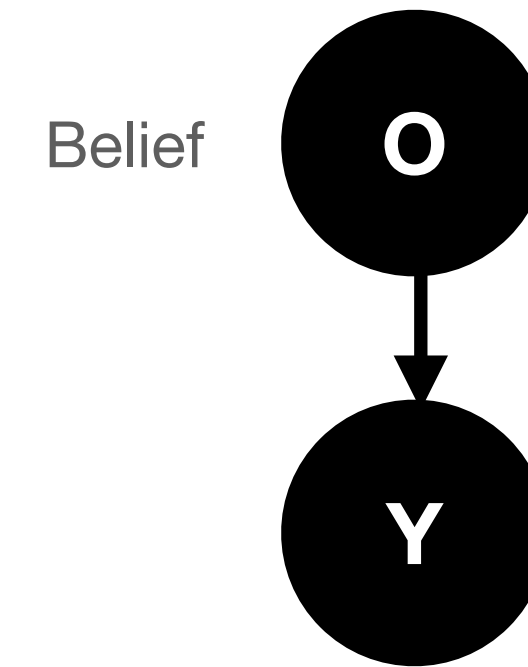IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

Evaluators produce replicates

# Focus efforts on replication
## Eval as improvement to the science

The value of <u>replicates</u>…

    Specify as a causal graphical model

Belief   **O**

**Y**

**Population of all possible programs on <u>my</u> machine**

**Population of all possible programs <u>on all suitable</u> machines**

**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz   Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

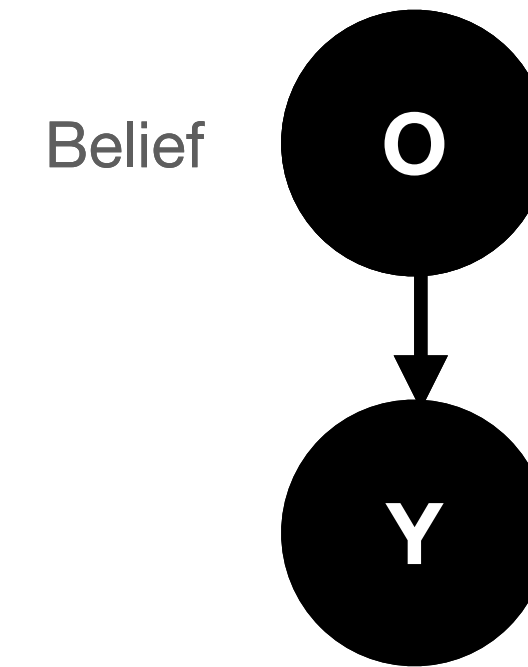IBM Research
Hawthorne, NY, USA
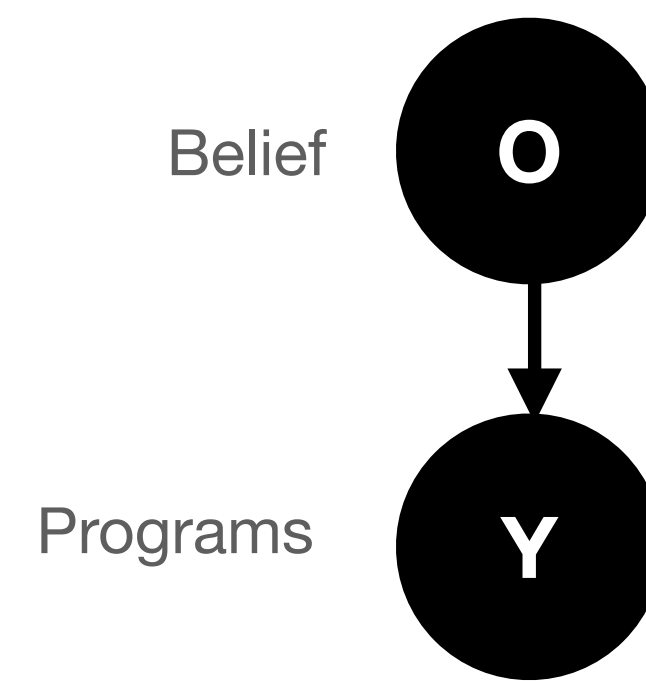pfs@us.ibm.com

Evaluators produce replicates

# Focus efforts on replication
## Eval as improvement to the science

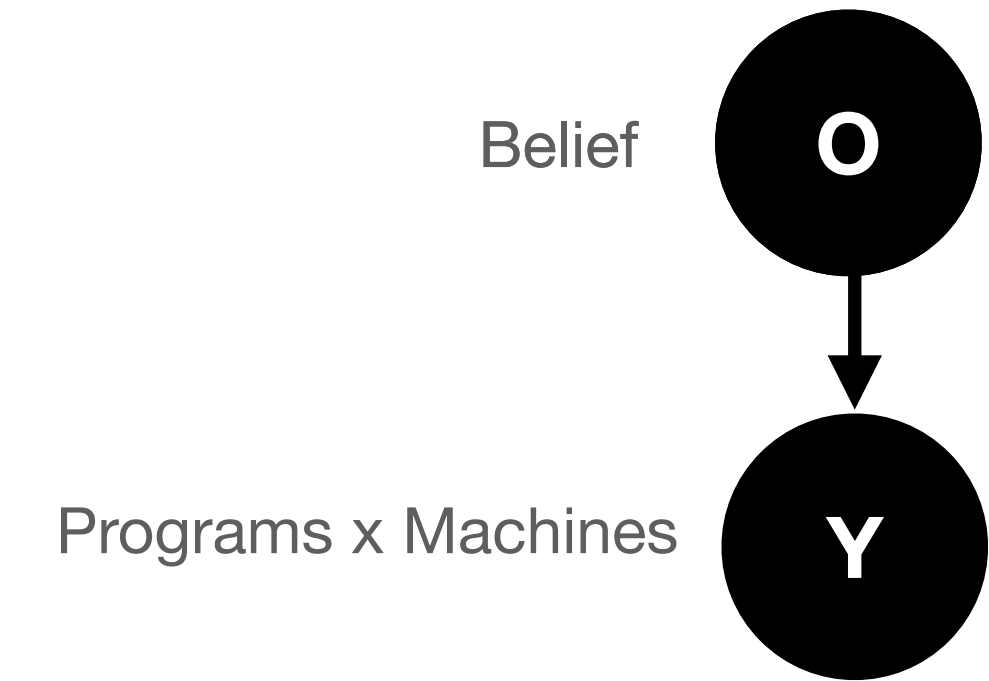The value of <u>replicates</u>…

Specify as a causal graphical model

Belief **O**

Programs **Y**

Belief **O**

Programs x Machines **Y**

**Population of all possible programs on <u>my</u> machine**

**Population of all possible programs <u>on all suitable</u> machines**

**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

Evaluators produce replicates
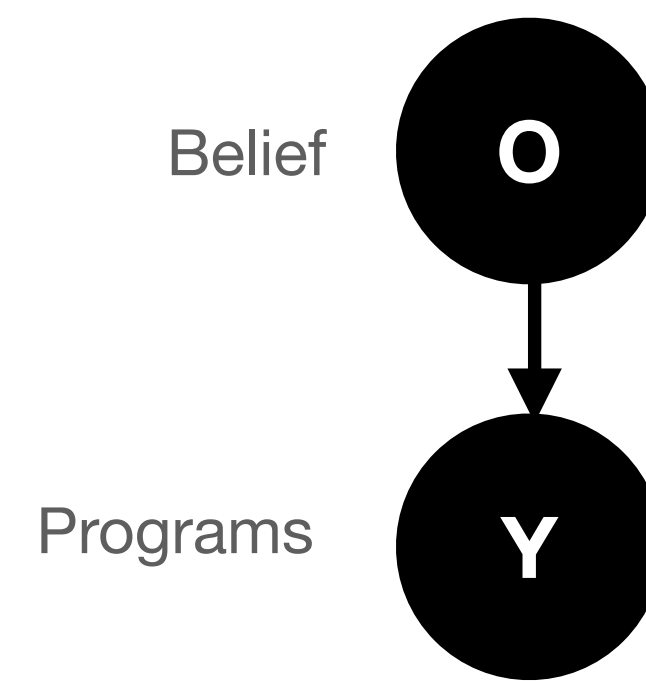
# Focus efforts on replication
## Eval as improvement to the science

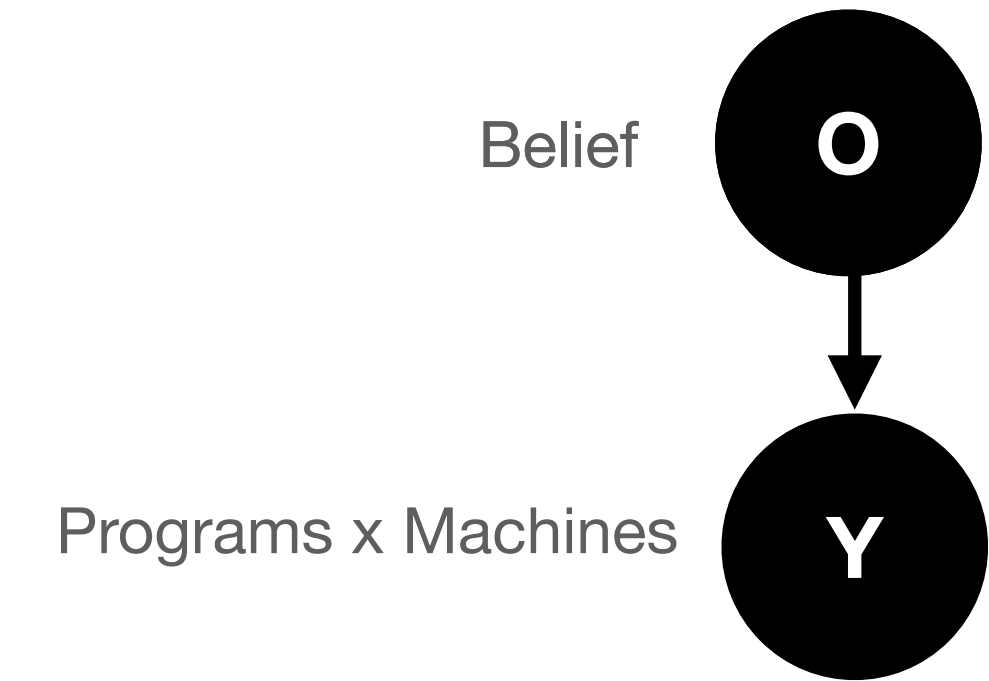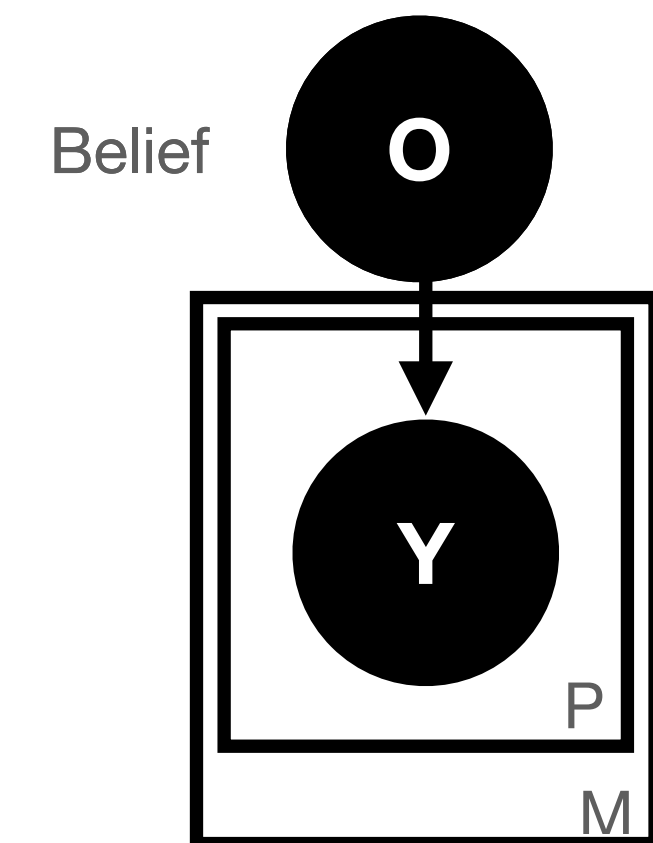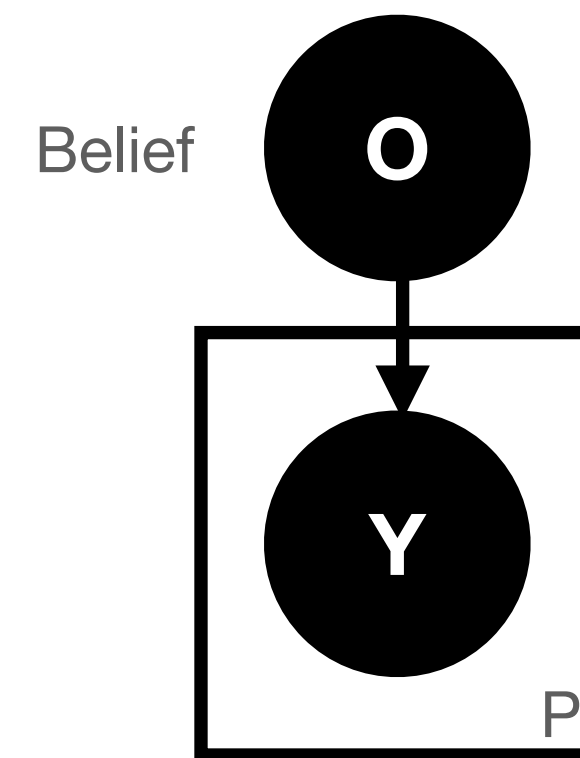The value of <u>replicates</u>…

Specify as a causal graphical model



Belief ● O

Programs ● Y

Belief ● O

Programs x Machines ● Y

**Population of all possible programs on <u>my</u> machine**

**Population of all possible programs <u>on all suitable machines</u>**

**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz  Amer Diwan
Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth
Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney
IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

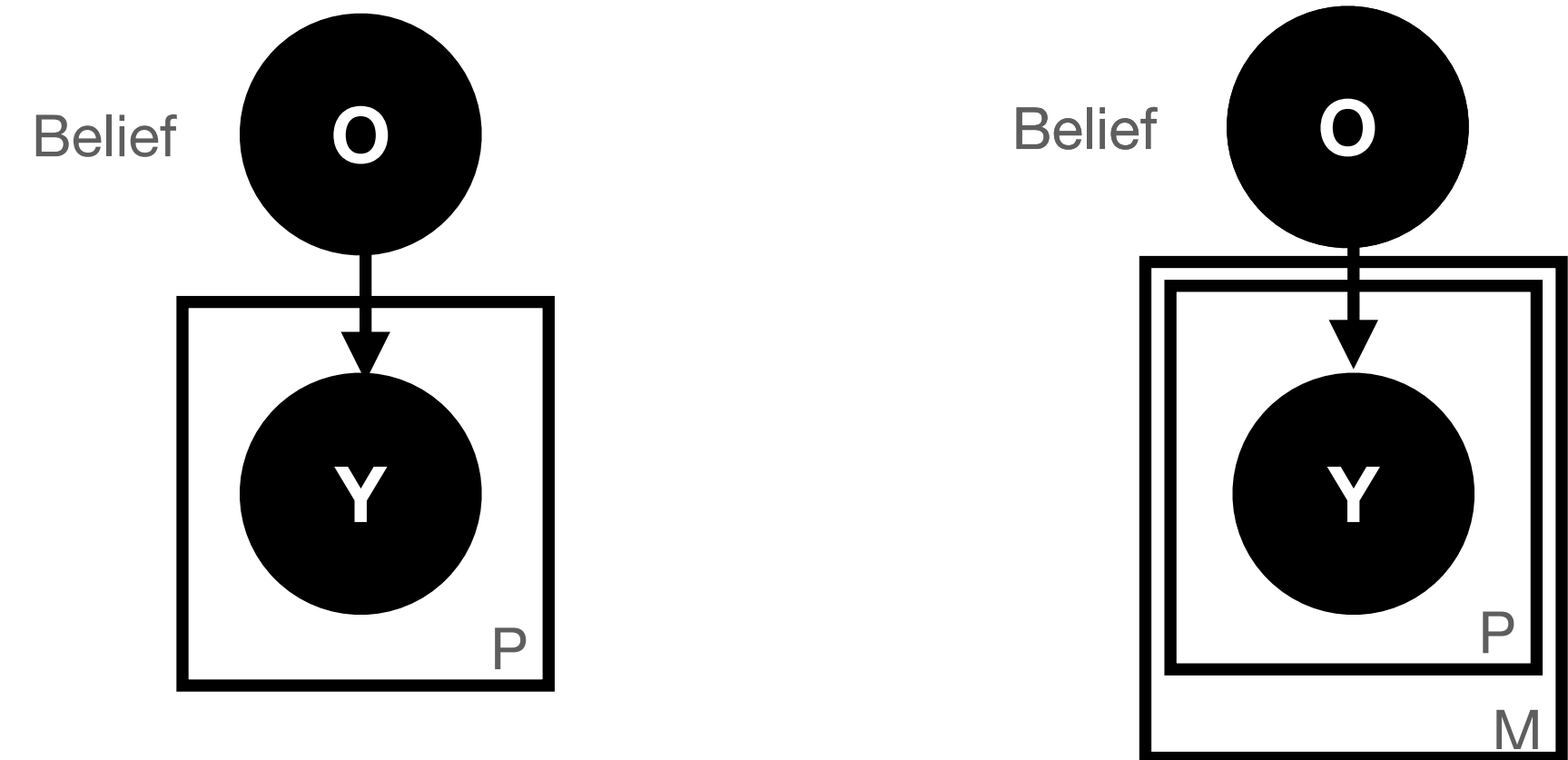Evaluators produce replicates

# Focus efforts on replication
## Eval as improvement to the science

The value of <u>replicates</u>…

Specify as a causal graphical model

Belief **O**

**Y** P

Belief **O**

**Y** P

M

**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

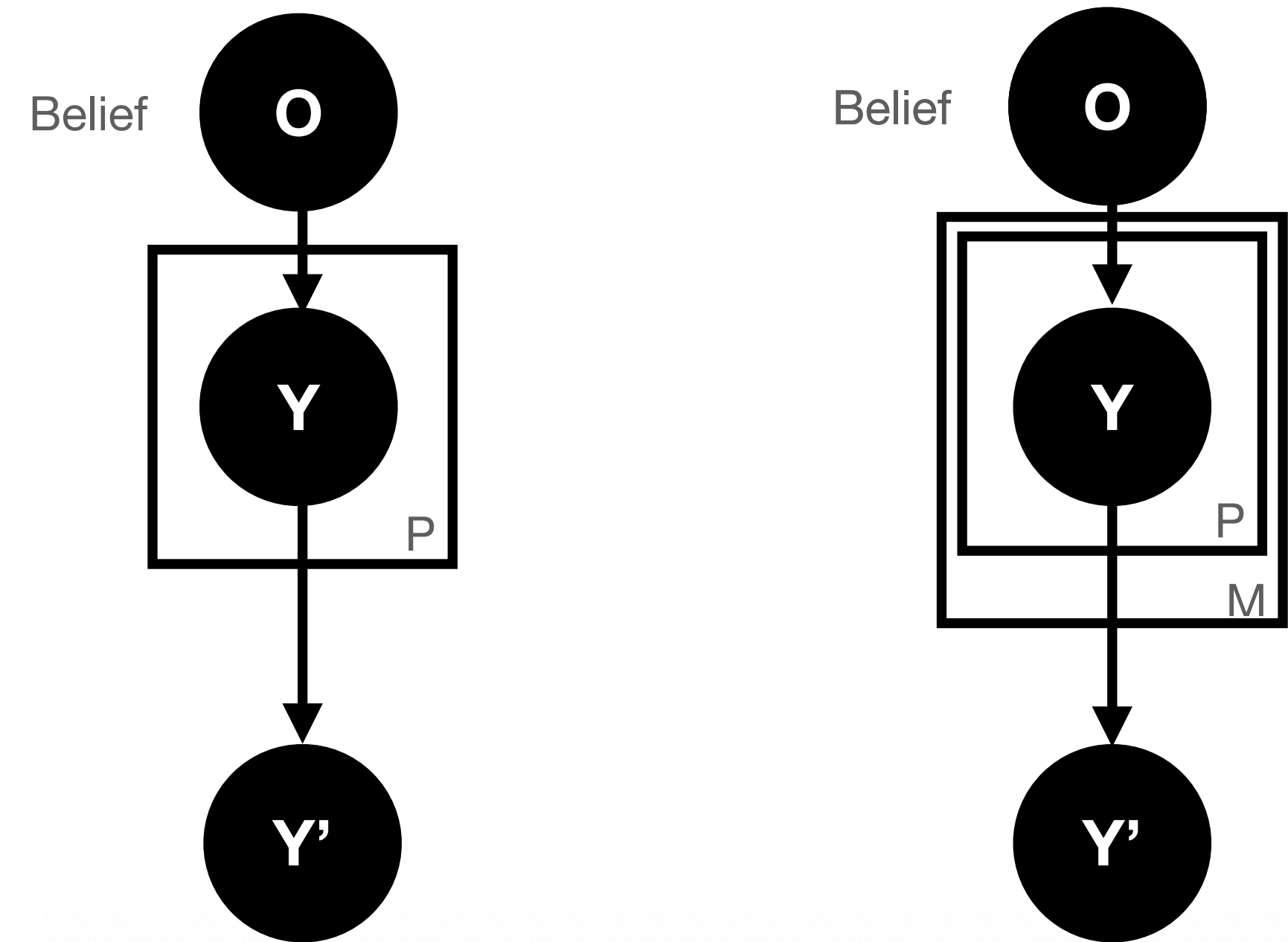Evaluators produce replicates
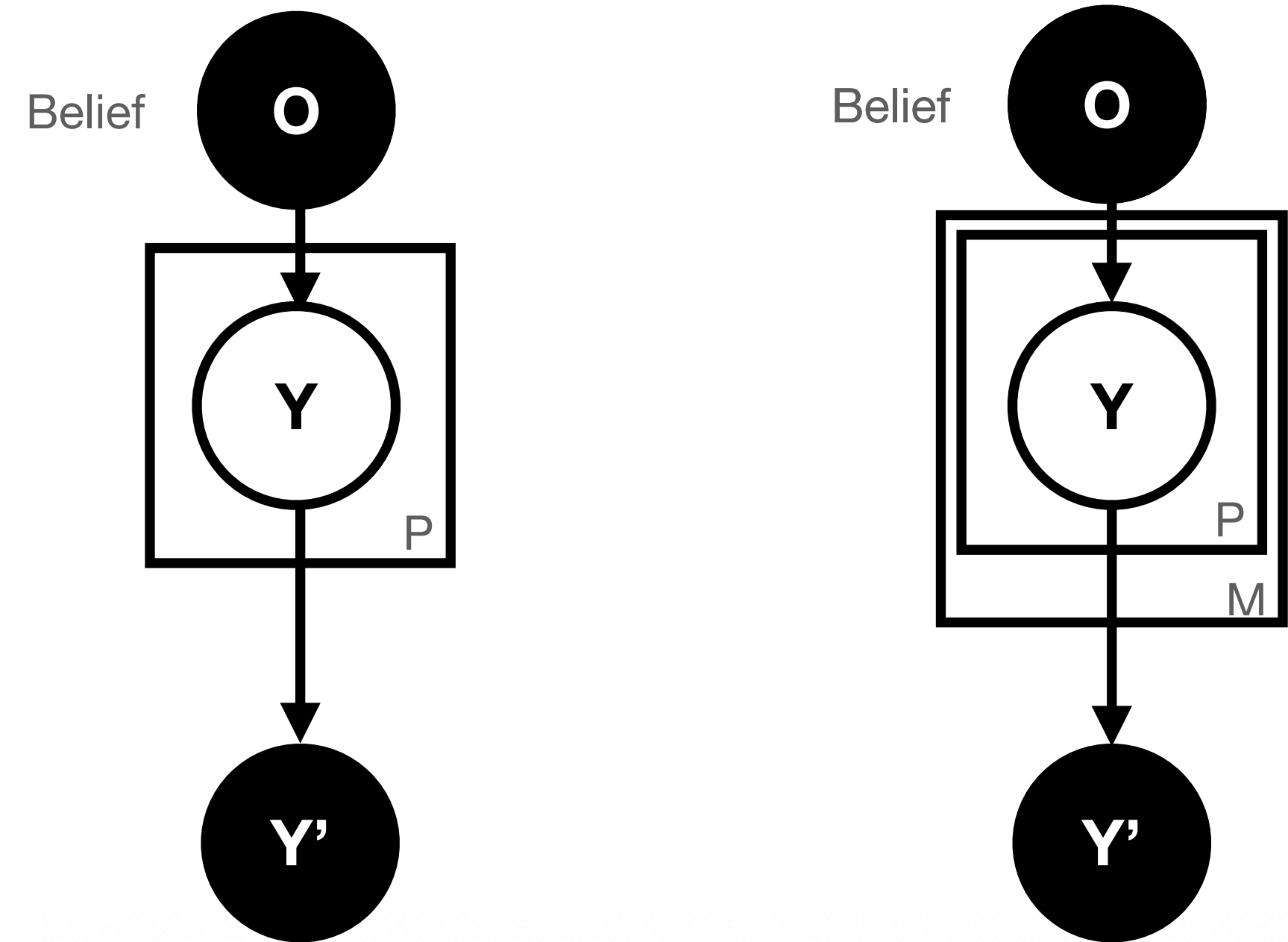
# Focus efforts on replication

## Eval as improvement to the science

<span style="color: #990066; font-weight: bold;">Sample is fixed! (Benchmark)</span>

The value of <u>replicates</u>…

Specify as a causal graphical model



Belief — O → Y (P)

Belief — O → Y (P) (M)



**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

IBM Research
Hawthorne, NY, USA
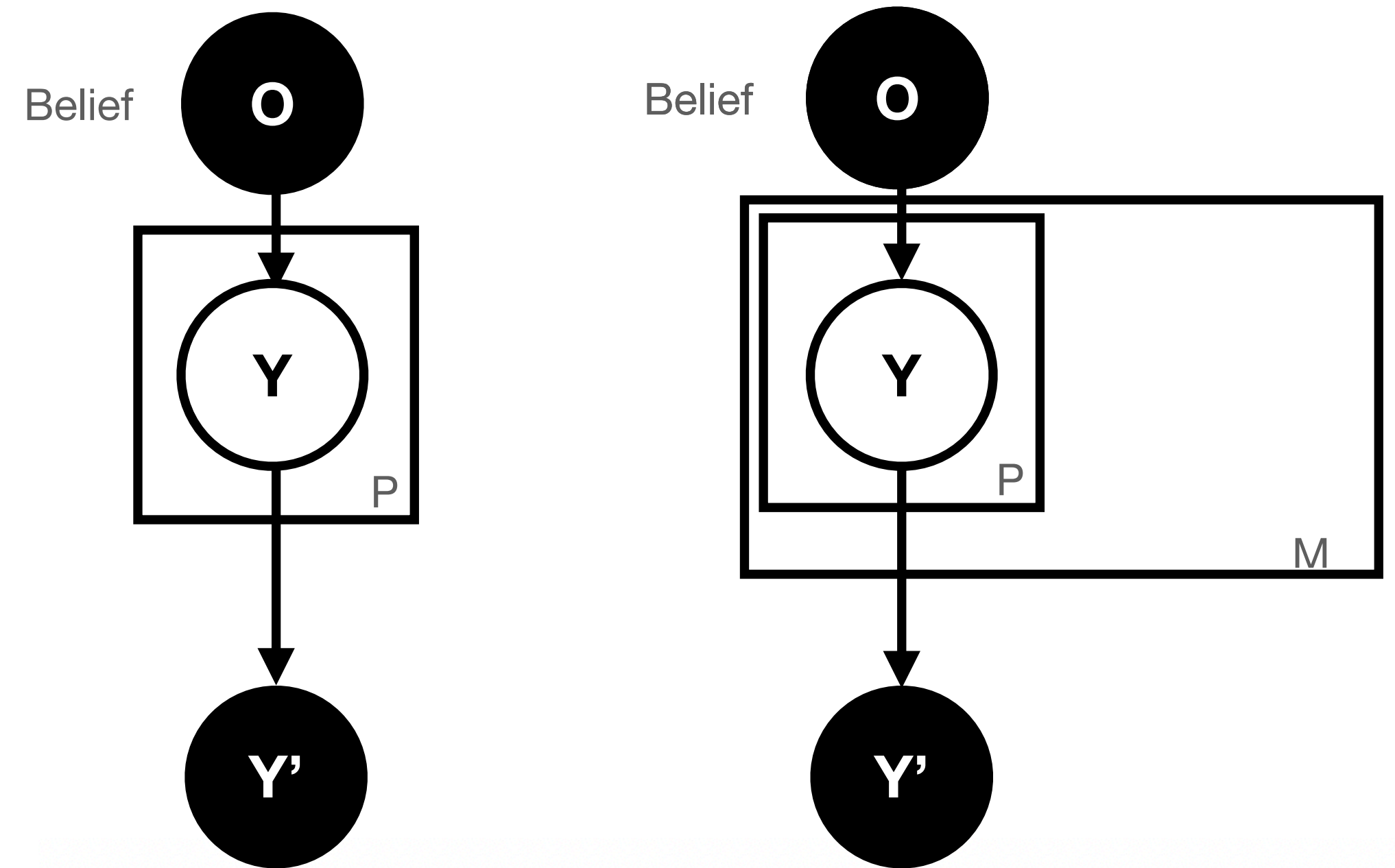pfs@us.ibm.com

Evaluators produce replicates

# Focus efforts on replication

## Eval as improvement to the science

**Sample is fixed! (Benchmark)**

The value of <u>replicates</u>…

Specify as a causal graphical model



**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

Evaluators produce replicates
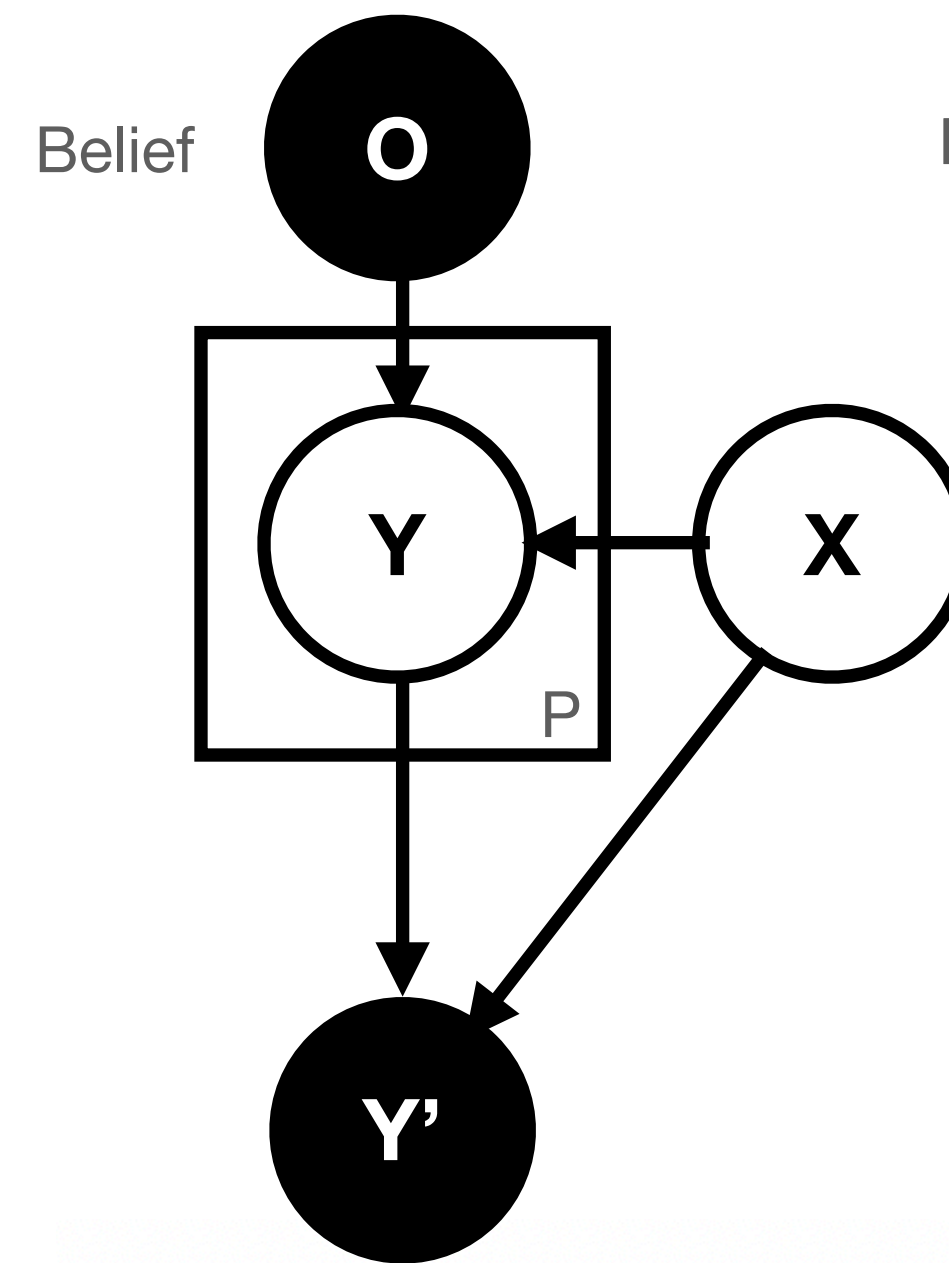
# Focus efforts on replication

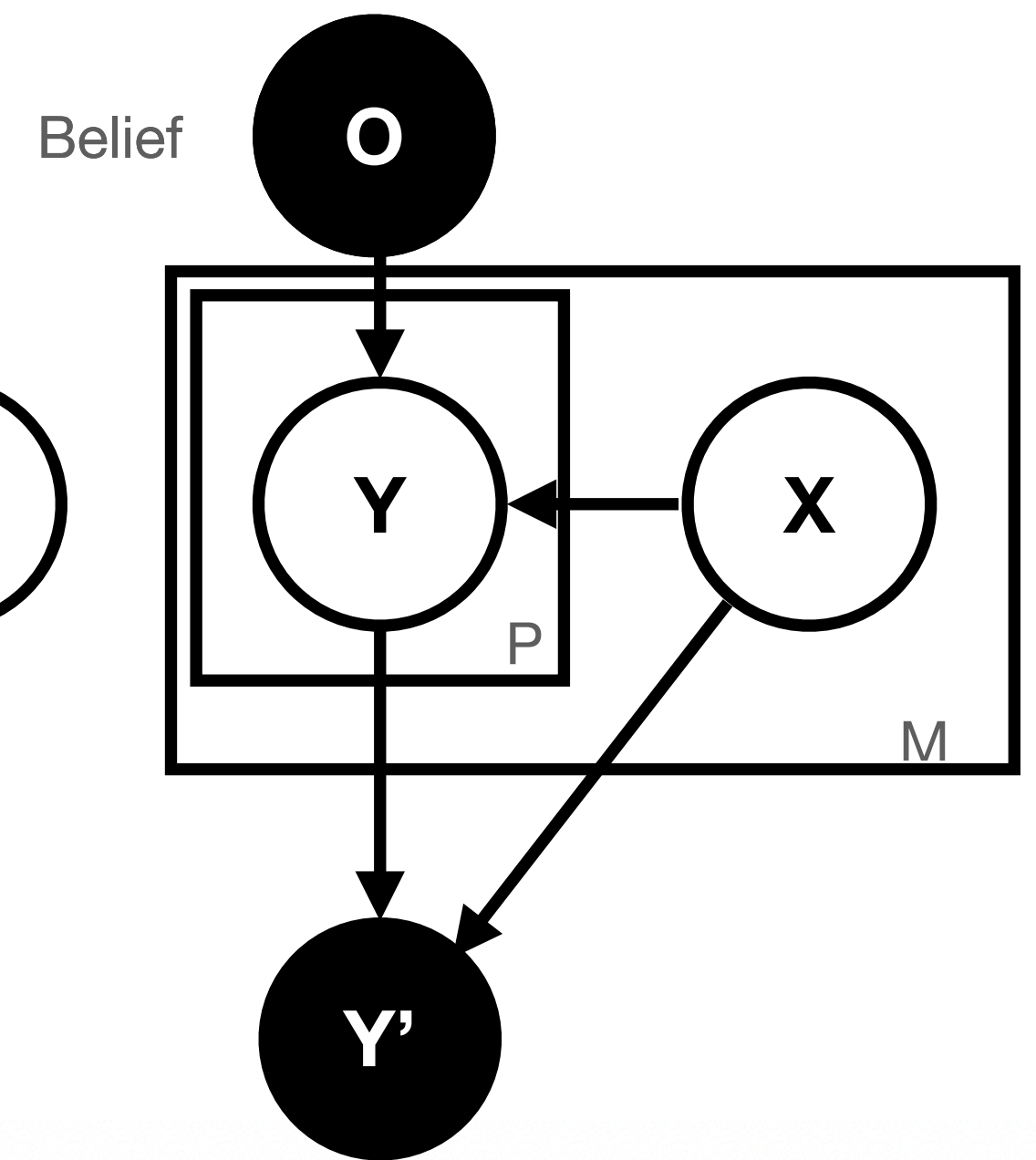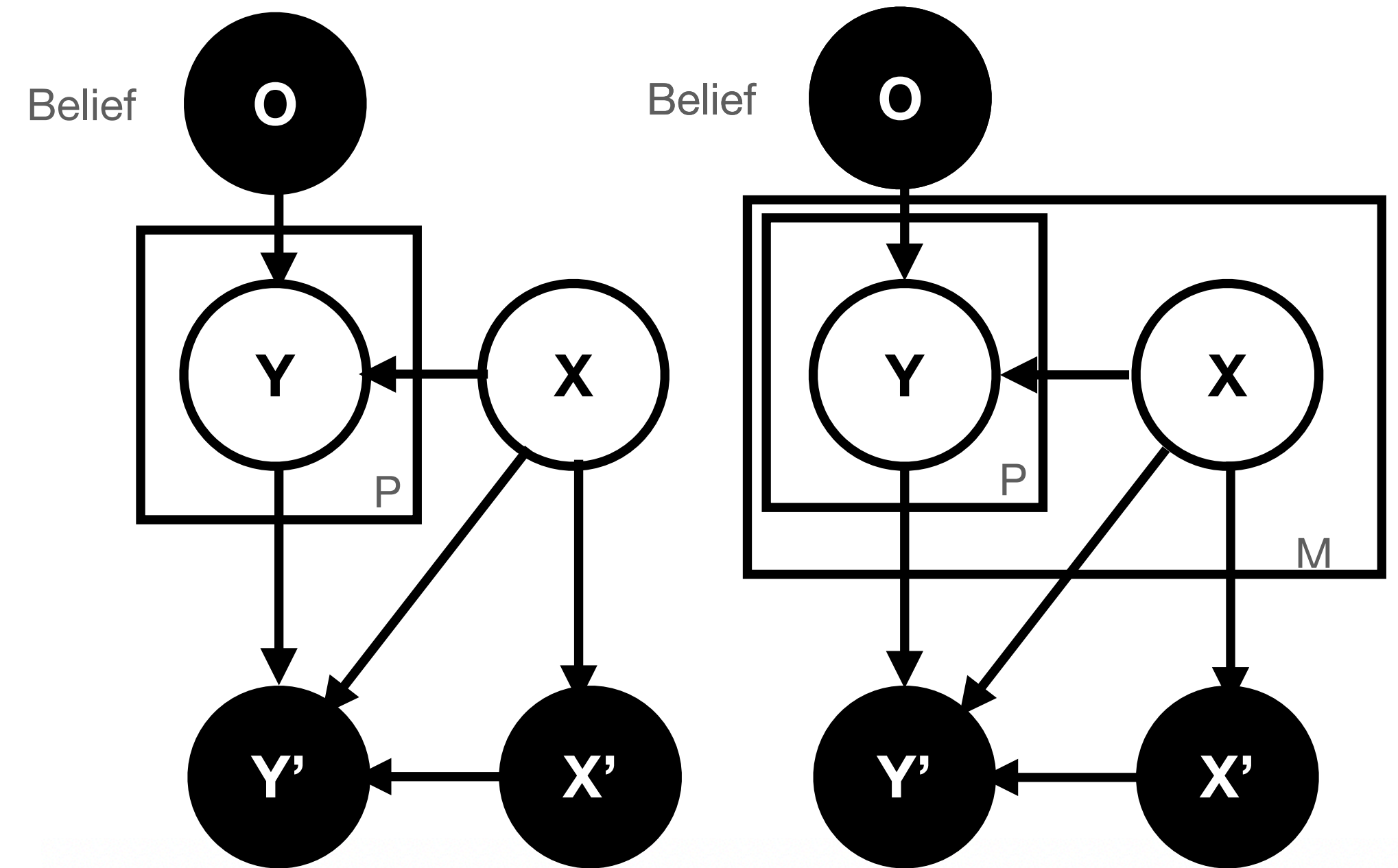## Eval as improvement to the science

The value of <u>replicates</u>…

Specify as a causal graphical model



**Sample is fixed! (Benchmark)**

Belief

O

Y

P

Y'

Belief

O

Y

P

M

Y'

**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

Evaluators produce replicates

# Focus efforts on replication
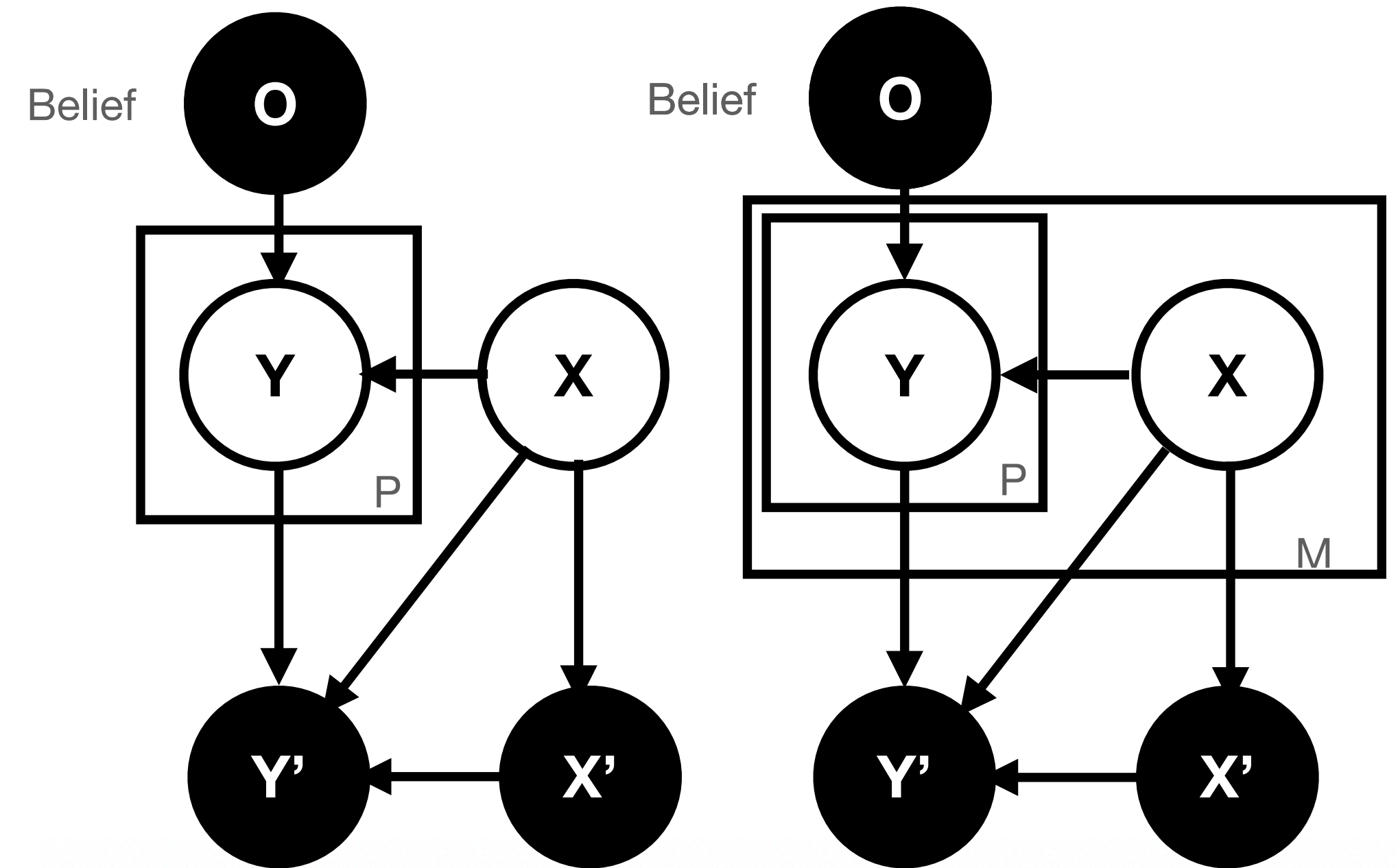## Eval as improvement to the science

The value of <u>replicates</u>…

Specify as a causal graphical model



Belief

O

Y

P

Y'

Belief

O

Y

P

M

Y'

**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz   Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney
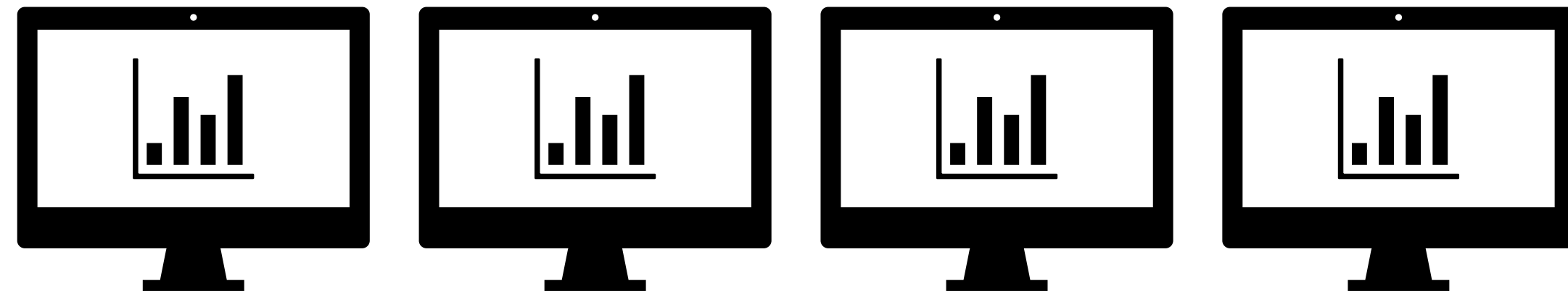
IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

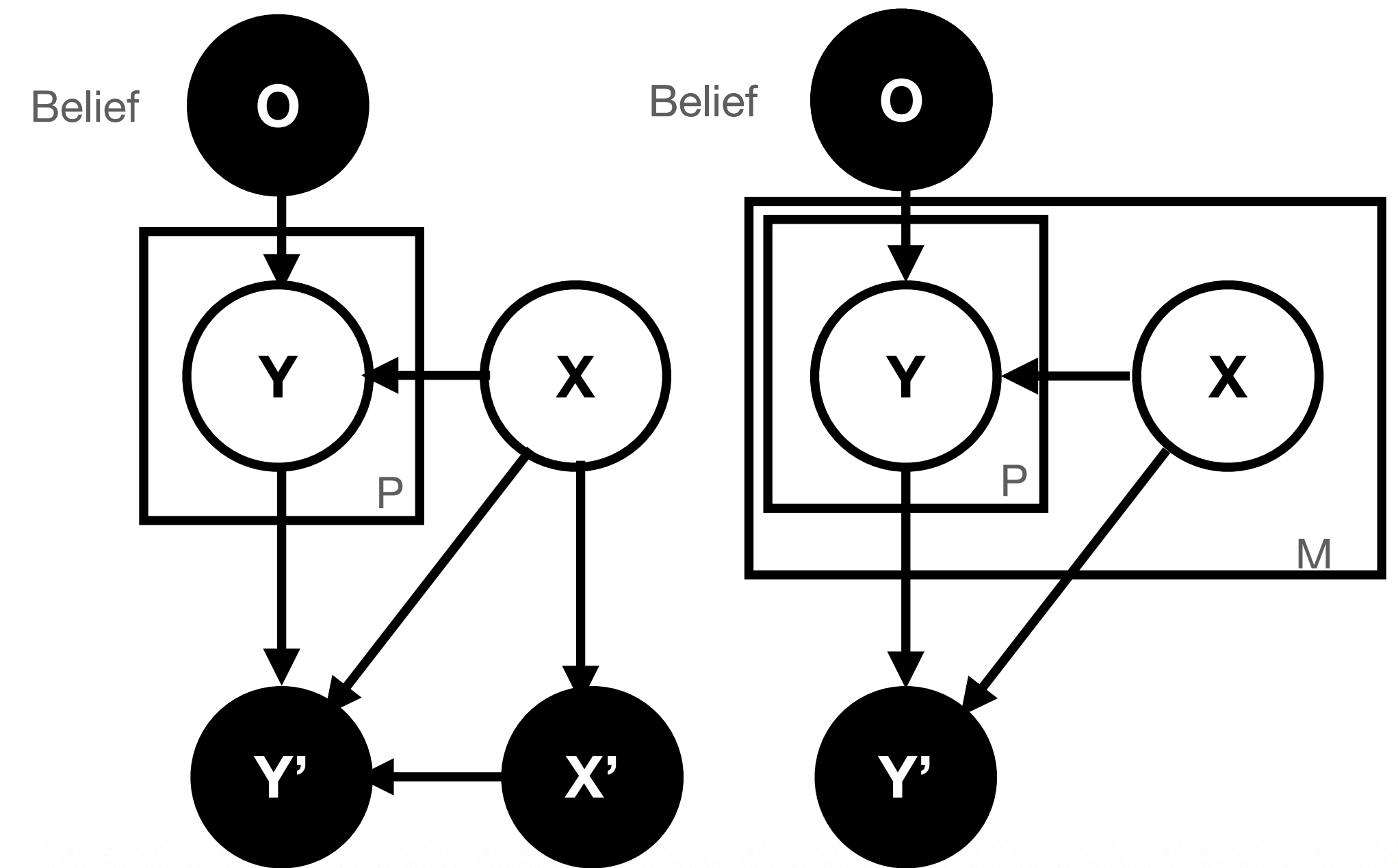Evaluators produce replicates

# Focus efforts on replication

## Eval as improvement to the science

The value of <u>replicates</u>…

Specify as a causal graphical model

**X has epistemic uncertainty**

**X has epistemic *and* aleatory uncertainty**



Belief

O

Y ← X

P

Y'



Belief

O

Y ← X

P

M

Y'

**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

Evaluators produce replicates

# Focus efforts on replication

## Eval as improvement to the science

The value of <u>replicates</u>…

Specify as a causal graphical model



**X has epistemic uncertainty**

**X has epistemic *and* aleatory uncertainty**

Belief

Belief

**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

Evaluators produce replicates

# Focus efforts on replication

## Eval as improvement to the science

The value of <u>replicates</u>…

Specify as a causal graphical model



**X has epistemic uncertainty**

**X has epistemic *and* aleatory uncertainty**

**Producing Wrong Data Without Doing Anything Obviously Wrong!**

Todd Mytkowicz  Amer Diwan
Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth
Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney
IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

Evaluators produce replicates

# Focus efforts on replication
## Eval as improvement to the science

The value of <u>replicates</u>…

Specify as a causal graphical model



**Elements of M chosen arbitrarily (good enough)**

Todd Mytkowicz  Amer Diwan

Department of Computer Science
University of Colorado
Boulder, CO, USA
{mytkowit,diwan}@colorado.edu

Matthias Hauswirth

Faculty of Informatics
University of Lugano
Lugano, CH
Matthias.Hauswirth@unisi.ch

Peter F. Sweeney

IBM Research
Hawthorne, NY, USA
pfs@us.ibm.com

Evaluators produce replicates

# Outline

Oral History of Artifact Evaluation (student perspective)

Evaluators produce replicates

Language design for reproducibility

# Why this is interesting
## CGMs are hard to get right

- Abuse of plate notation?

    - Y' is *not* randomly sampled

    - Should X' be a random variable?

    - Should we have a separate value for P?

# Better: state assumptions in a language
## Specifically, a hypothesis language

```
1    O : { "O2", "O3" }
2    Y : nat
3    (progid) Y <- O
4    sharp (progid) assert (Y > 0)
5    Y_A = Y | O = "O2"
6    Y_B = Y | O = "O3"
7    (progid) assert (Y_A > Y_B)
```



```
1    for trialid in repeat(progid, machineid) under complete:
2        measure(Y)
```

Language design for reproducibility

# Better: state assumptions in a language
## Specifically, a hypothesis language

```
1    O : { "O2", "O3" }
2    Y : nat
3    E : nat
4    P : { "Pentium4", "Core2", "m5O3CPU" }
5    C : { "gcc", "intel" }
6    L : nat
7    (progid) Y <- O, L, E, C, P
8    Y_A = Y | O = "O2", L
9    Y_B = Y | O = "O3", L
10   (progid) assert (Y_A > Y_B)
11   (progid) Y_B >--> E
```

# HyPL

$$op ::= \ = \ | \ > \ | \ <$$

$$coef ::= \ ? \ | \ n$$

$$sup ::= \text{nat} \ | \ \text{bool} \ | \ \{\mathbf{str}_1, \mathbf{str}_2, \ldots, \mathbf{str}_n\} \ | \ \text{real}$$

$$decl ::= X : \langle sup \rangle \ | \ X : \langle sup \rangle \text{ of } (\mathbf{unitid}_i) \ | \ Y' = Y|(X_1 \ \langle op \rangle \ v_1, \ldots, X_n \ \langle op \rangle \ v_n)$$

$$hfn ::= \langle coef \rangle \ | \ \langle coef \rangle \ X \ | \ \langle coef \rangle \ X_1 \ X_2 \ | \ \langle coef \rangle \ \exp(\langle hfn \rangle) \ | \ \langle hfn \rangle + \langle hfn \rangle$$

$$htype ::= \text{sharp } (\mathbf{unitid}_i) \ | \ (\mathbf{unitid}_i) \ | \ \text{belief}$$

$$bexp ::= \top \ | \ \bot \ | \ X \ | \ ! \langle bexp \rangle \ | \ \langle bexp \rangle \ \&\& \ \langle bexp \rangle \ | \ \langle bexp \rangle \ || \ \langle bexp \rangle \ | \ \langle hfn \rangle \ \langle op \rangle \ \langle hfn \rangle \ | \ \langle hyp \rangle$$

$$hyp ::= \underbrace{\langle htype \rangle \ Y := \langle hfn \rangle}_{SEM \ (strong \ causal)} \ | \ \underbrace{\langle htype \rangle \ Y <\!- X}_{weak \ causal} \ | \ \underbrace{\langle htype \rangle \ Y <\!- X|Z}_{weak \ causal \ conditional} \ | \ \underbrace{\langle htype \rangle \ Y >\!-\!> X}_{monotonic \ (associational)}$$

$$| \ \langle htype \rangle \ \text{assert} \ \langle bexp \rangle \ | \ \text{when } \langle bexp \rangle \text{ then } \langle hyp+ \rangle \text{ end}$$

$$stmt ::= \langle decl \rangle \ | \ \langle hyp \rangle$$

$$model ::= \langle stmt+ \rangle$$

Language design for reproducibility

# Why another PPL?

It's *not* all about the parameters

# Additional affordances via language-based approach

# Enables: Structured Search

## …or, search beyond keywords

```
1   O : { "O2", "O3" }
2   Y : nat
3   E : nat
4   P : { "Pentium4", "Core2", "m5O3CPU" }
5   C : { "gcc", "intel" }
6   L : nat
7   (progid) Y <- O, L, E, C, P
8   Y_A = Y | O = "O2", L
9   Y_B = Y | O = "O3", L
10  (progid) assert (Y_A > Y_B)
11  (progid) Y_B >--> E
```

# Enables: Continuous Auditing
## …or, regression testing for past studies

# Enables: Onboarding neophytes
## Make adhering to best practices easier!



SIGPLAN Empirical Evaluation Checklist

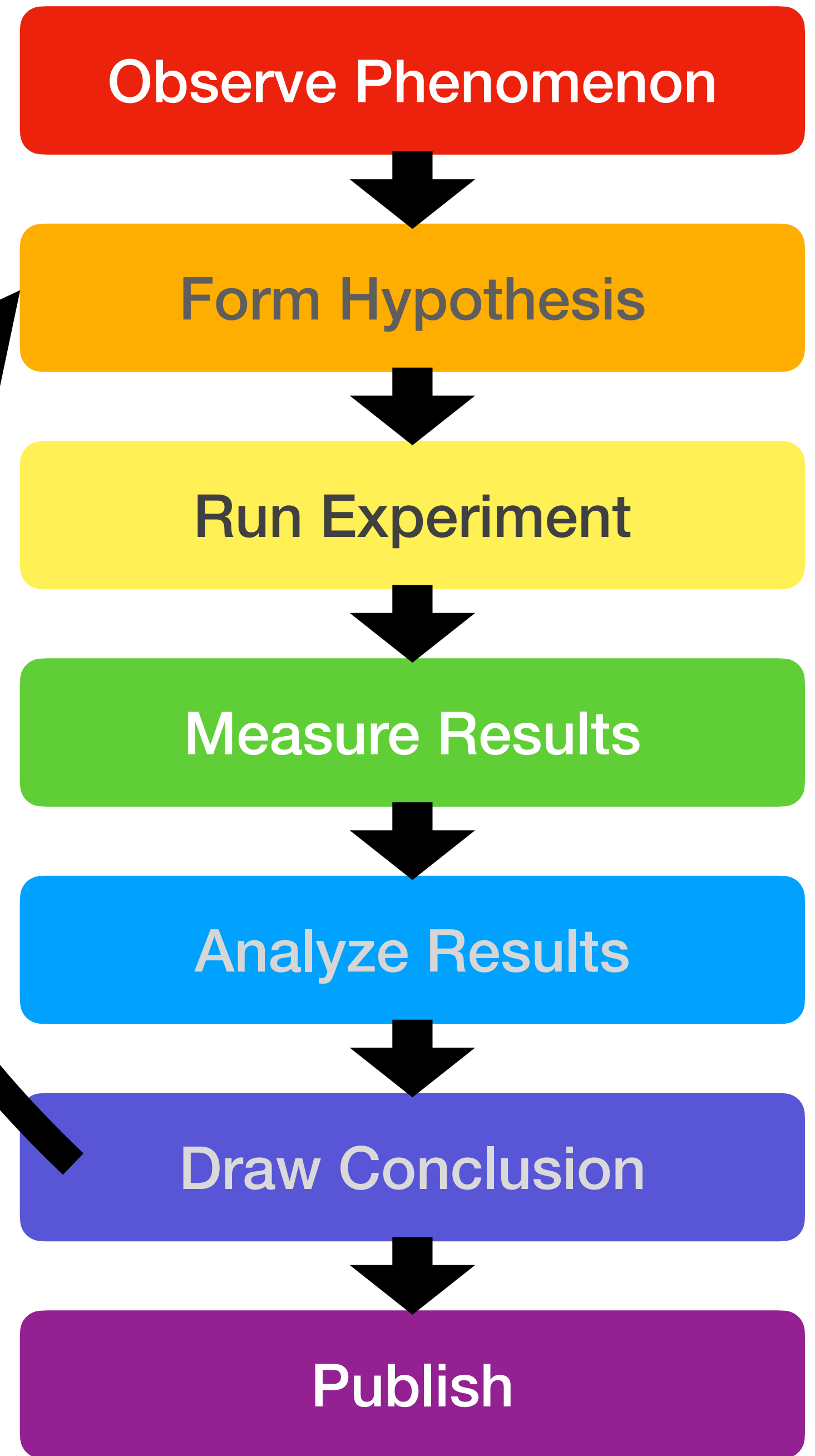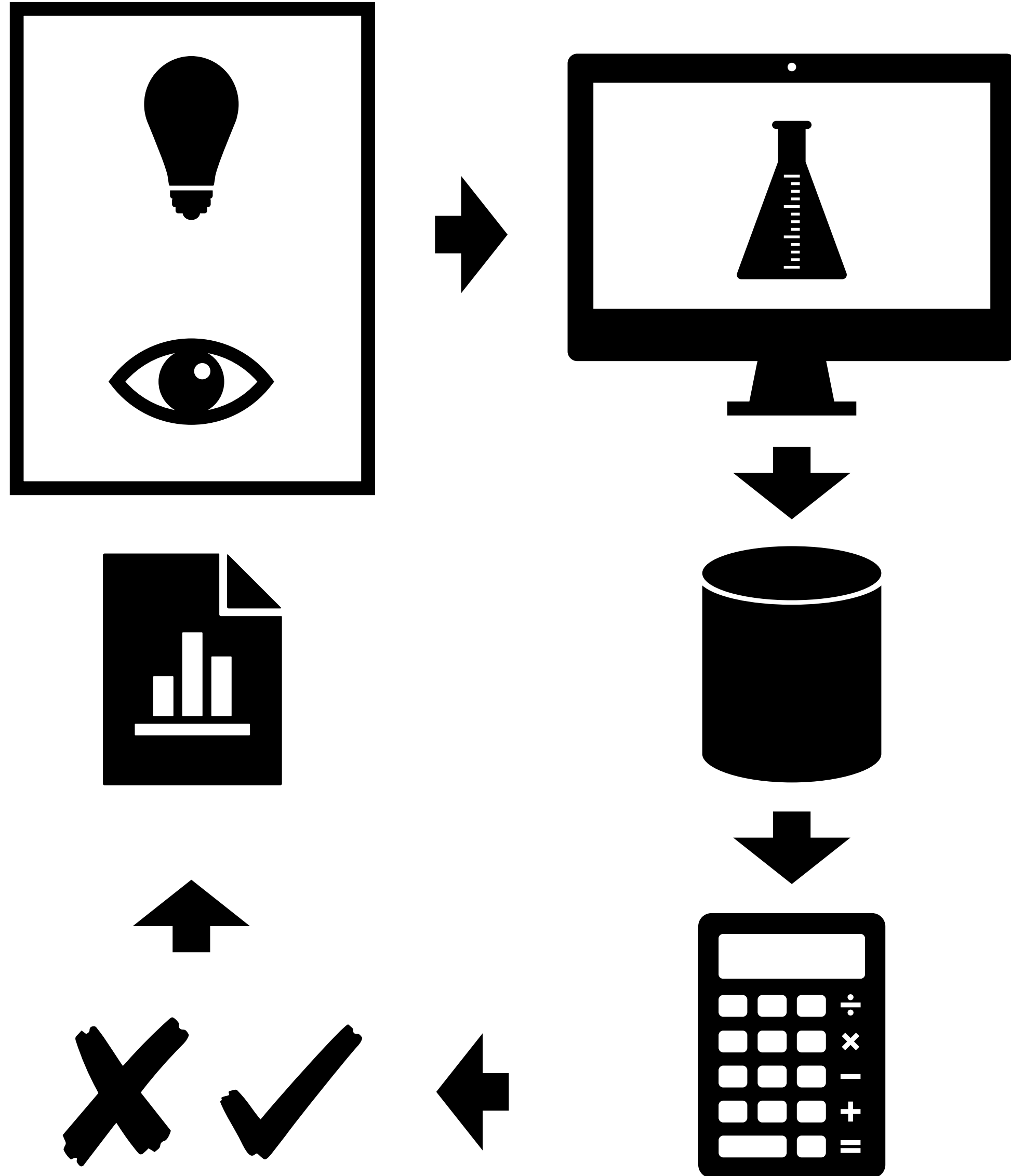Language design for reproducibility

# Challenges in application to cybersecurity

# Extreme values

Interested in maxima or the long tail?

Need different methods!

**Extreme values &
Non-scientific knowledge**



Observe Phenomenon

Form Hypothesis

Run Experiment

Measure Results

Analyze Results

Draw Conclusion

Publish

Not an end, but hopefully a 🥕