A CONVERSATION WITH

# NIHAT AY

@KyleCranmer

New York University

Department of Physics

Center for Data Science

CILVR Lab

# Yesterday's talk

θ Theory

**PREDICTION** →

← **INFERENCE**

x Data

# The Scientific Method as an Ongoing Process

**Make Observations**
What do I see in nature?
This can be from one's own experiences, thoughts, or reading.

**Think of Interesting Questions**
Why does that pattern occur?

**Develop General Theories**
General theories must be consistent with most or all available data and with other current theories.

**Formulate Hypotheses**
What are the general causes of the phenomenon I am wondering about?

**Refine, Alter, Expand, or Reject Hypotheses**

**Gather Data to Test Predictions**
Relevant data can come from the literature, new observations, or formal experiments. Thorough testing requires replication to verify results.

**Develop Testable Predictions**
If my hypotesis is correct, then I expect a, b, c,...
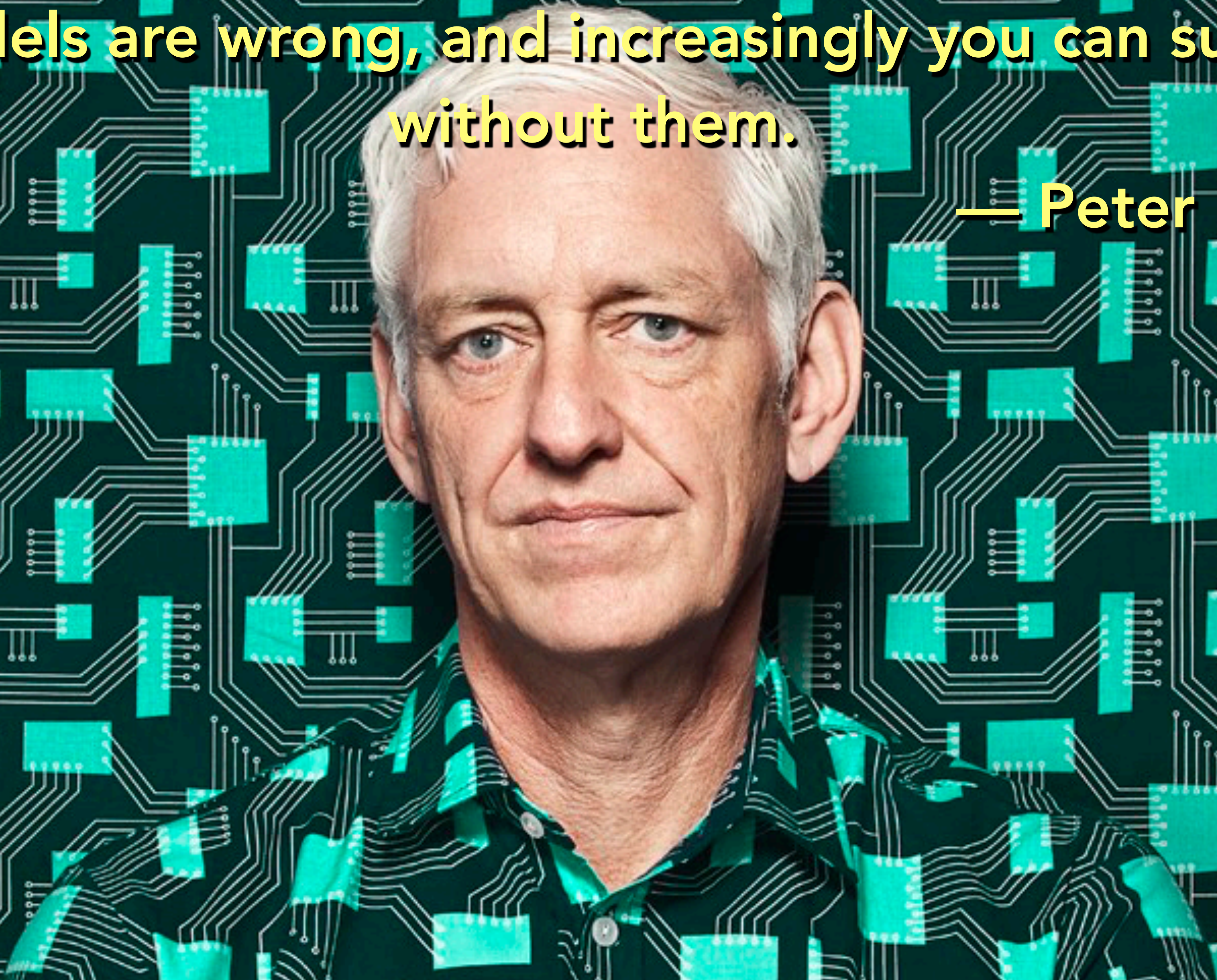
CHRIS ANDERSON  SCIENCE  06.23.08  12:00 PM

# THE END OF THEORY: THE DATA DELUGE MAKES THE SCIENTIFIC METHOD OBSOLETE

All models are wrong, and increasingly you can succeed without them.

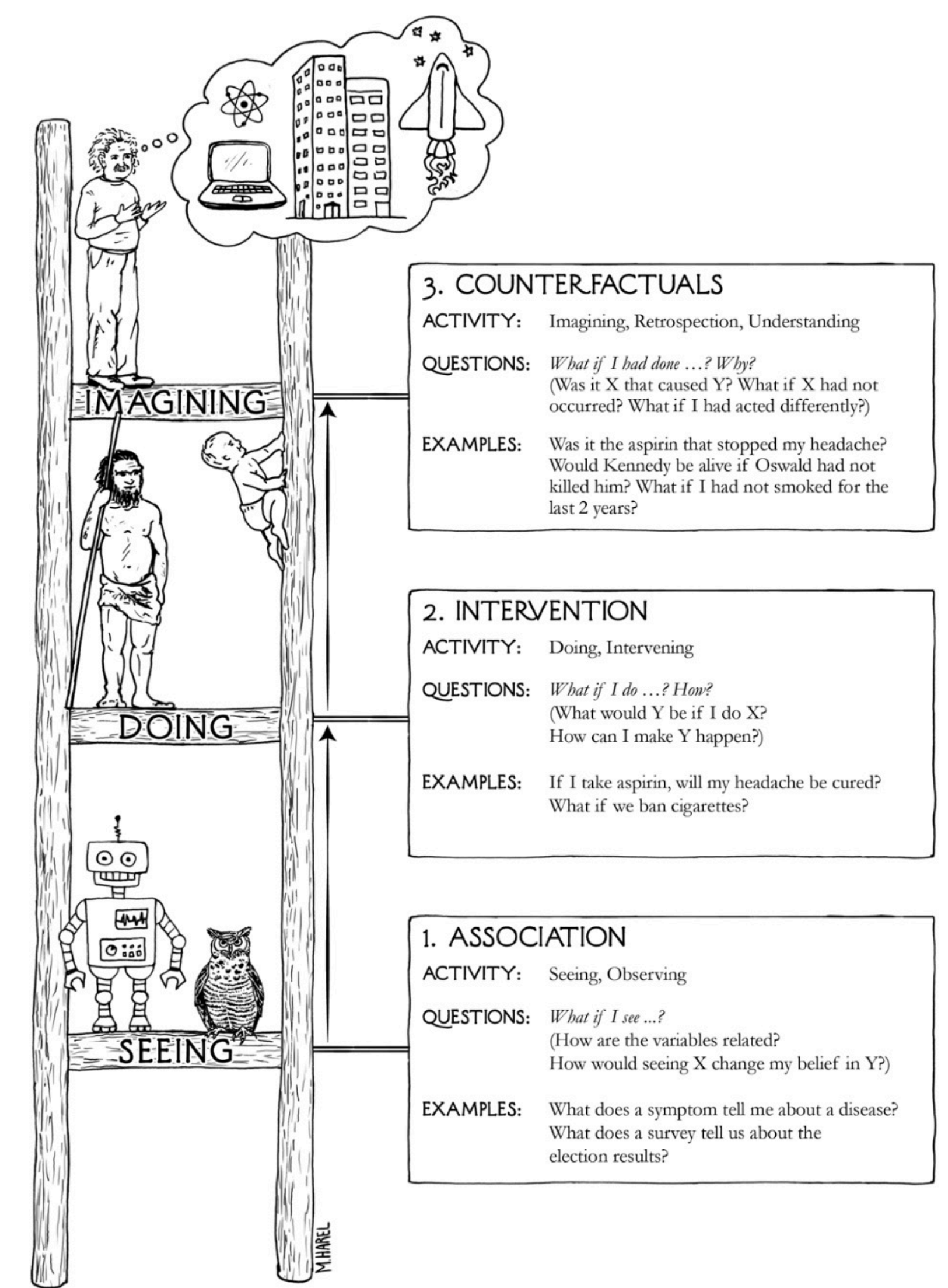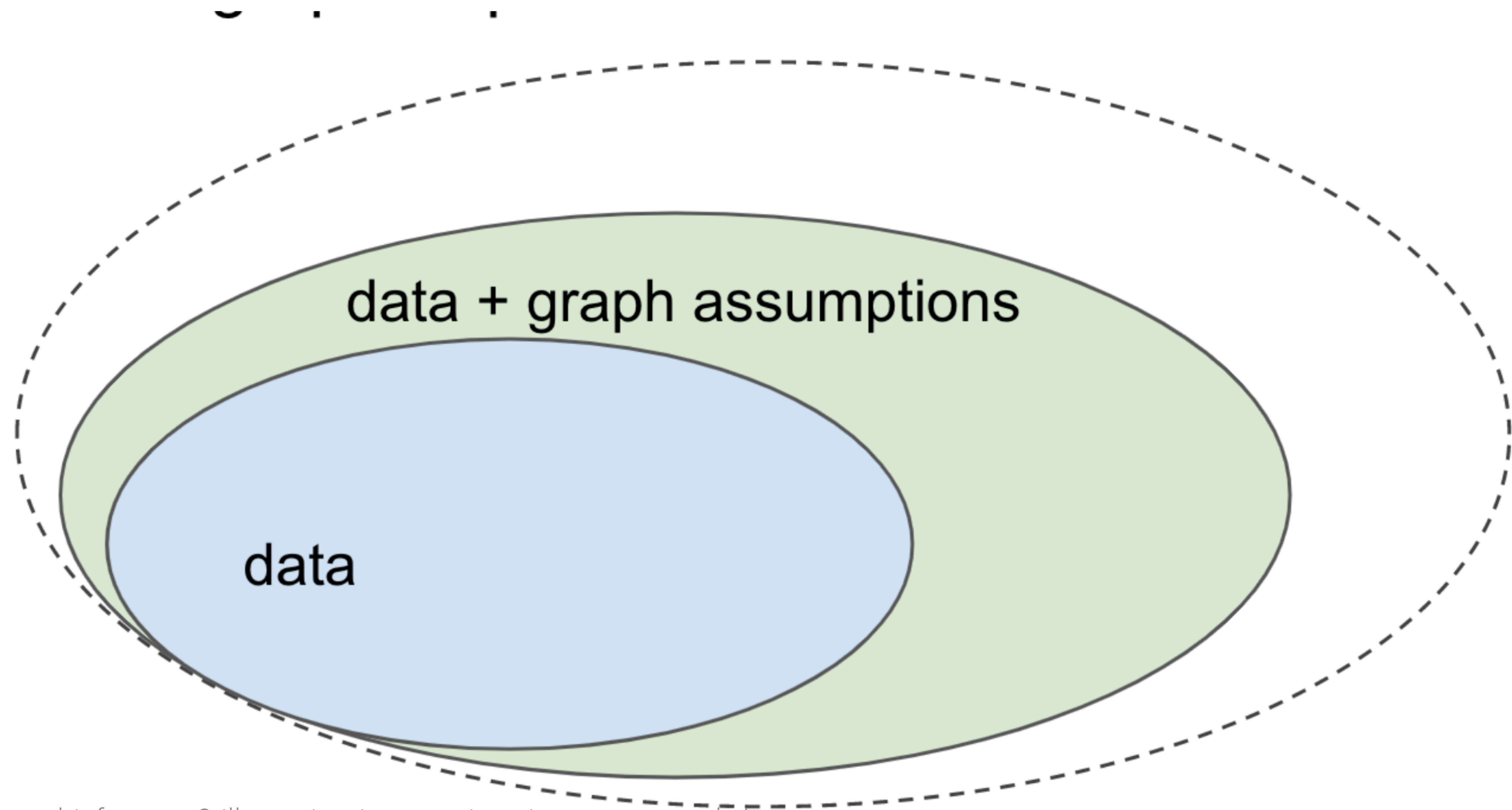— Peter Norvig

# The ladder of Causation



FIGURE 1.2. The Ladder of Causation, with representative organisms at each level. Most animals, as well as present-day learning machines, are on the first rung, learning from association. Tool users, such as early humans, are on the second rung if they act by planning and not merely by imitation. We can also use experiments to learn the effects of interventions, and presumably this is how babies acquire much of their causal knowledge. Counterfactual learners, on the top rung, can imagine worlds that do not exist and infer reasons for observed phenomena. (*Source:* Drawing by Maayan Harel.)
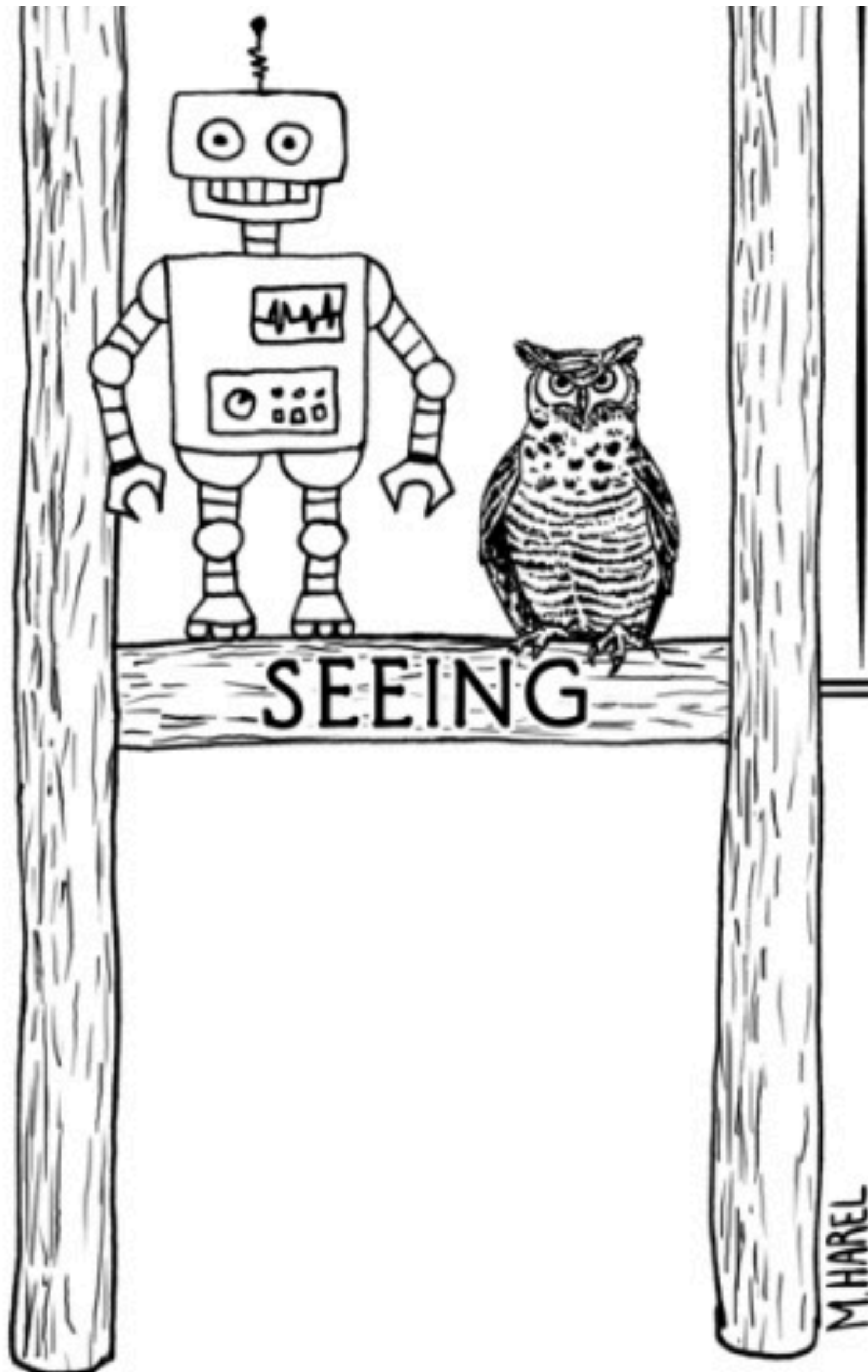
## The morale of the story

The morale of this story is summed up in the following picture:

1. ASSOCIATION
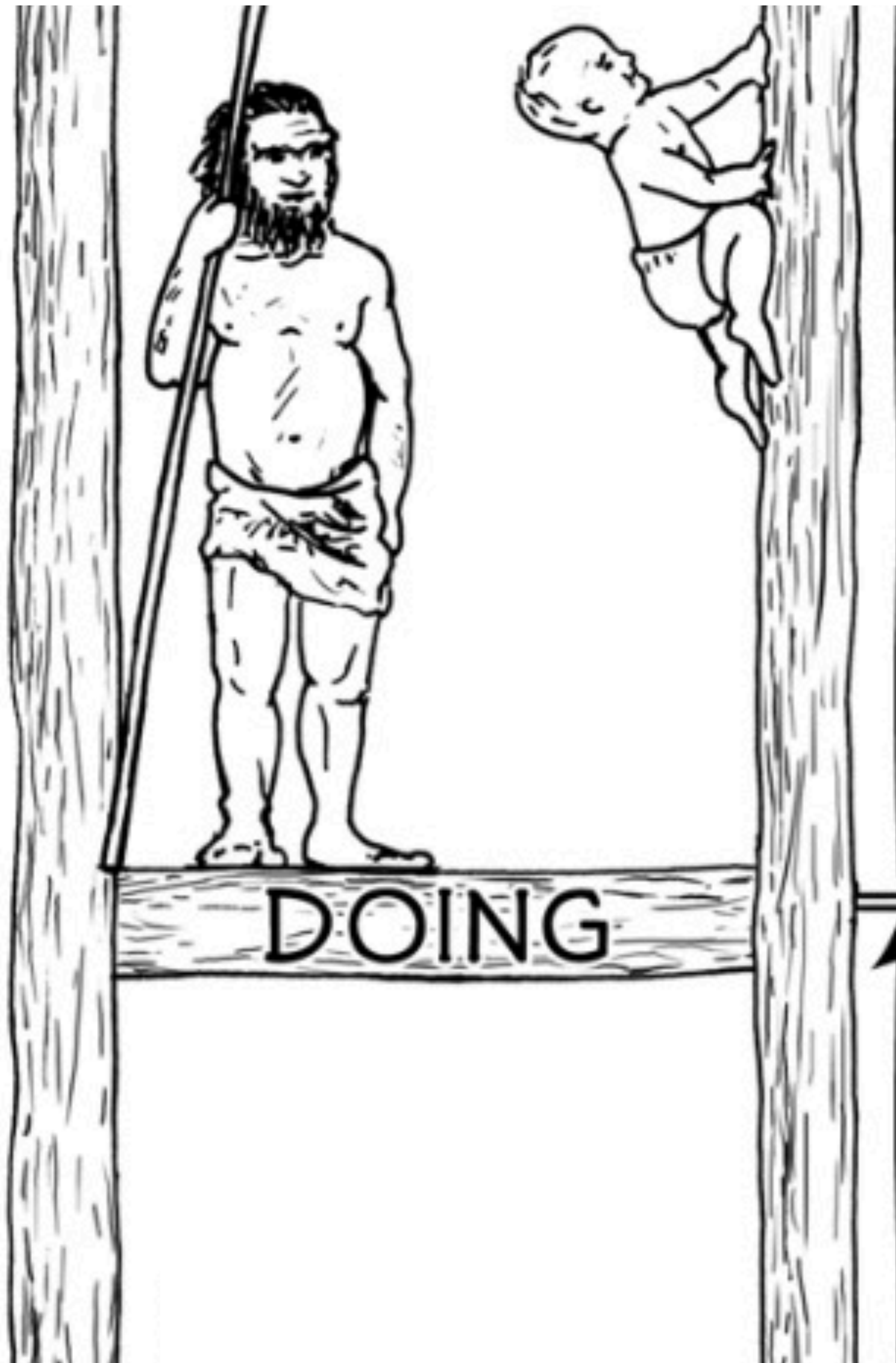
ACTIVITY: Seeing, Observing

QUESTIONS: *What if I see ...?*
(How are the variables related?
How would seeing X change my belief in Y?)

EXAMPLES: What does a symptom tell me about a disease?
What does a survey tell us about the
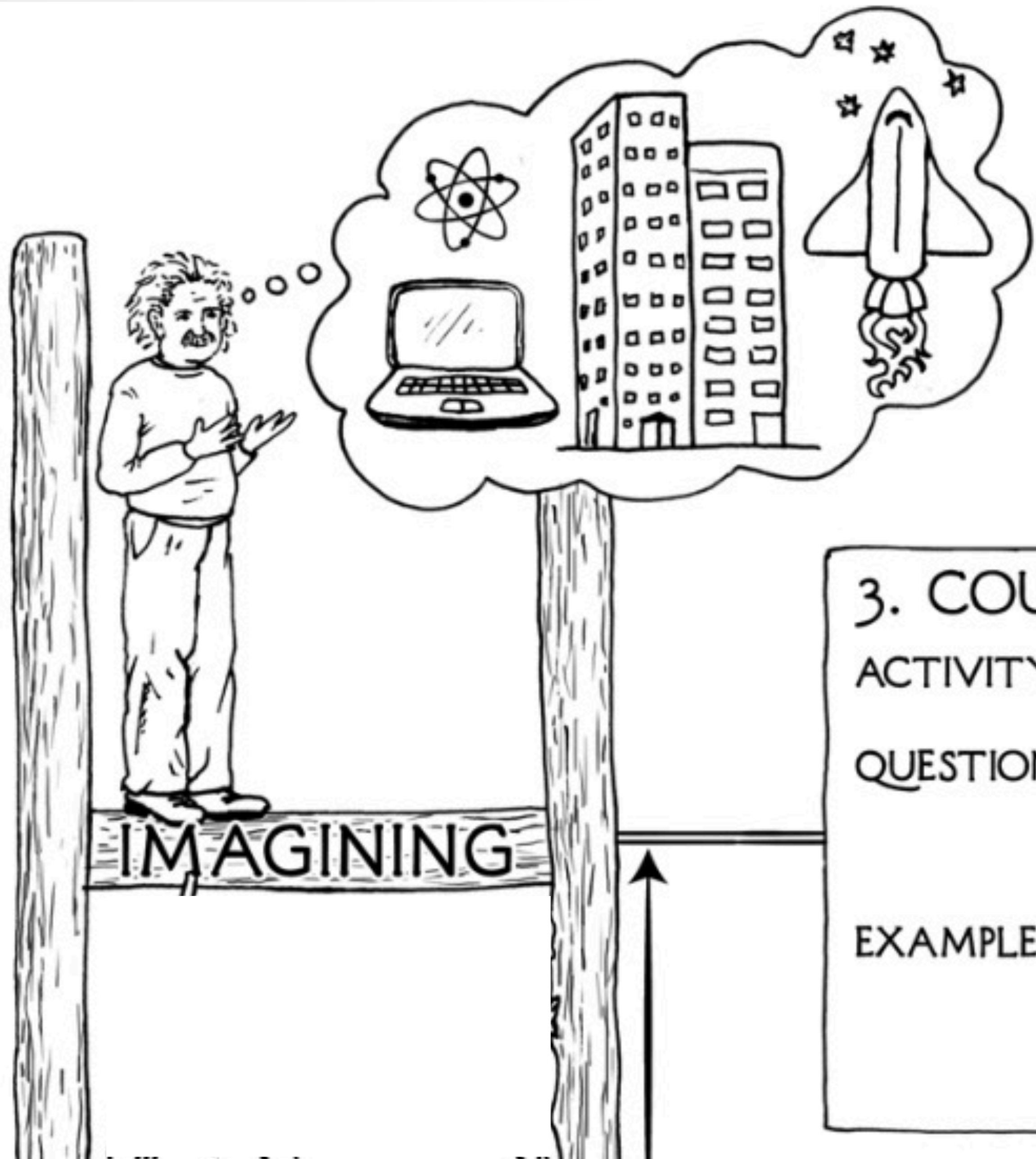election results?

**DOING**

### 2. INTERVENTION

**ACTIVITY:** Doing, Intervening

**QUESTIONS:** *What if I do …? How?*
(What would Y be if I do X?
How can I make Y happen?)

**EXAMPLES:** If I take aspirin, will my headache be cured?
What if we ban cigarettes?

## 3. COUNTERFACTUALS

ACTIVITY: Imagining, Retrospection, Understanding

QUESTIONS: *What if I had done …? Why?*
(Was it X that caused Y? What if X had not occurred? What if I had acted differently?)

EXAMPLES: Was it the aspirin that stopped my headache? Would Kennedy be alive if Oswald had not killed him? What if I had not smoked for the last 2 years?

Yoshua Bengio on [arXiv:1901.10912] and public FB discussion

# Causally interpreted Bayesian networks

**Definition:** A DAG $G = (V, E)$ together with a family $\kappa_v : \mathbb{X}_{pa(v)} \times \mathbb{X}_v \to [0,1]$, $v \in V$, of Markov kernels is called *Bayesian network*. The interpretation of the Markov kernels $\kappa_v$ as mechanisms of the nodes implies a causal nature of the Bayesian network.

Consider two distinct nodes $v, w \in V$.
- If $v \to w$ we call $v$ a *(pot.)* *direct cause* of $w$ and $w$ a *(pot.)* *direct effect* of $v$.
- If $v \rightsquigarrow w$ we call $v$ a *(pot.)* *cause* of $w$ and $w$ an *(pot.)* *effect* of $v$.



• J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press 2000, 2009.

```
x = randn()
y = x + 1 + sqrt(3)*randn()
```

```
y = 1 + 2*randn()
x = (y-1)/4 + sqrt(3)*randn()/2
```

```
z = randn()
y = z + 1 + sqrt(3)*randn()
x = z
```

# Intervention and the *do*-operation



$$\hat{\alpha}(a) := \begin{cases} 1, & \text{if } a = \hat{a} \\ 0, & \text{if } a \neq \hat{a} \end{cases}$$

$$p(u, a, b, c) = \varphi(u)\, \alpha(u; a)\, \beta(u; b)\, \gamma(a, b; c)$$

$$p(u, a, b, c \,|\, do(\hat{a})) := \varphi(u)\, \hat{\alpha}(a)\, \beta(u; b)\, \gamma(a, b; c)$$

## Intervention and the *do*-operation

$$p(b \,|\, do(\hat{a})) \;=\; \sum_{u,a,c} p(u,a,b,c \,|\, do(\hat{a}))$$

$$= \sum_{u,a,c} \varphi(u)\,\hat{\alpha}(a)\,\beta(u;b)\,\gamma(a,b;c)$$

$$= \sum_{u,c} \varphi(u)\,\beta(u;b)\,\gamma(\hat{a},b;c)$$

$$= \sum_{u} \varphi(u)\,\beta(u;b) \sum_{c} \gamma(\hat{a},b;c)$$

$$= p(b) \;\neq\; p(b \,|\, \hat{a})$$

## Intervention and the *do*-operation

$$p(c \,|\, do(\hat{u})) \;=\; \sum_{u,a,b} p(u,a,b,c \,|\, do(\hat{u}))$$

$$= \sum_{u,a,b} \hat{\varphi}(u)\,\alpha(u;a)\,\beta(u;b)\,\gamma(a,b;c)$$

$$= \sum_{a,b} \alpha(\hat{u};a)\,\beta(\hat{u};b)\,\gamma(a,b;c)$$

$$= \frac{\sum_{a,b} \varphi(\hat{u})\,\alpha(\hat{u};a)\,\beta(\hat{u};b)\,\gamma(a,b;c)}{\varphi(\hat{u})}$$

$$= \frac{\sum_{a,b} p(\hat{u},a,b,c)}{p(\hat{u})}$$

$$= p(c \,|\, \hat{u})$$

## When do they coincide?

**Theorem:** *Let $G = (V,E)$ be a DAG, and let $A$ and $B$ be two disjoint subsets of $V$. Then the following two statements are equivalent:*
**(1)** *For every Bayesian network $\mathfrak{B}$ with underlying DAG $G$ we have*

$$p(x_B \,|\, do(x_A)) = p(x_B \,|\, x_A).$$

**(2)** *$B$ is not a cause of $A$ and there is no common cause of $A$ and $B$.*

$$p(b \,|\, do(a)) = p(b \,|\, a)$$
$$p(c \,|\, do(b)) \neq p(c \,|\, b)$$
$$p(c \,|\, do(a)) \neq p(c \,|\, a)$$

$$P(y|do(X)) = p(y|x)$$

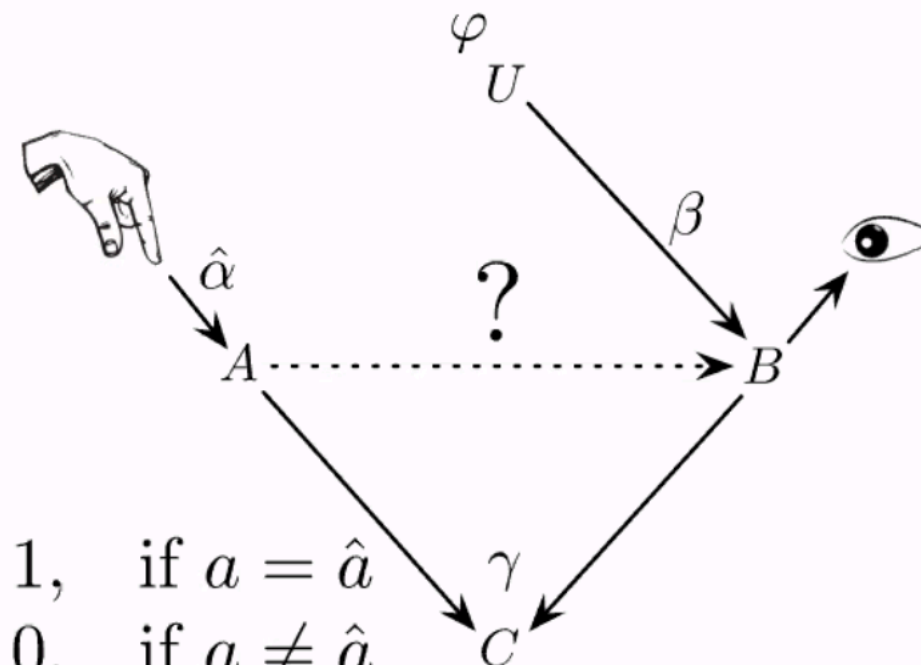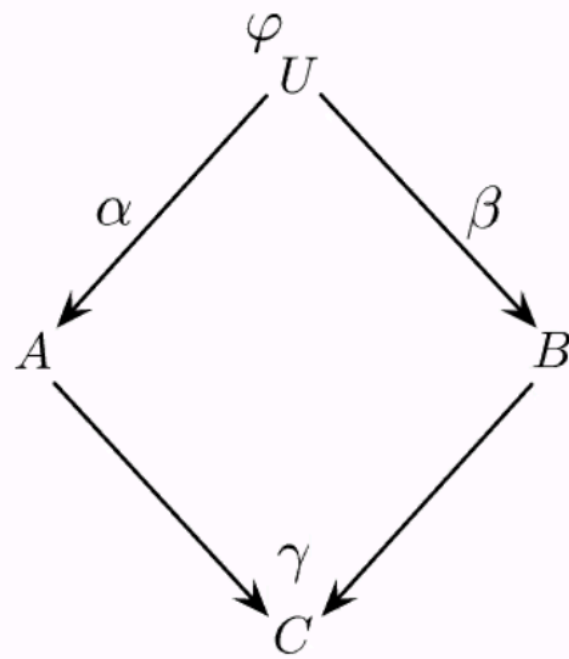$$P(y|do(X)) = p(y)$$

$$P(y|do(X)) = p(y)$$

```
x = randn()
x = 3
y = x + 1 + sqrt(3)*randn()
x = 3
```

```
y = 1 + 2*randn()
x = 3
x = (y-1)/4 + sqrt(3)*randn()/2
x = 3
```

```
z = randn()
x = 3
x = z
x = 3
y = z + 1 + sqrt(3)*randn()
x = 3
```

# Identifiability of causal effects

**Theorem:** *Let $\mathfrak{B}$ be a Bayesian network with DAG $G = (V, E)$. For two distinct nodes $v$ and $w$ the following holds:*

$$p(x_w \mid do(x_v)) = \sum_{x_{pa(v)}} p(x_{pa(v)}) \, p(x_w \mid x_v, x_{pa(v)}) \,.$$

$$p(x_w \mid do(x_v)) =$$

$$\sum_{x_i, x_j} p(x_i, x_j) \, p(x_w \mid x_v, x_i, x_j)$$

3. COUNTERFACTUALS

ACTIVITY: Imagining, Retrospection, Understanding

QUESTIONS: *What if I had done …? Why?*
(Was it X that caused Y? What if X had not occurred? What if I had acted differently?)

EXAMPLES: Was it the aspirin that stopped my headache? Would Kennedy be alive if Oswald had not killed him? What if I had not smoked for the last 2 years?

### Example 1: David Blei's election example

This is an example David brought up during the **Causality Panel** and I referred back to this in my talk. I'm including it here for the benefit of those who attended my MLSS talk:

*Given that Hilary Clinton did not win the 2016 presidential election, and given that she did not visit Michigan 3 days before the election, and given everything else we know about the circumstances of the election, what can we say about the probability of Hilary Clinton winning the election, had she visited Michigan 3 days before the election?*

Let's try to unpack this. We are are interested in the probability that:

- she *hypothetically* wins the election

conditionied on four sets of things:

- she lost the election
- she did not visit Michigan
- any other relevant an observable facts
- she *hypothetically* visits Michigan

It's a weird beast: you're simultaneously conditioning on her visiting Michigan and not visiting Michigan. And you're interested in the probability of her winning the election given that she did not. WHAT?

Why would quantifying this probability be useful? Mainly for credit assignment. We want to know why she lost the election, and to what degree the loss can be attributed to her failure to visit Michigan three days before the election. Quantifying this is useful, it can help political advisors make better decisions next time.

# Identifiability of causal effects (front-door example)

$$p(c \mid do(a)) = \sum_{u,b} \varphi(u)\, \beta(a;b)\, \gamma(u,b;c)$$

$$= \sum_{u,b} p(u)\, p(b \mid a)\, p(c \mid u,b)$$

$$= \sum_{u,b} \left( \sum_{a'} p(a')\, p(u \mid a') \right) p(b \mid a)\, p(c \mid u,b)$$

$$= \sum_{b} p(b \mid a) \sum_{a'} p(a') \sum_{u} p(u \mid a')\, p(c \mid u,b)$$

$$= \sum_{b} p(b \mid a) \sum_{a'} p(a') \sum_{u} p(u \mid a',b)\, p(c \mid u,a',b)$$

$$= \sum_{b} p(b \mid a) \sum_{a'} p(a')\, p(c \mid a',b)$$

# In robotics



## Causal effects in the sensorimotor loop



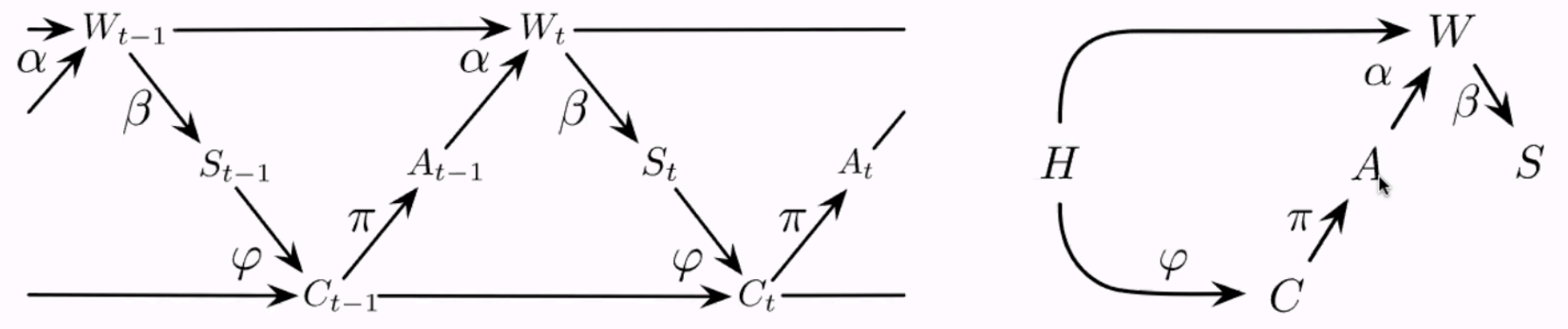$$p(s \mid do(a)) \quad = \quad \sum_c p(s \mid c, a)\, p(c)$$

• N. Ay, K. Zahedi. *On the Causal Structure of the Sensorimotor Loop*. GSO 2014.

• N. Ay, K. Zahedi. *An Information-Theoretic Approach to Prediction and Deliberative Decision Making of Embodied Systems*. Proc. of ICCN 2011. Advances in Cognitive Neurodynamics, Springer 2012.

**Causally Correct Partial Models for Reinforcement Learning**

Danilo J. Rezende [*1]   Ivo Danihelka [*12]   George Papamakarios [1]   Nan Rosemary Ke [3]   Ray Jiang [1]
Theophane Weber [1]   Karol Gregor [1]   Hamza Merzic [1]   Fabio Viola [1]   Jane Wang [1]   Jovana Mitrovic [1]
Frederic Besse [1]   Ioannis Antonoglou [12]   Lars Buesing [1]
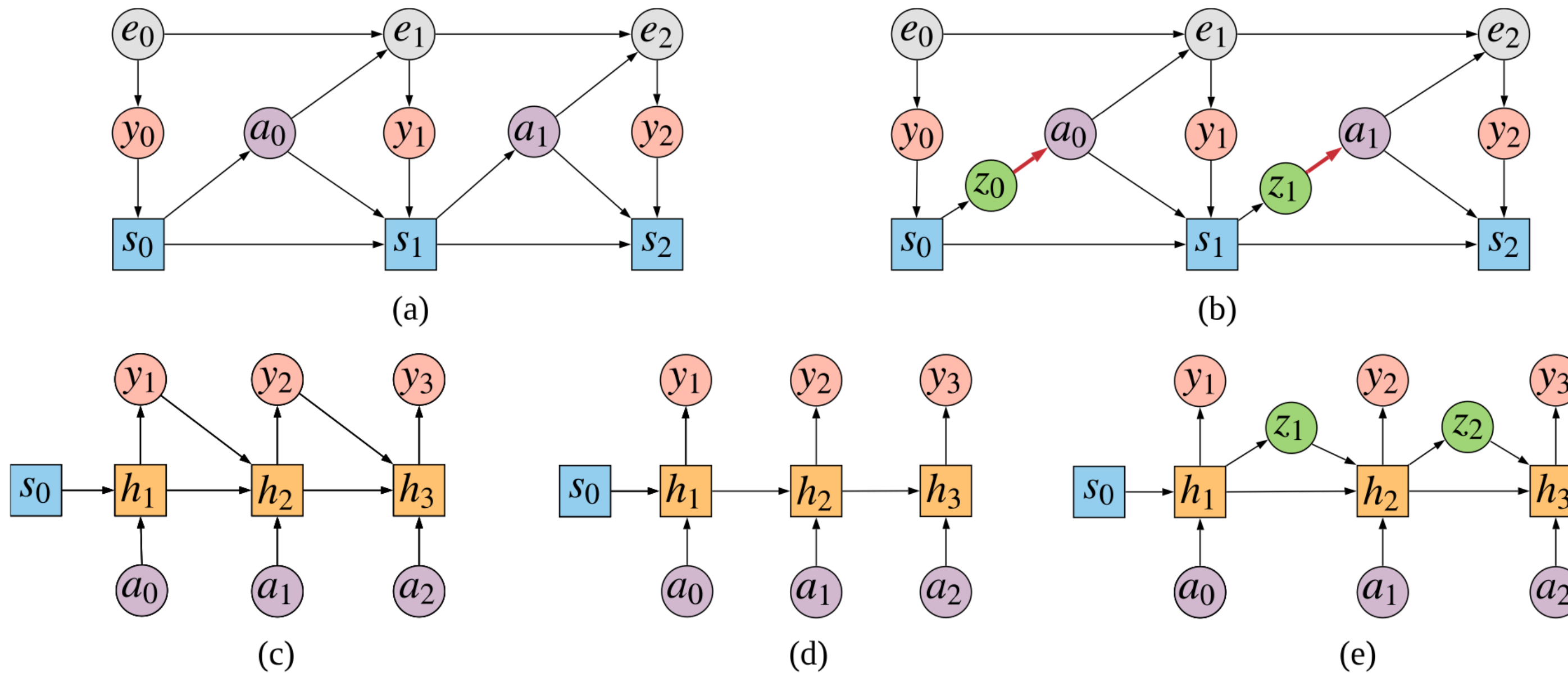
**Causally Correct Partial Models**



*Figure 3.* Graphical representations of the environment, the agent, and the various models. Circles are stochastic nodes, rectangles are deterministic nodes. (a) Agent interacting with the environment, generating a trajectory $\{y_t, a_t\}_{t=0}^{T}$. These trajectories are the training data for the models. (b) Same as (a) but also including the backdoor $z_t$ in the generated trajectory. The red arrows indicate the locations of the interventions. (c) Standard autoregressive generative model of observations. The model predicts the observation $y_t$ which it then feeds into $h_{t+1}$. (d) Example of a Non-Causal Partial Model (NCPM) that predicts the observation $y_t$ without feeding it into $h_{t+1}$. (e) Proposed Causal Partial Model (CPM), with a backdoor $z_t$ for the actions.