

# Causal discovery and Deep Learning

Faria et. al. 2022; Pearl et. al. 2018; Neal et. al. 2022; Schölkopf et. al. 2019; Bengio et. al. 2020; Peters et. al. 2017, Ke et. al. 2020, Brouillard et. al. 2020

Gonçalo R. A. Faria  
@goncalorafaria



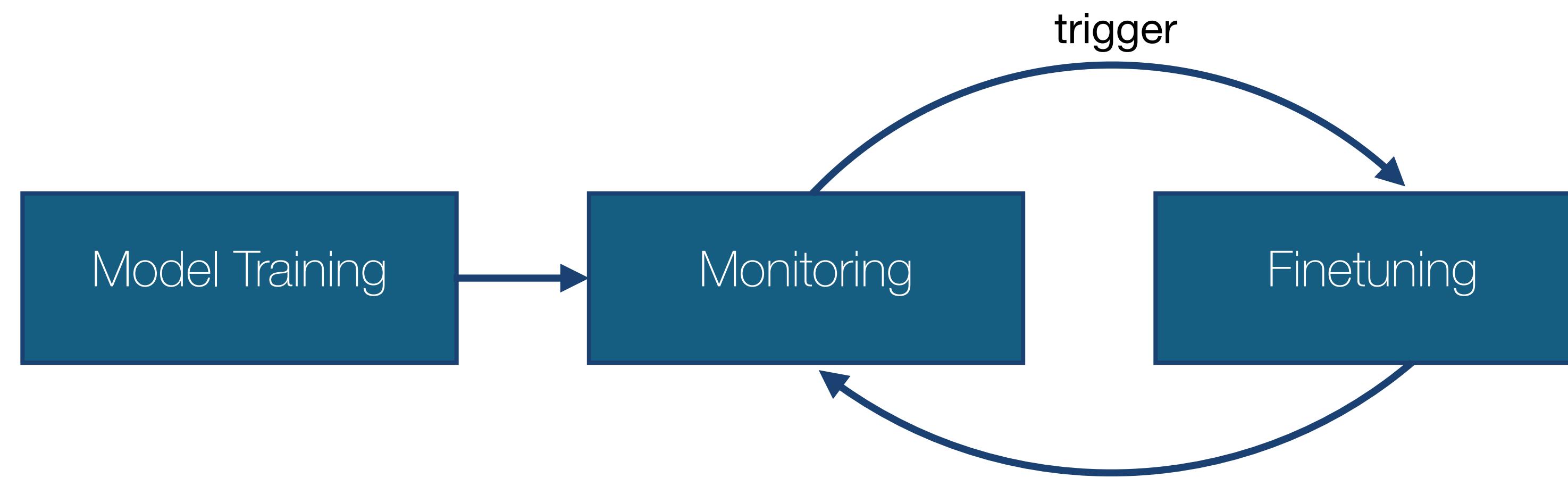
Deep Learning Sessions Portugal

# Overview

- Limitations of current Deep Learning systems
- Causal Discovery
- Identifiability
- Interventions and Causal Discovery
- Latent interventions and Causal Discovery

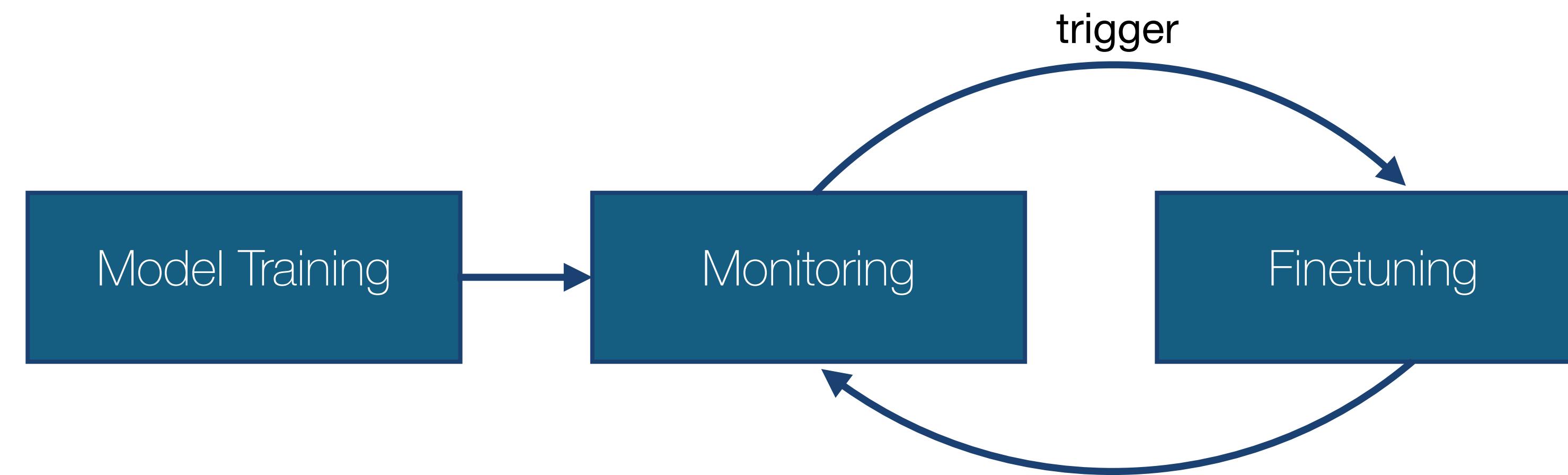
# Limitations of current Deep Learning

- Distribution shifts.



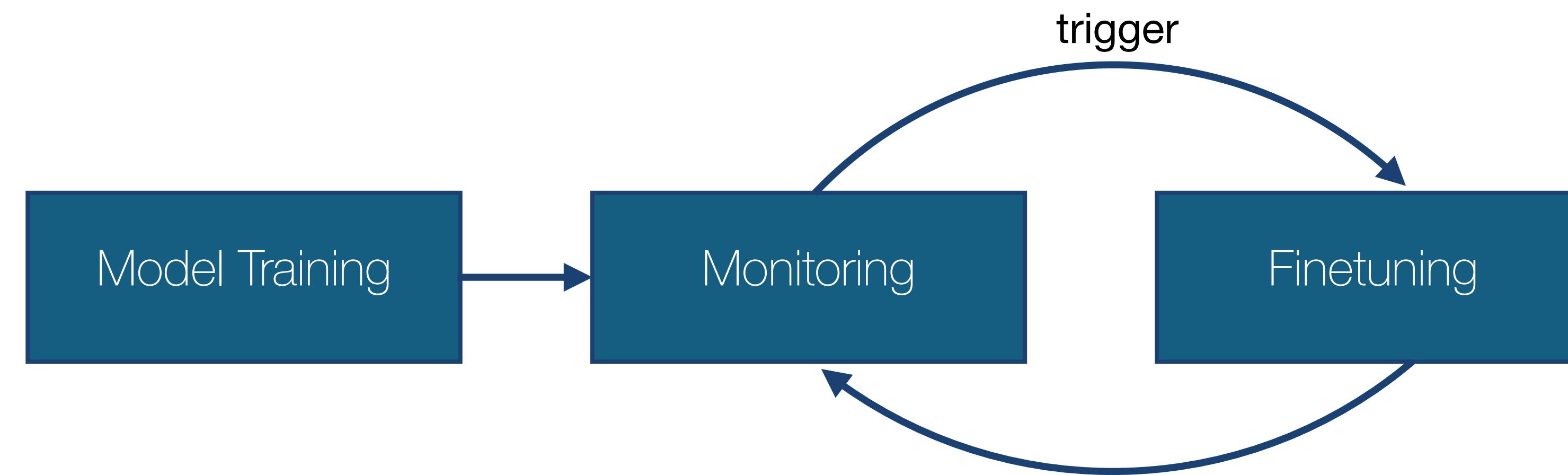
# Limitations of current Deep Learning

- Distribution shifts.
- Model deployment drives distribution shifts.



# Limitations of current Deep Learning

- Distribution shifts.
- Model deployment drives distribution shifts.
- Test-set performance vs generalization.
- Interpretability and scientific understanding



# Benefits of causal predictors

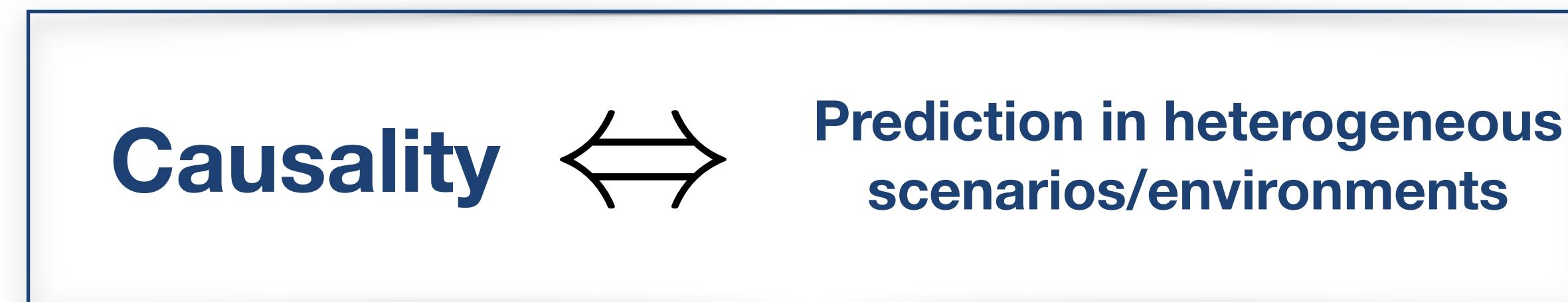
- Causal predictors are optimally robust.

Consider the setting :

$$(X^e, Y^e) \sim F^e, \quad e \in \mathcal{E} \subset \mathcal{F}$$

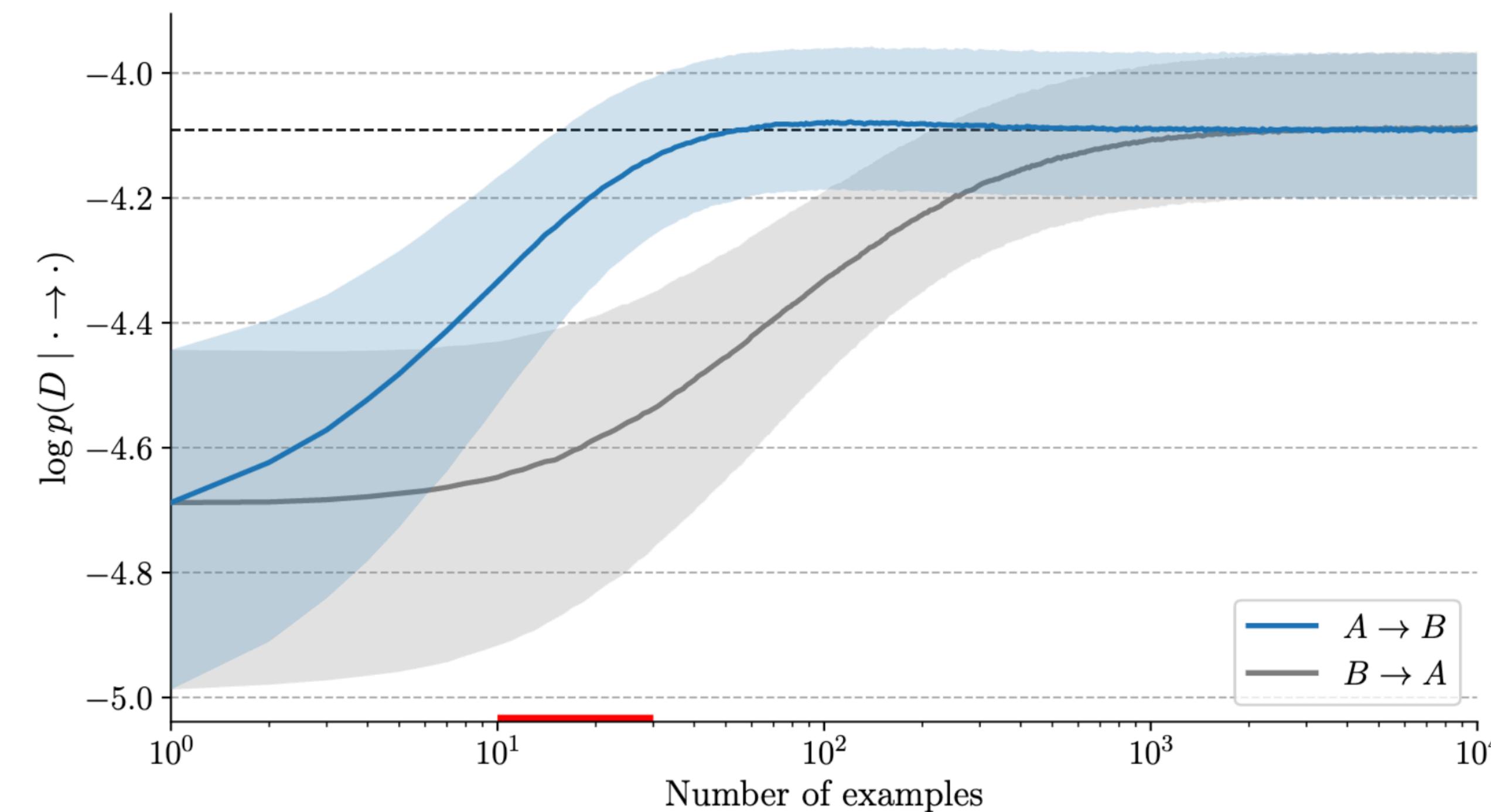
$$\hat{\beta} = \operatorname{argmin}_{\beta} \max_{e \in \mathcal{F}} \mathbb{E} |Y^e - (X^e)^T \beta|^2$$

**causal predictor**



# Benefits of causal predictors

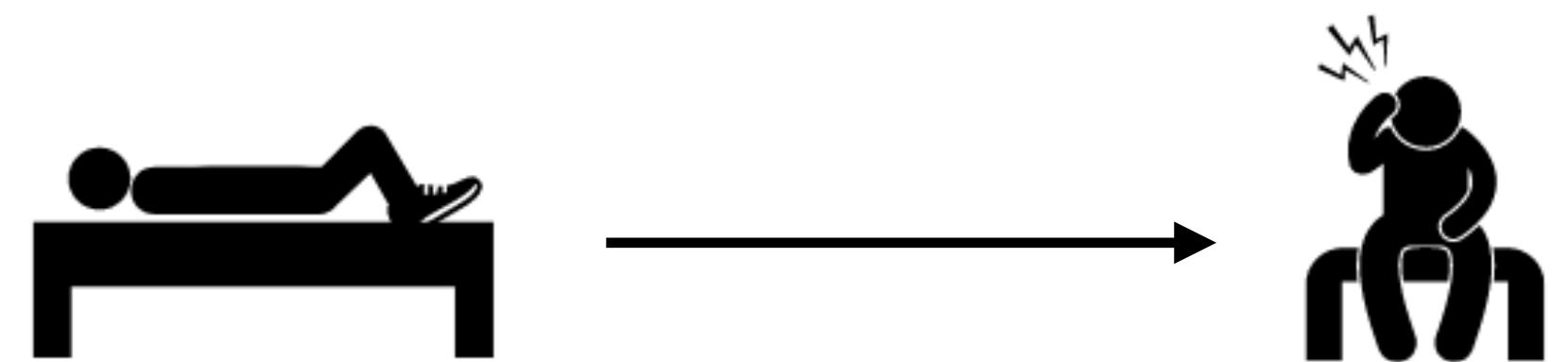
- Causal predictors are optimally robust.
- Fast and data efficient adaptation to transfer distribution.



**What if we do not know the causal model ?**

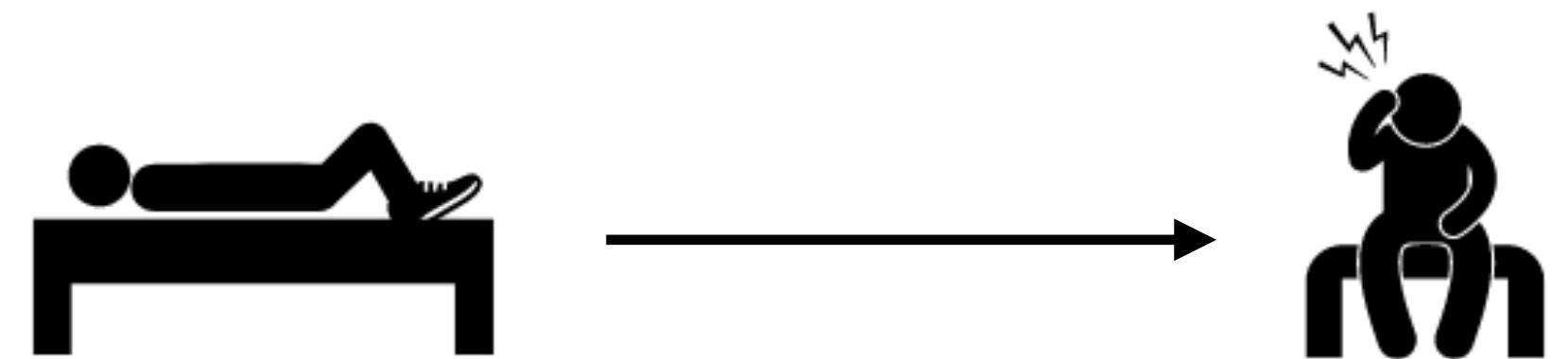
# Causal Discovery

- Say we have a dataset.



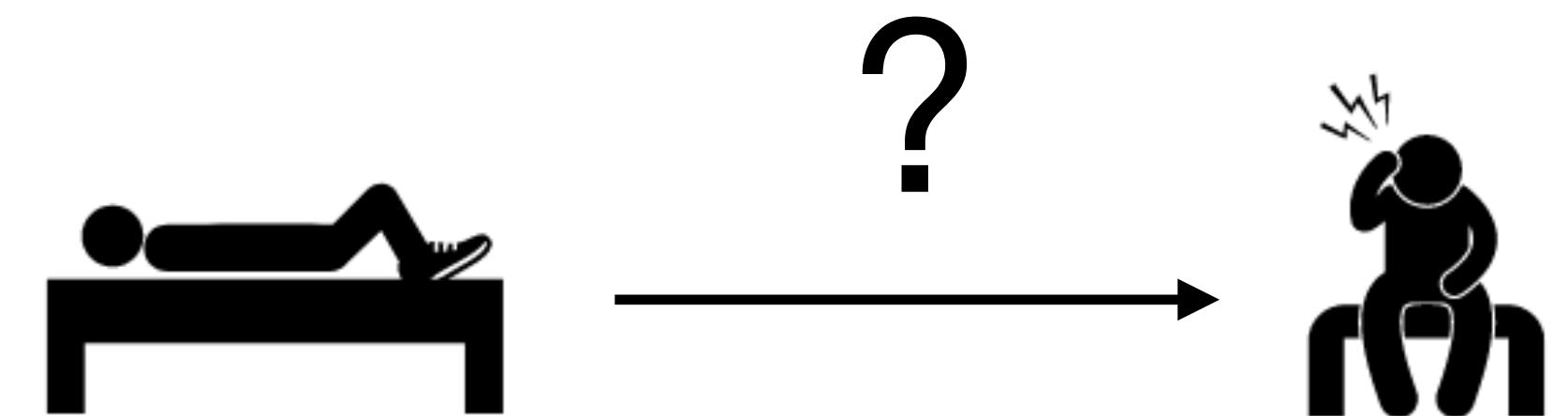
# Causal Discovery

- Say we have a dataset.
- Cause-effect relation?



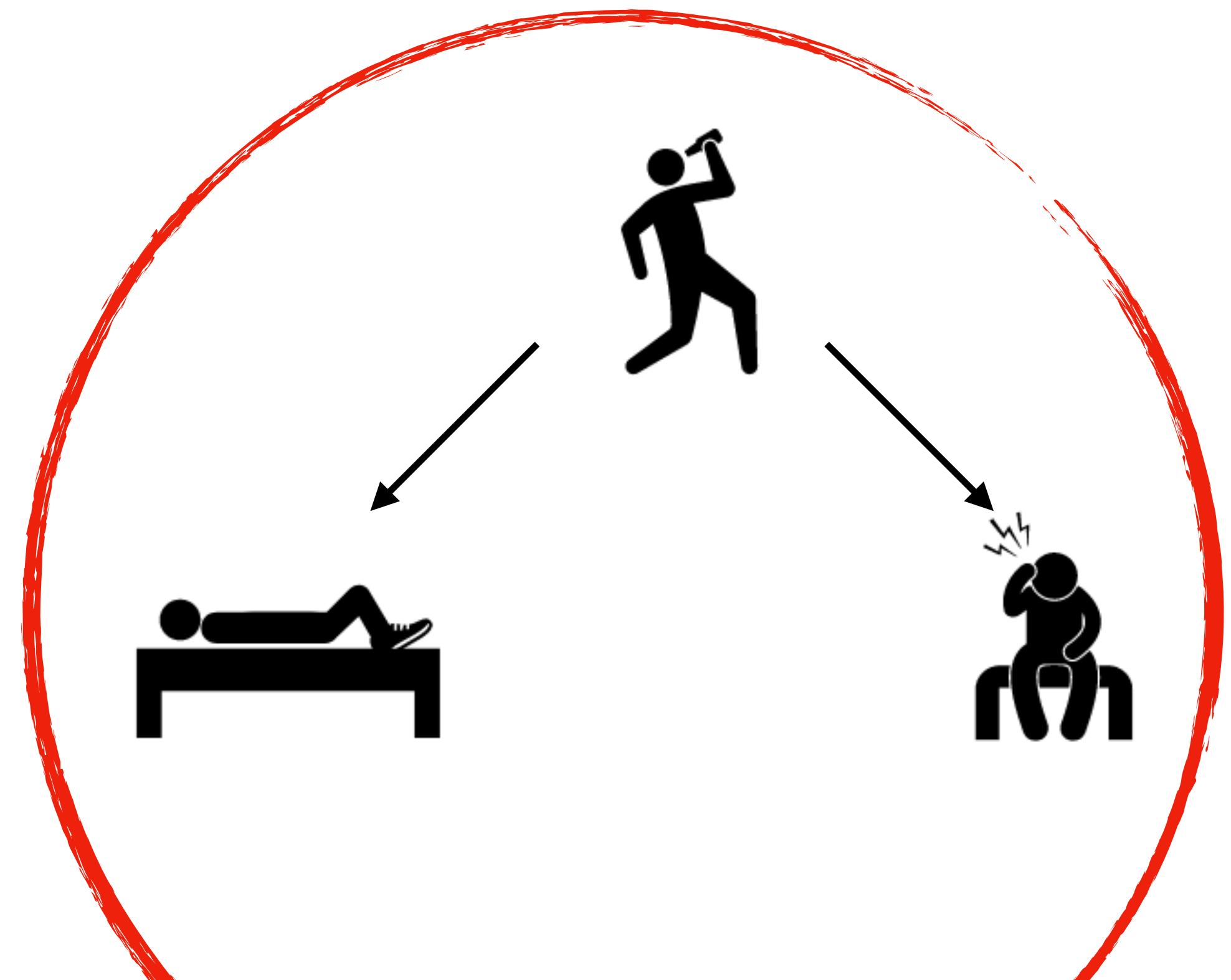
# Causal Discovery

- Say we have a dataset.
- Cause-effect relation?
- **Correlation is not causation**



# Causal Discovery

- Say we have a dataset.
- Cause-effect relation?
- **Correlation is not causation**
- Domain-expertise
- Hidden Confounding



# Data-driven Causal Discovery



# Identifiability



independent from



given



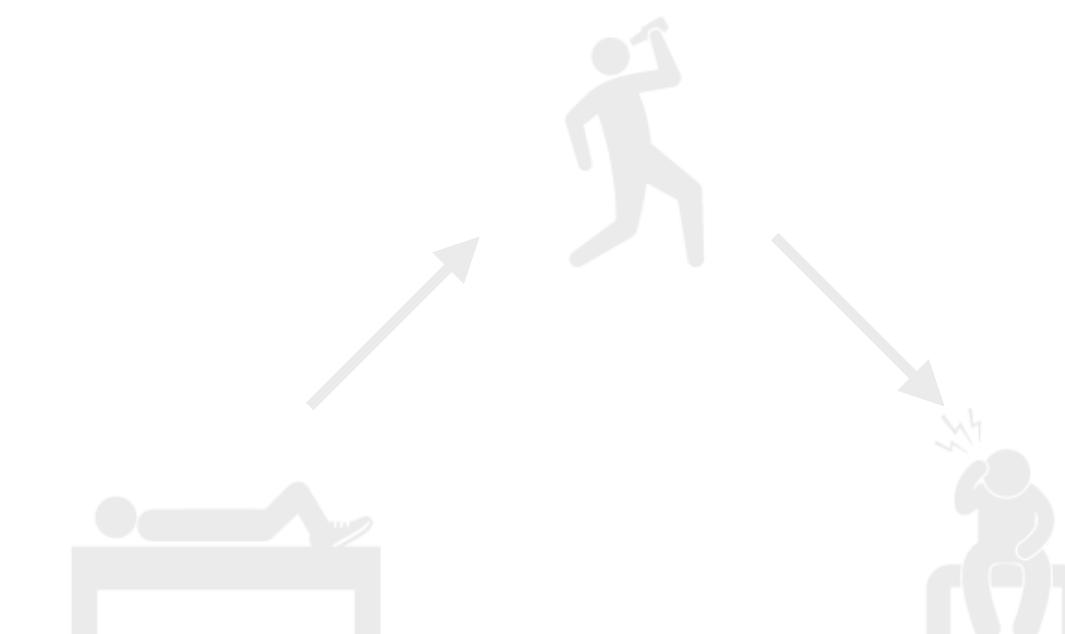
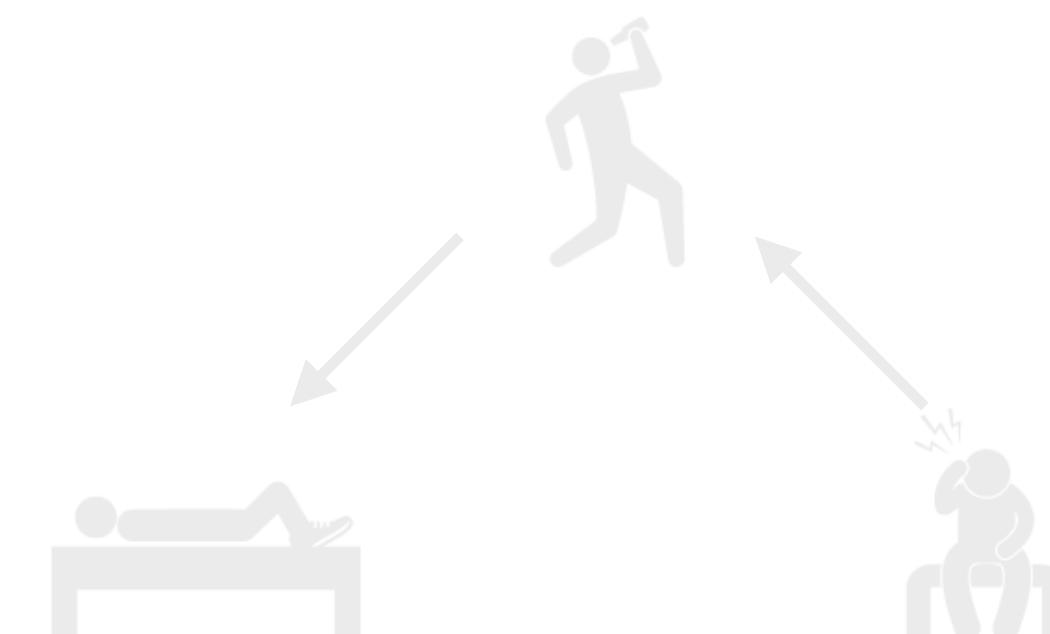
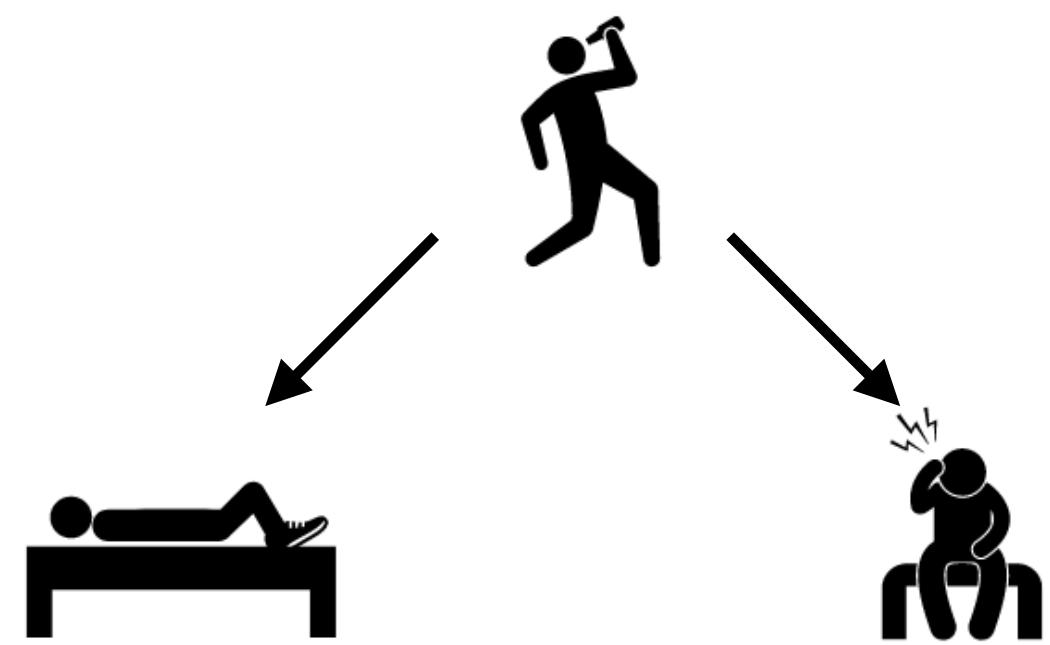
# Identifiability



independent from



given



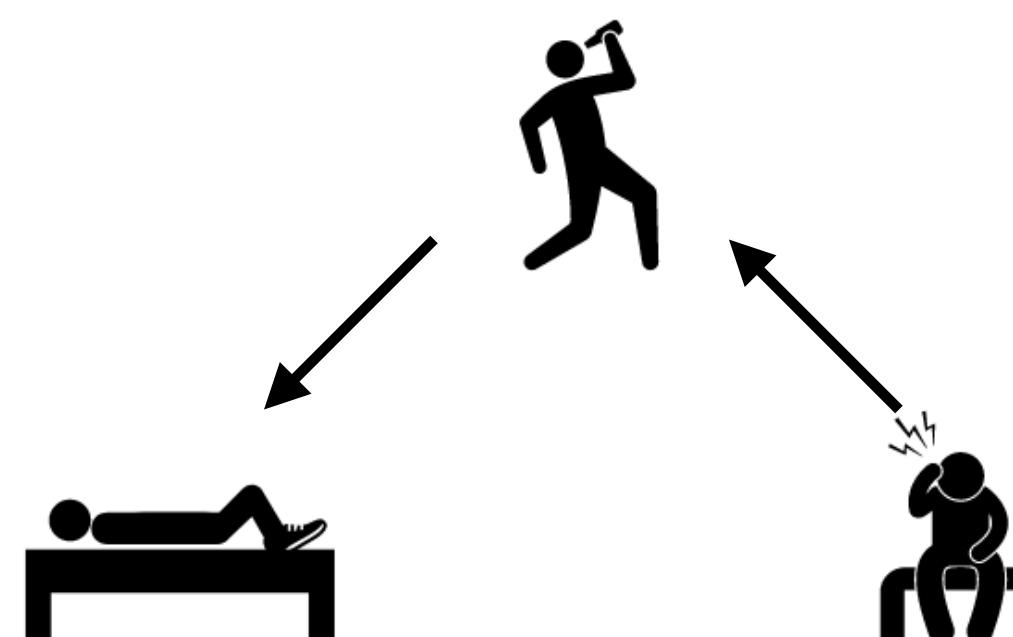
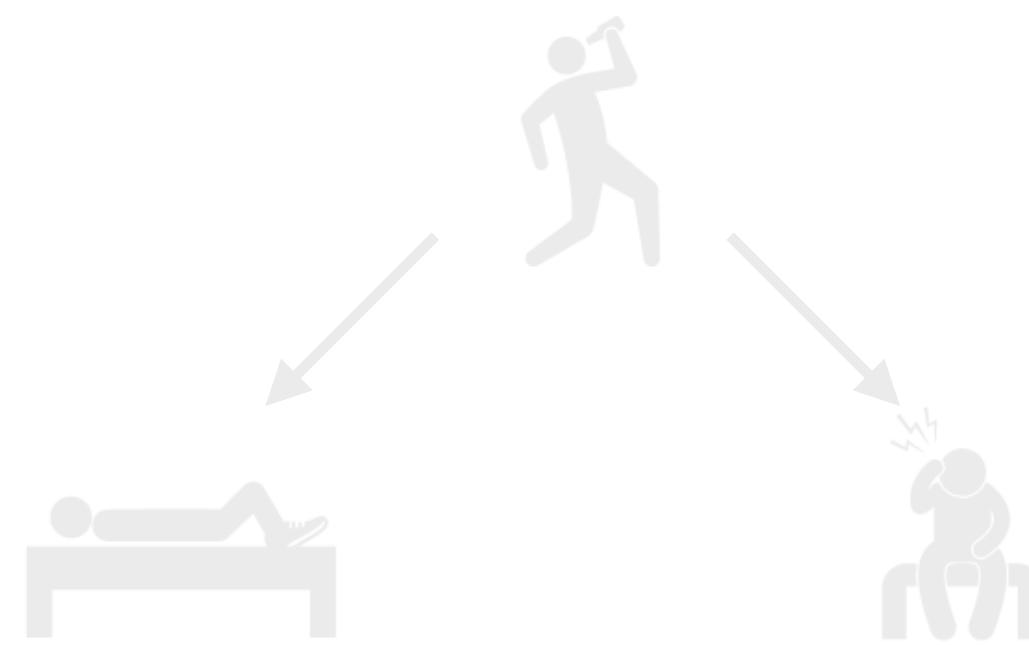
# Identifiability



independent from



given



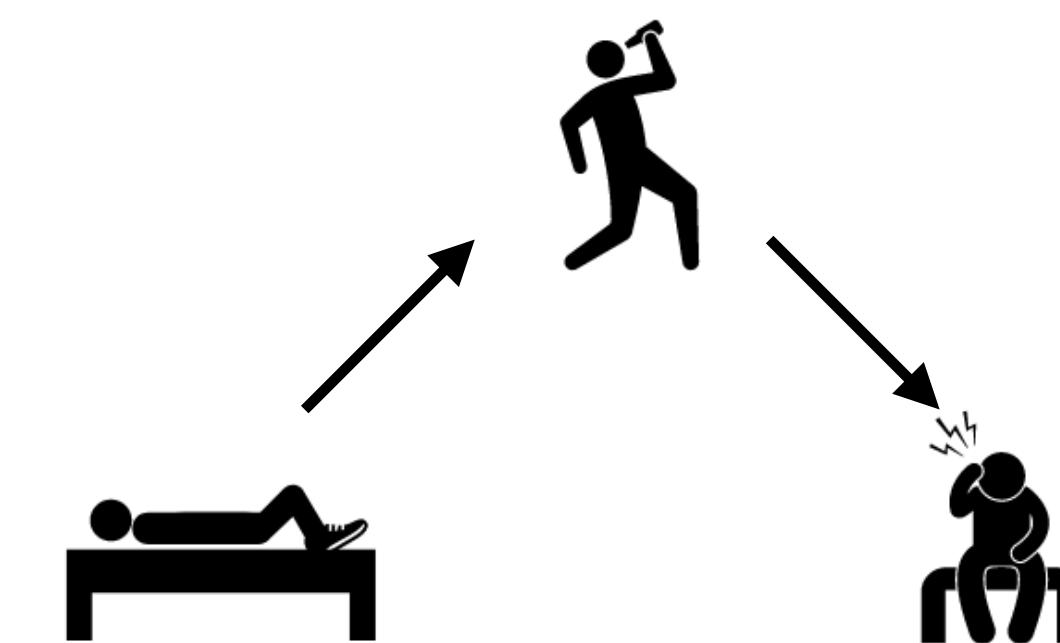
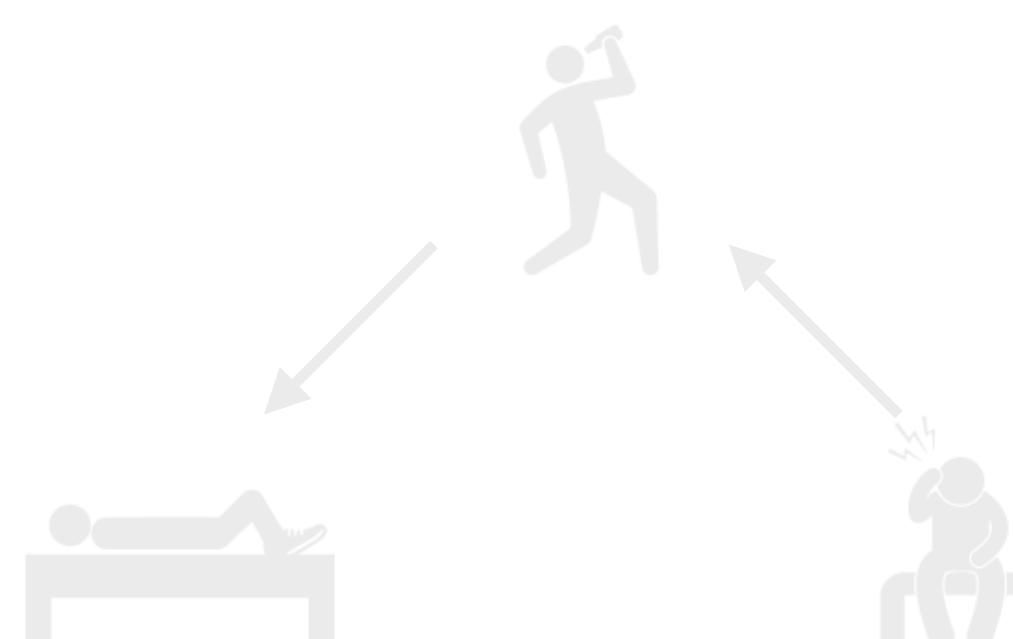
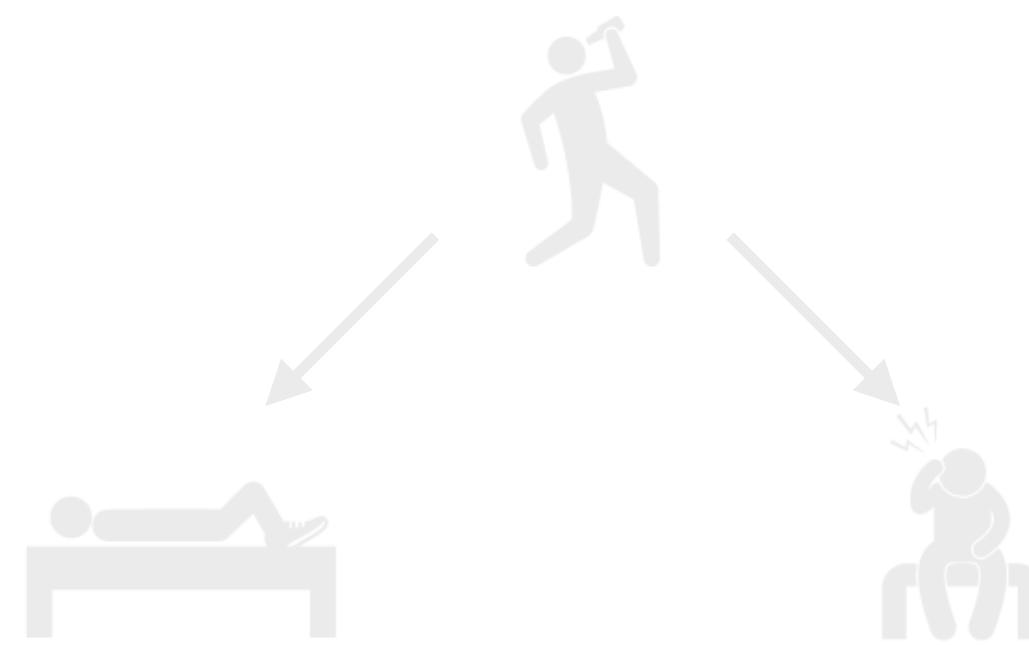
# Identifiability



independent from



given



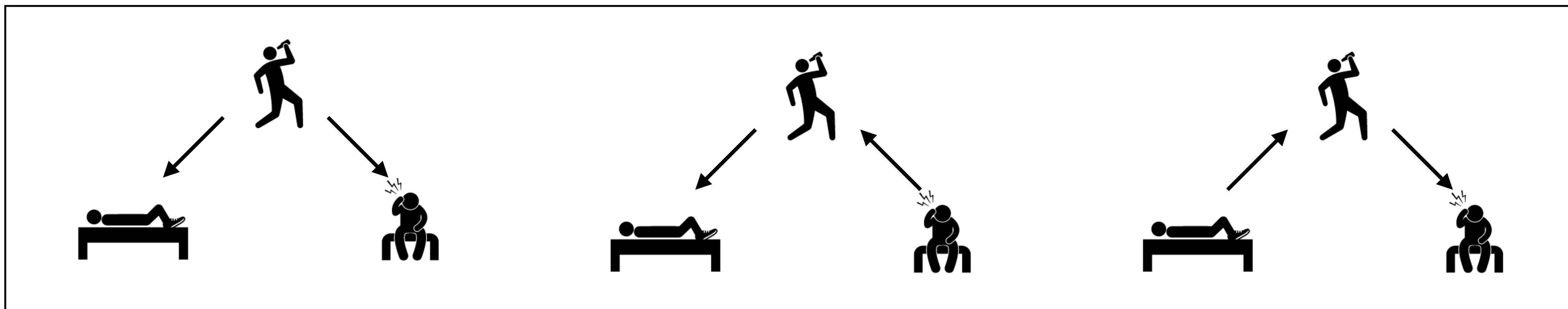
# Identifiability



independent from



given



Markov Equivalence class

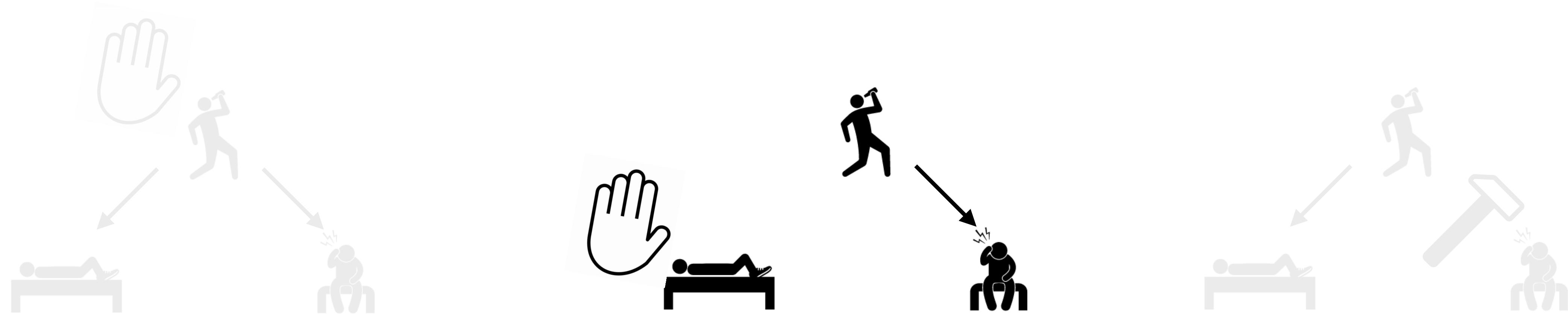
# Interventions

- Science; hypotheses; experiments (interventions)



# Interventions

- Science; hypotheses; experiments (interventions)



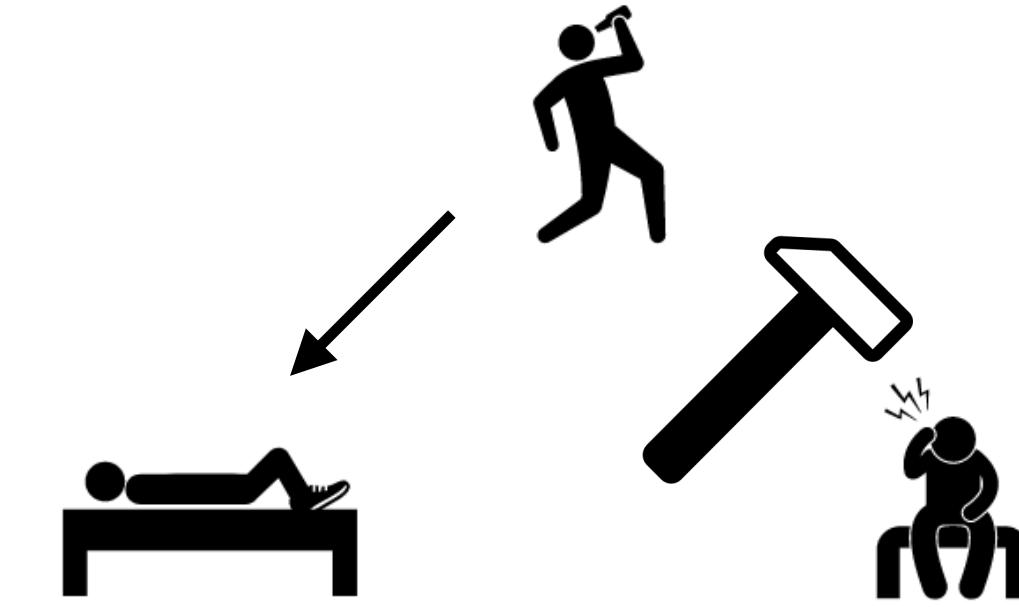
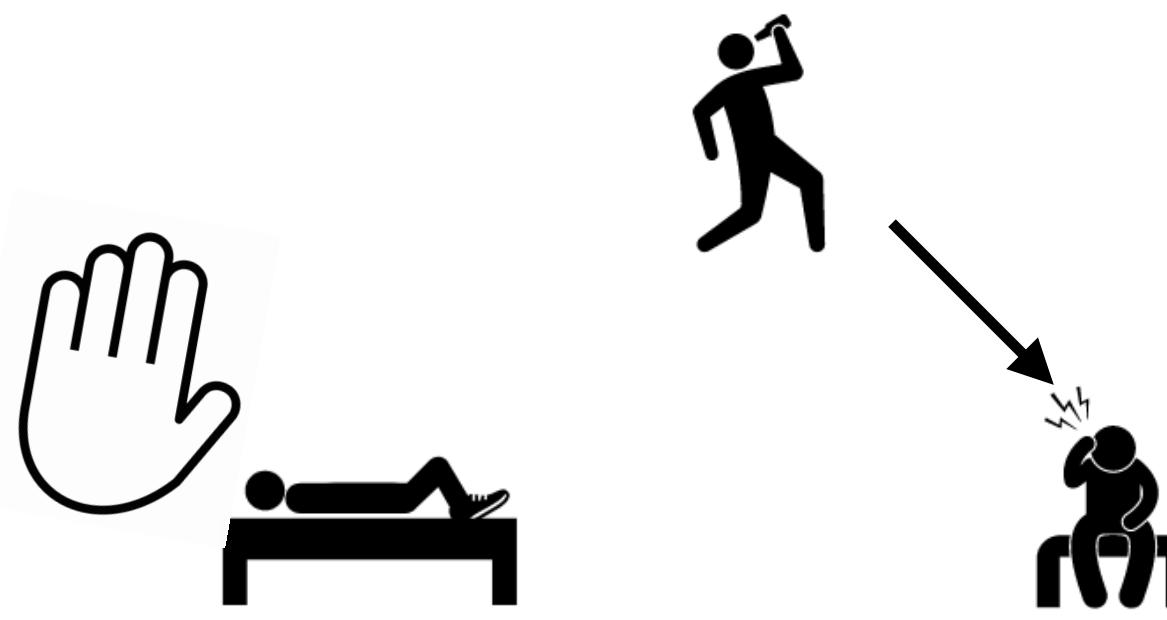
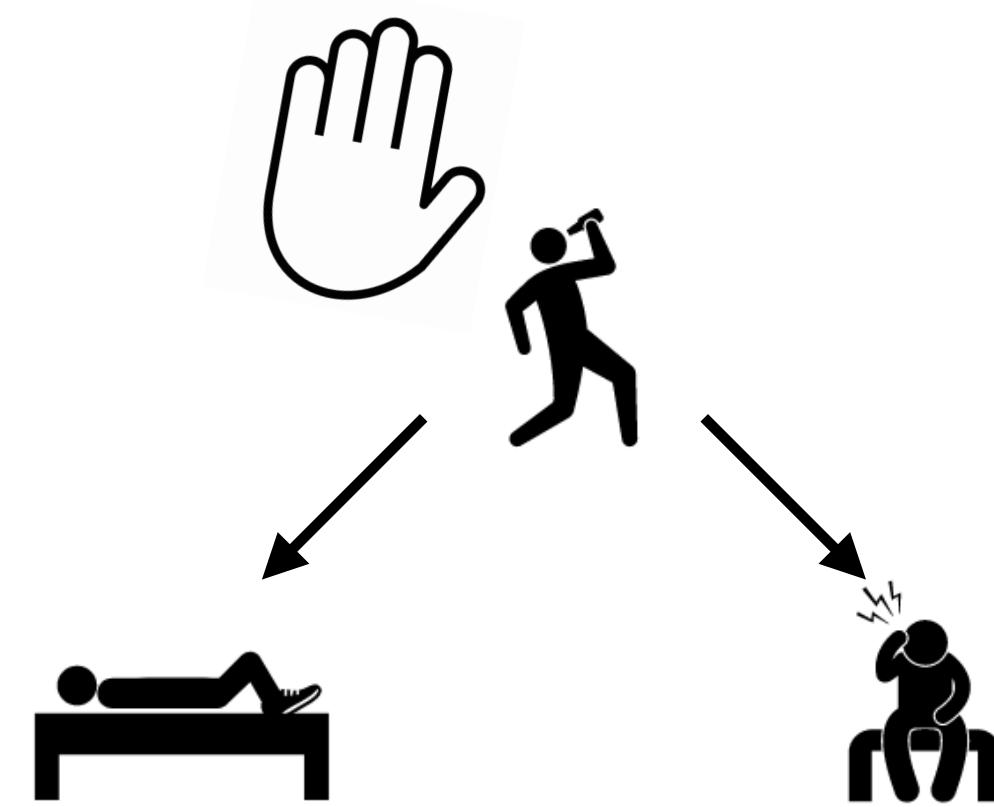
# Interventions

- Science; hypotheses; experiments (interventions)



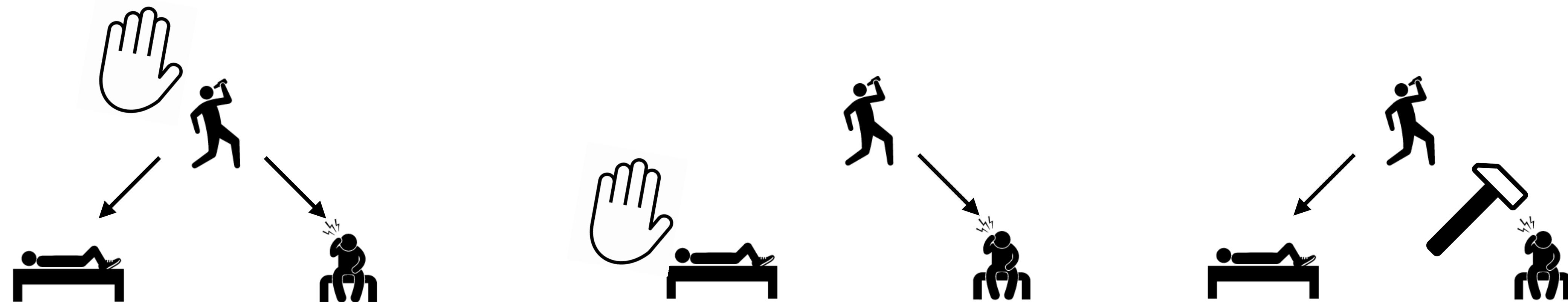
# Interventions

- Science; hypotheses; experiments (interventions)
- Causal discovery with intervention data



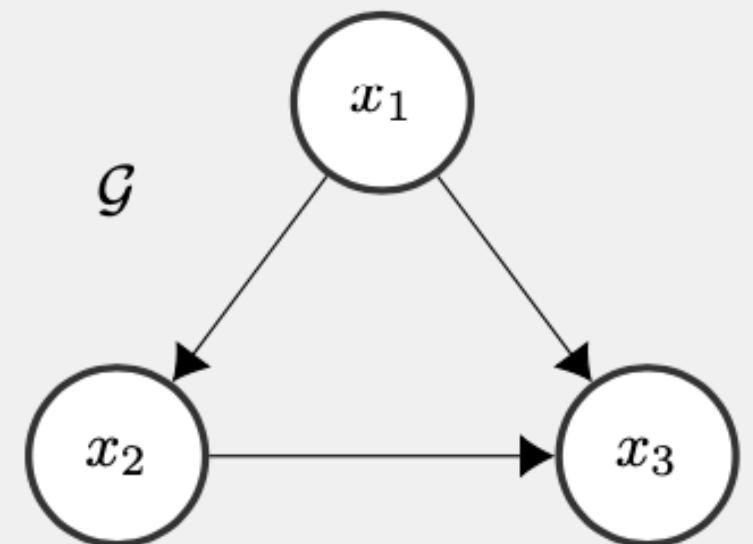
# Interventions

- Science; hypotheses; experiments (interventions)
- Causal discovery with intervention data
- **Complex; expensive; unethical.**



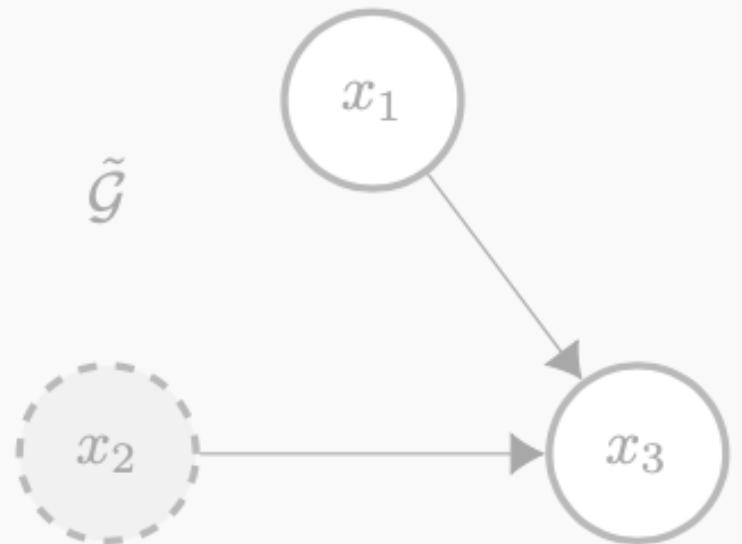
# Structural Causal Models

# Structural causal models



$$\begin{cases} x_1 := \varepsilon_1 \\ x_2 := 4x_1 + 1 + \varepsilon_2 \\ x_3 := 5 \sin(x_1) + x_2 + \varepsilon_3 \end{cases}, \quad \varepsilon_1, \varepsilon_2, \varepsilon_3 \sim \mathcal{N}(0, 1)$$

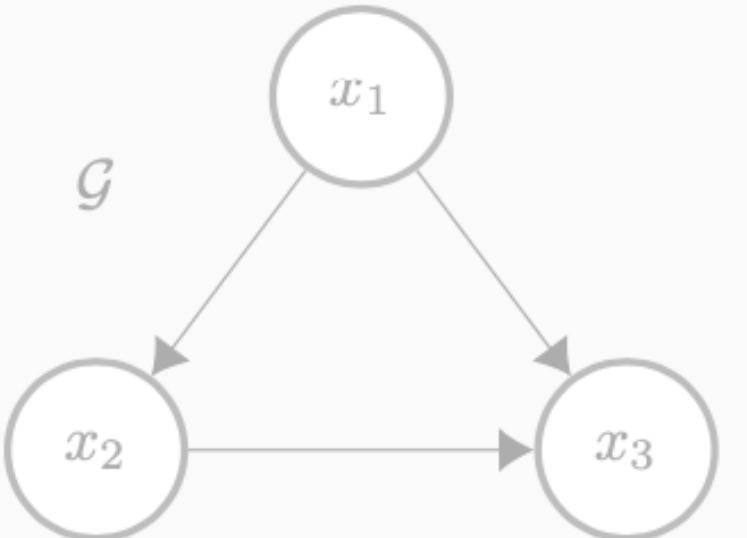
$$p_{\mathcal{M}}(X) = p_{\mathcal{M}}(x_3|x_2, x_1)p_{\mathcal{M}}(x_2|x_1)p_{\mathcal{M}}(x_1)$$



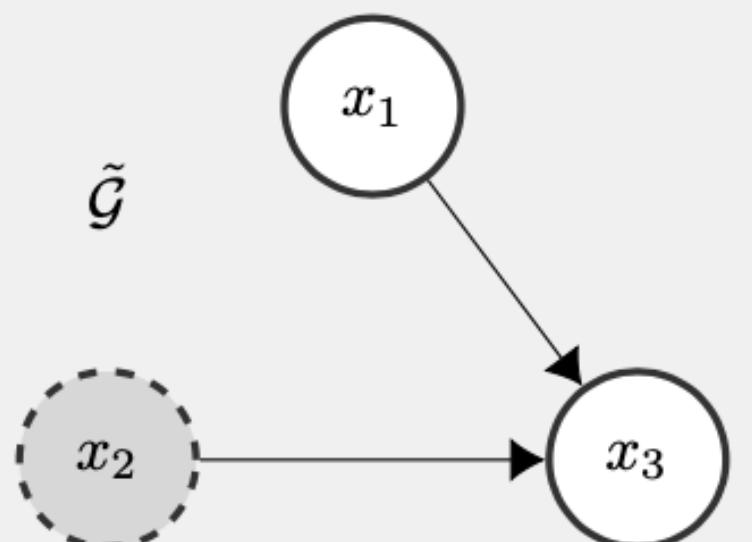
$$\begin{cases} x_1 := \varepsilon_1 \\ x_2 := 2 \\ x_3 := 5 \sin(x_1) + x_2 + \varepsilon_3 \end{cases}, \quad \varepsilon_1, \varepsilon_3 \sim \mathcal{N}(0, 1)$$

$$p_{\tilde{\mathcal{M}}}(X) = p_{\mathcal{M}; do(x_2:=2)}(X), \quad I = \{2\}$$

# Structural causal models



$$\begin{cases} x_1 := \varepsilon_1 \\ x_2 := 4x_1 + 1 + \varepsilon_2 \\ x_3 := 5 \sin(x_1) + x_2 + \varepsilon_3 \end{cases}, \quad \varepsilon_1, \varepsilon_2, \varepsilon_3 \sim \mathcal{N}(0, 1)$$

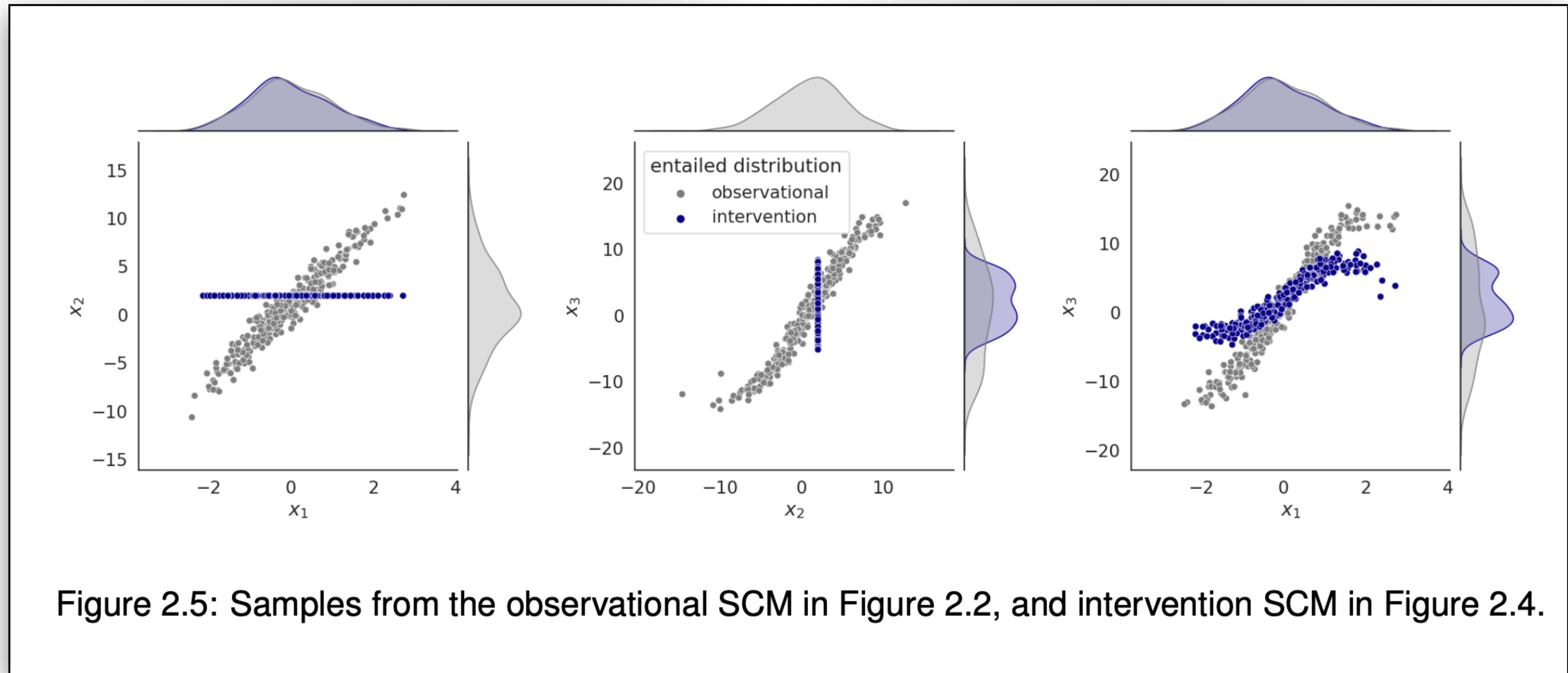


$$\begin{cases} x_1 := \varepsilon_1 \\ x_2 := 2 \\ x_3 := 5 \sin(x_1) + x_2 + \varepsilon_3 \end{cases}, \quad \varepsilon_1, \varepsilon_3 \sim \mathcal{N}(0, 1)$$

$$p_{\tilde{\mathcal{M}}}(X) = p_{\mathcal{M}; do(x_2:=2)}(X), \quad I = \{2\}$$

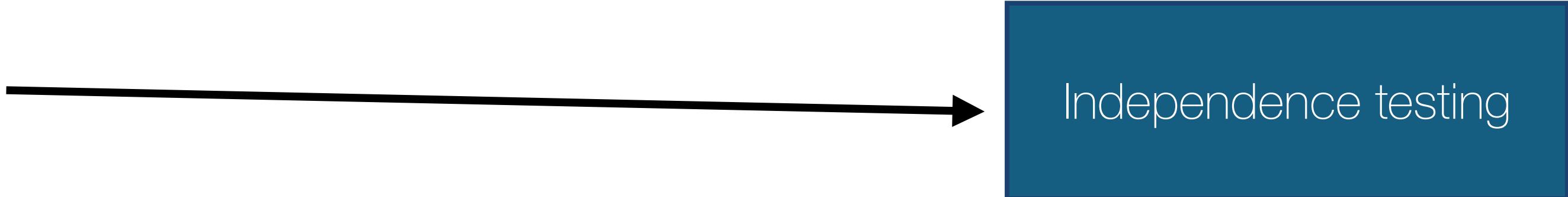
$$p_{\tilde{\mathcal{M}}}(X) = p_{\mathcal{M}; do(x_2:=2)}(X), \quad I = \{2\}$$

# Structural causal models



# Causal Discovery Methods

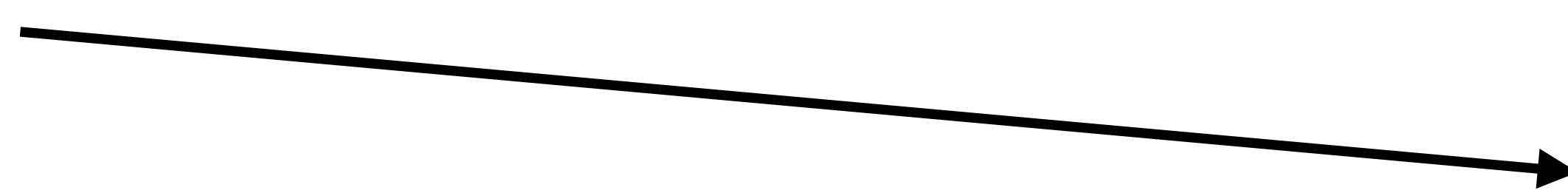
- Constrained-Based
- Score based
- Hybrid
- Others



Independence testing

# Causal Discovery Methods

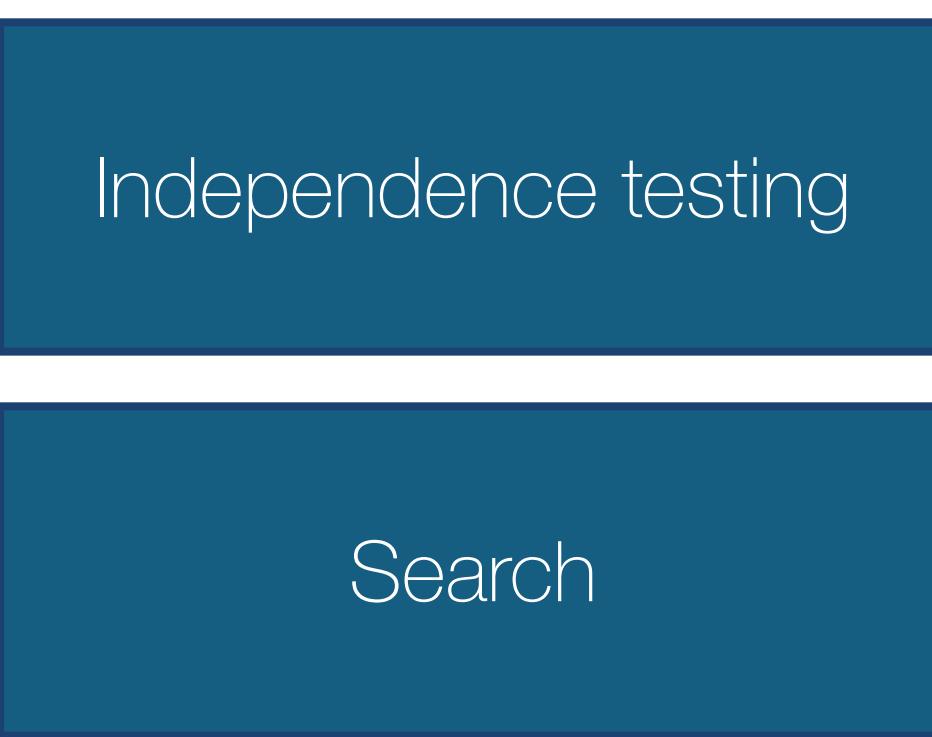
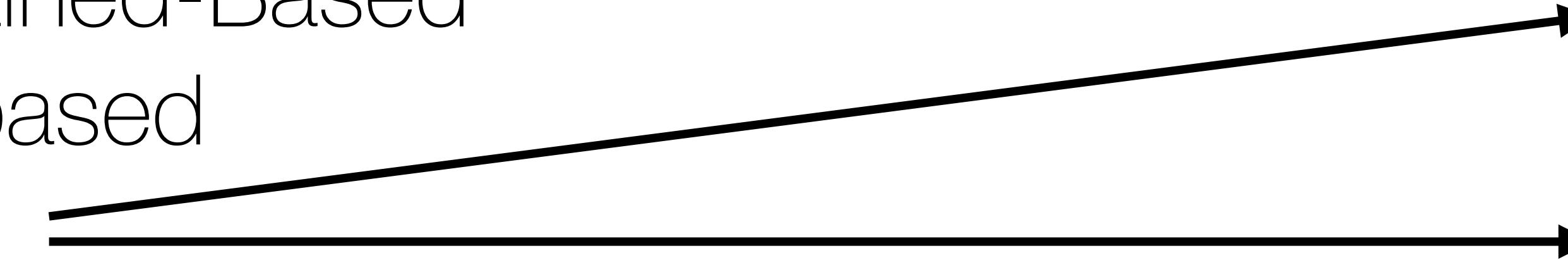
- Constrained-Based
- Score based
- Hybrid
- Others



Search

# Causal Discovery Methods

- Constrained-Based
- Score based
- Hybrid
- Others

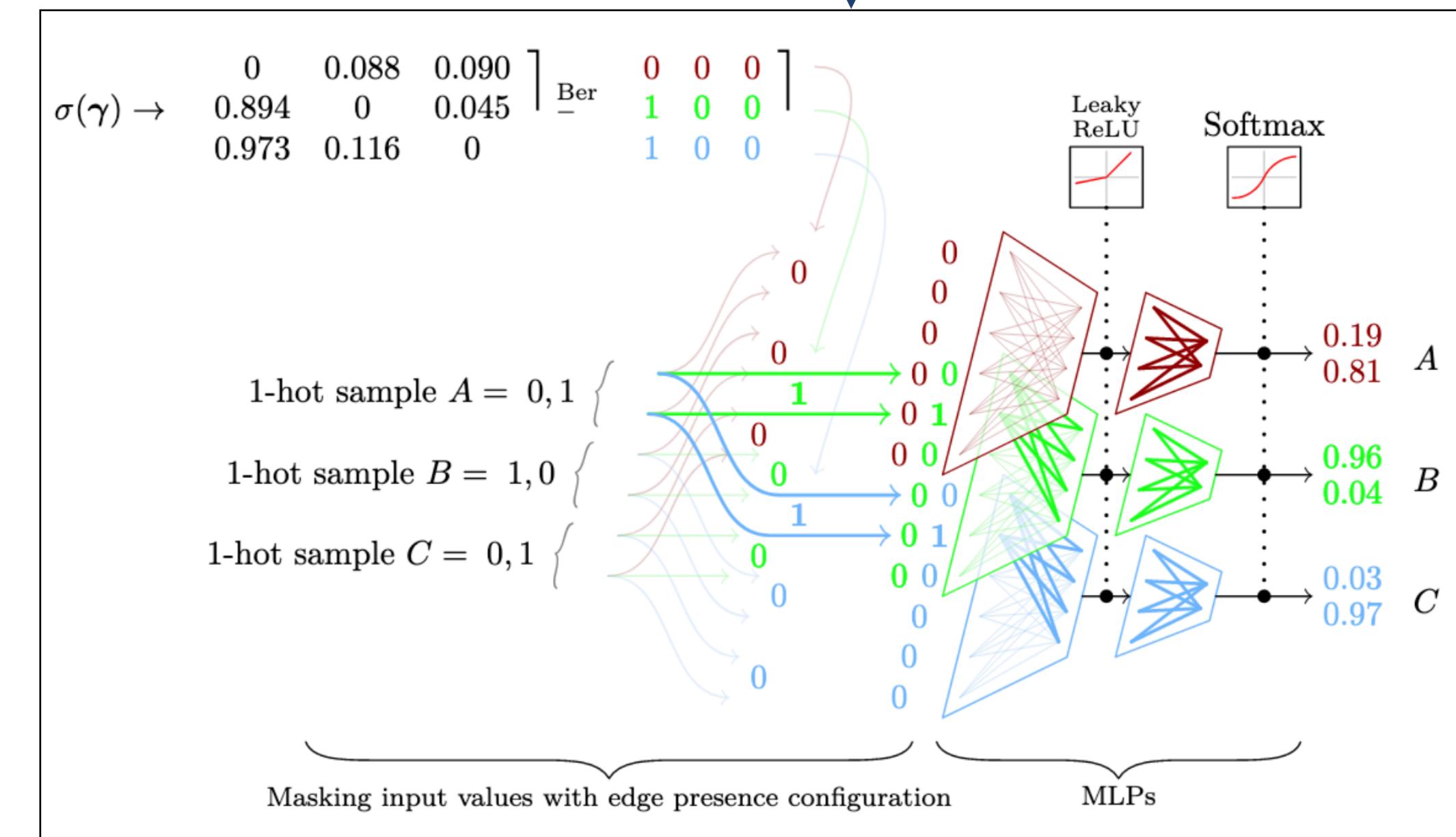


# Score based methods with deep learning

$$\Lambda^* = \arg \max_{\Lambda} S(\Lambda)$$

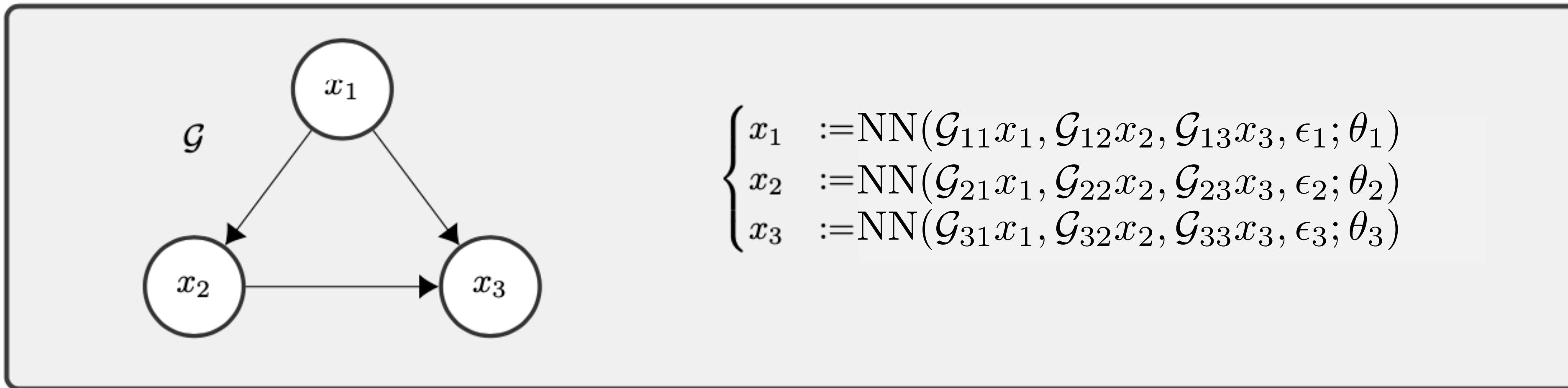
$$S(\Lambda) = \max_{\theta} \mathbb{E}_{\mathcal{G} \sim \text{DAG}(\mathcal{G}; \Lambda)} [\log p(\mathcal{D} | \mathcal{G}, \theta) + \log p(\mathcal{G})]$$

( most of the time )  
approximated w/ fully  
factorized Bernoulli



Penalize dense  
graphs/enforce  
acyclicity

# Score based methods with deep learning



Similar to doing LASSO for all of the variables at once w/ acyclicity constraints.

# Interventions and Causal Discovery

# Causal discovery with intervention data

$$\Lambda^*, \tilde{\mathcal{M}}^* = \arg \max_{\Lambda, \tilde{\mathcal{M}}} S(\Lambda, \tilde{\mathcal{M}})$$

$$S(\Lambda, \tilde{\mathcal{M}}) = \max_{\theta, \phi} \mathbb{E}_{\mathcal{G} \sim \text{DAG}(\mathcal{G}; \Lambda)} \left[ \mathbb{E}_{k \sim p(k)} \left[ \log p(\mathcal{D}^k | \mathcal{G}, \theta^k, \tilde{\mathcal{M}}^{(k)}) + \log p(\mathcal{G}) \right] \right]$$

Intervention  
specific  
parameters

↓  
Intervention  
variables

# Differentiable Causal Discovery Under Latent Interventions



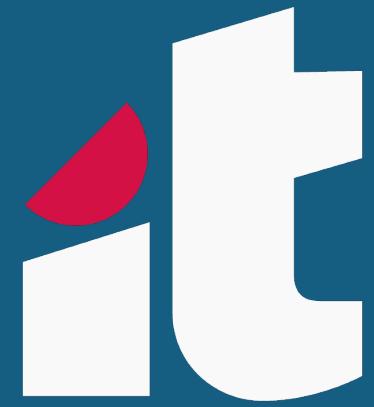
Gonçalo R. A. Faria



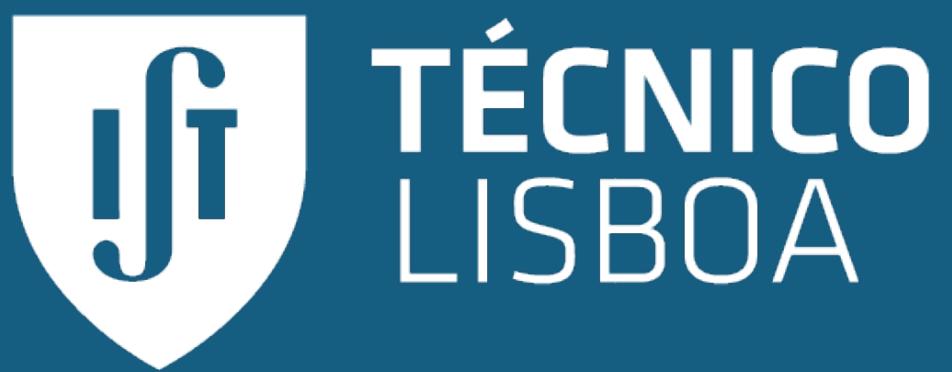
André F. T. Martins



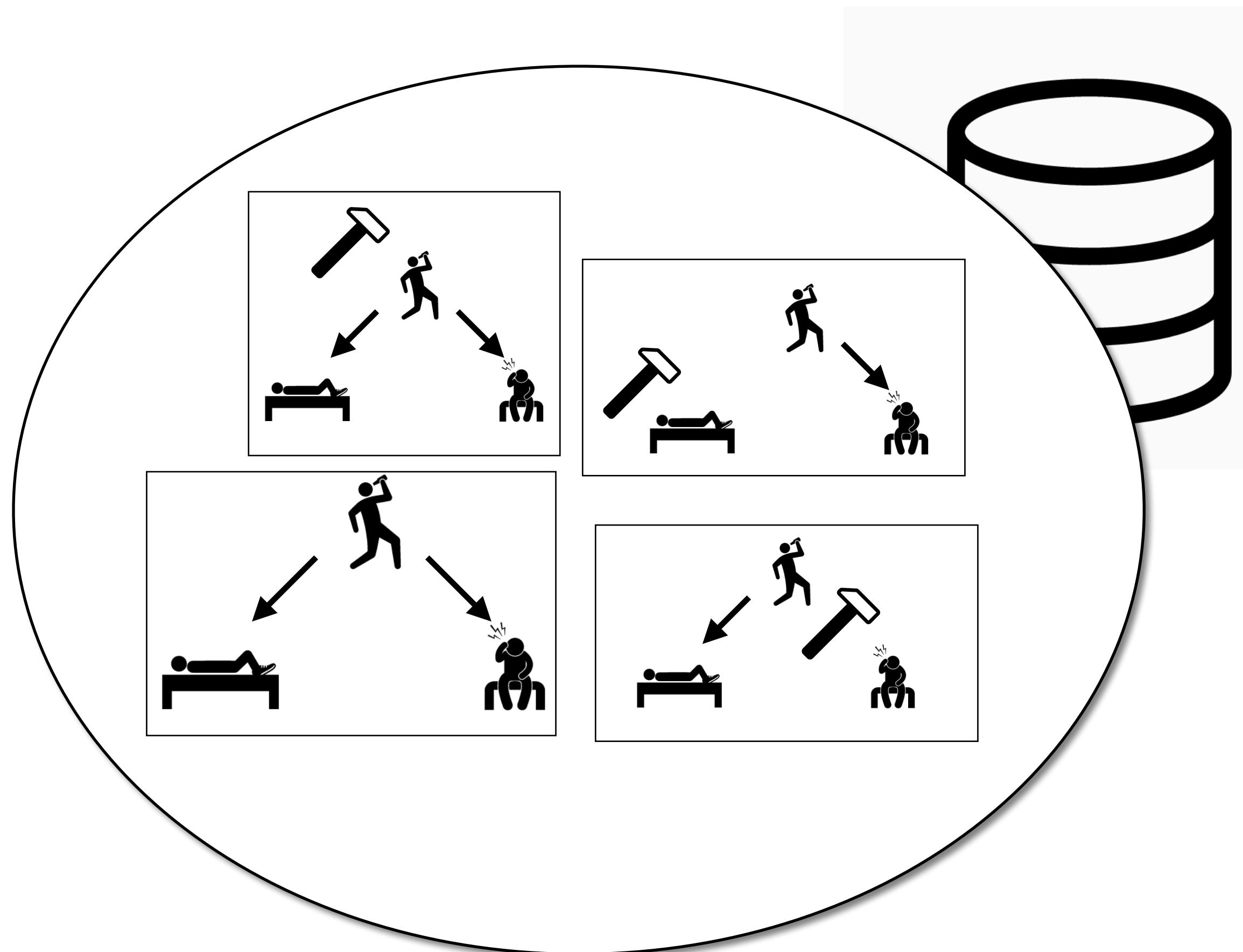
Mário A. T. Figueiredo



instituto de  
telecomunicações



# Latent interventions

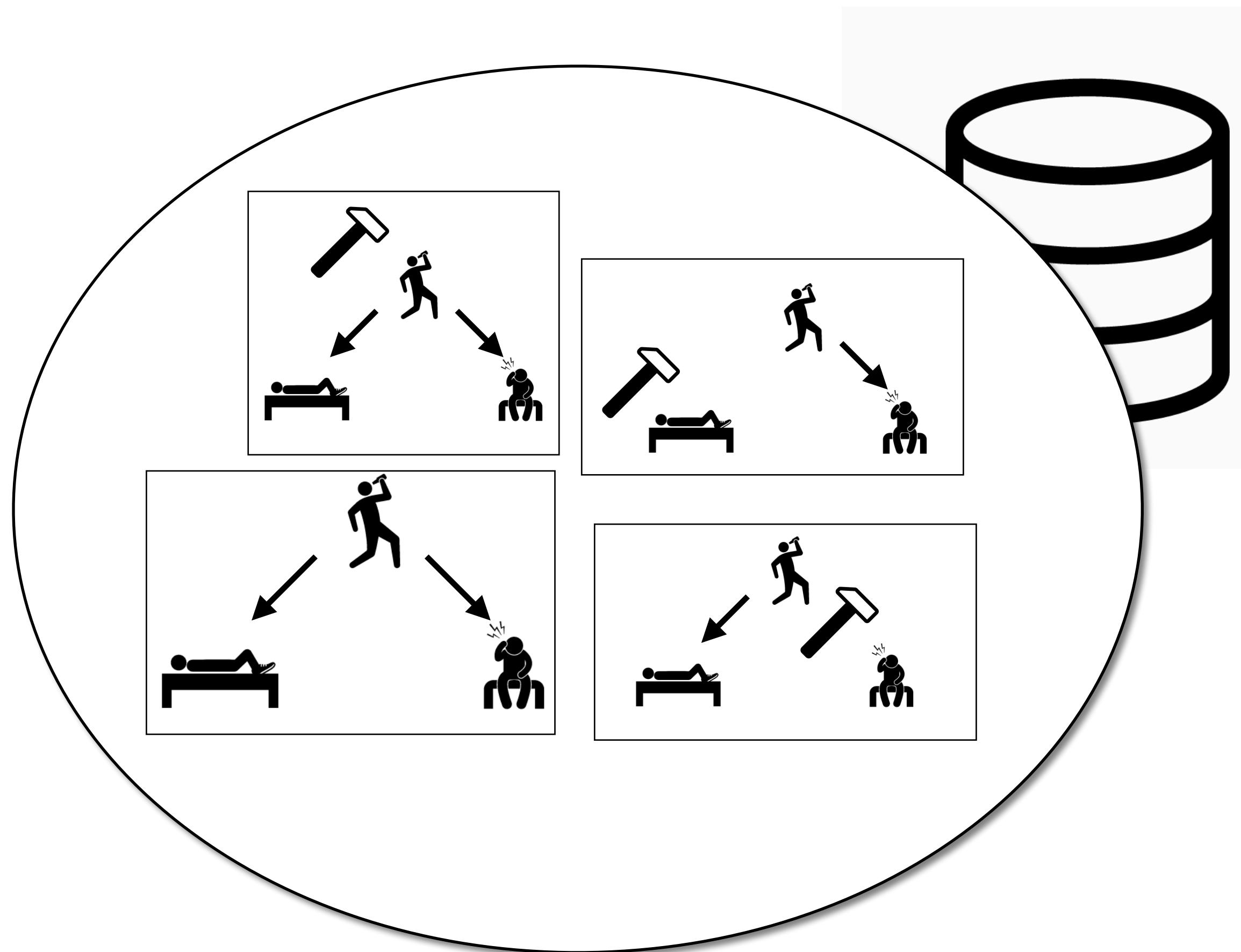


Mixture of experimental regimes.

**For each sample:**

- do not know correspondence to intervention regime;

# Latent interventions



Mixture of experimental regimes.

**For each sample:**

- do not know correspondence to intervention regime;

**For each intervention:**

- do not know experimental conditions

# Can we recover latent interventions?

# Intervention recovery

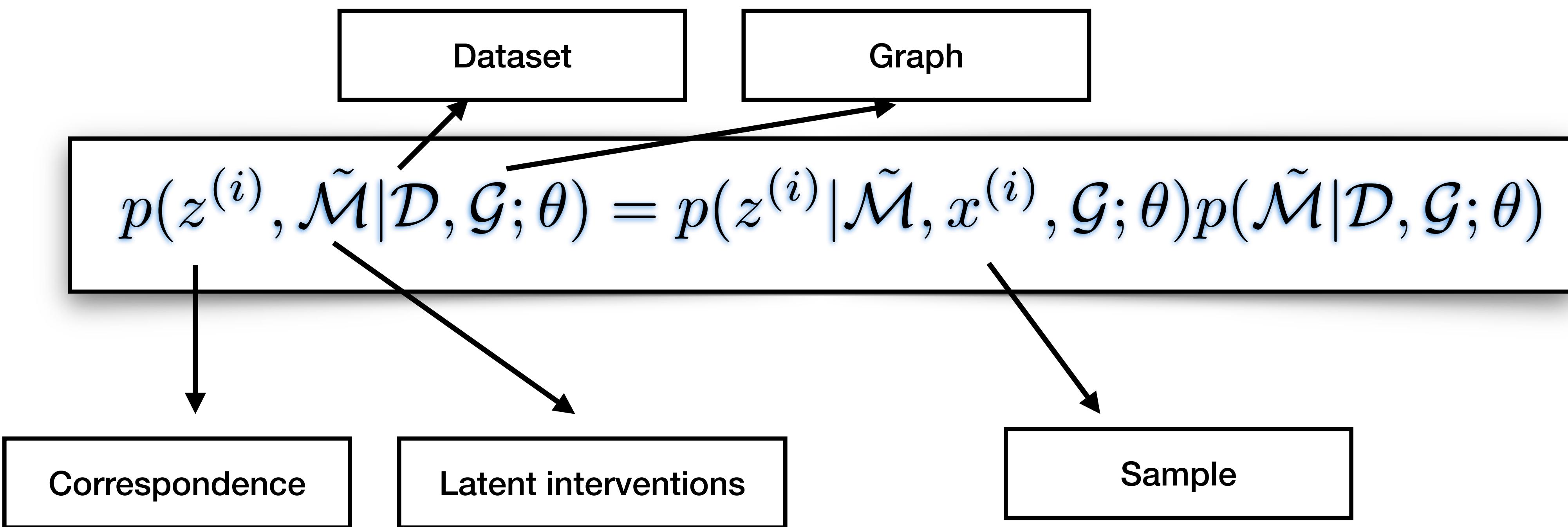
- Given the causal graph.
- Recover interventions and correspondences.
- Propose joint distribution

Infinite mixture of  
intervention SCMs

+

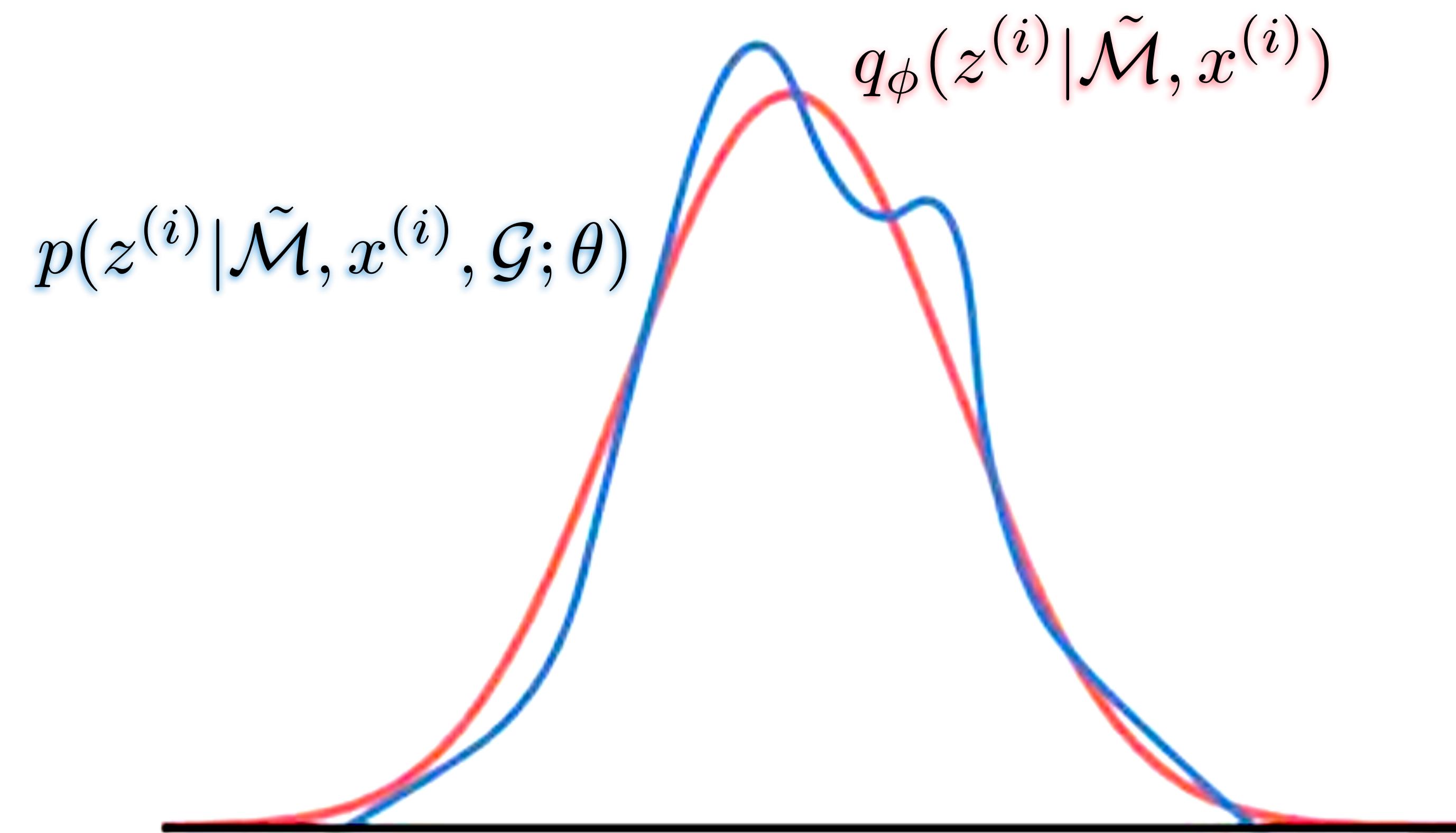
Prior distribution for the  
interventions  
( Dirichlet process)

# Approximate posterior inference



# Approximate posterior inference

- Variational inference

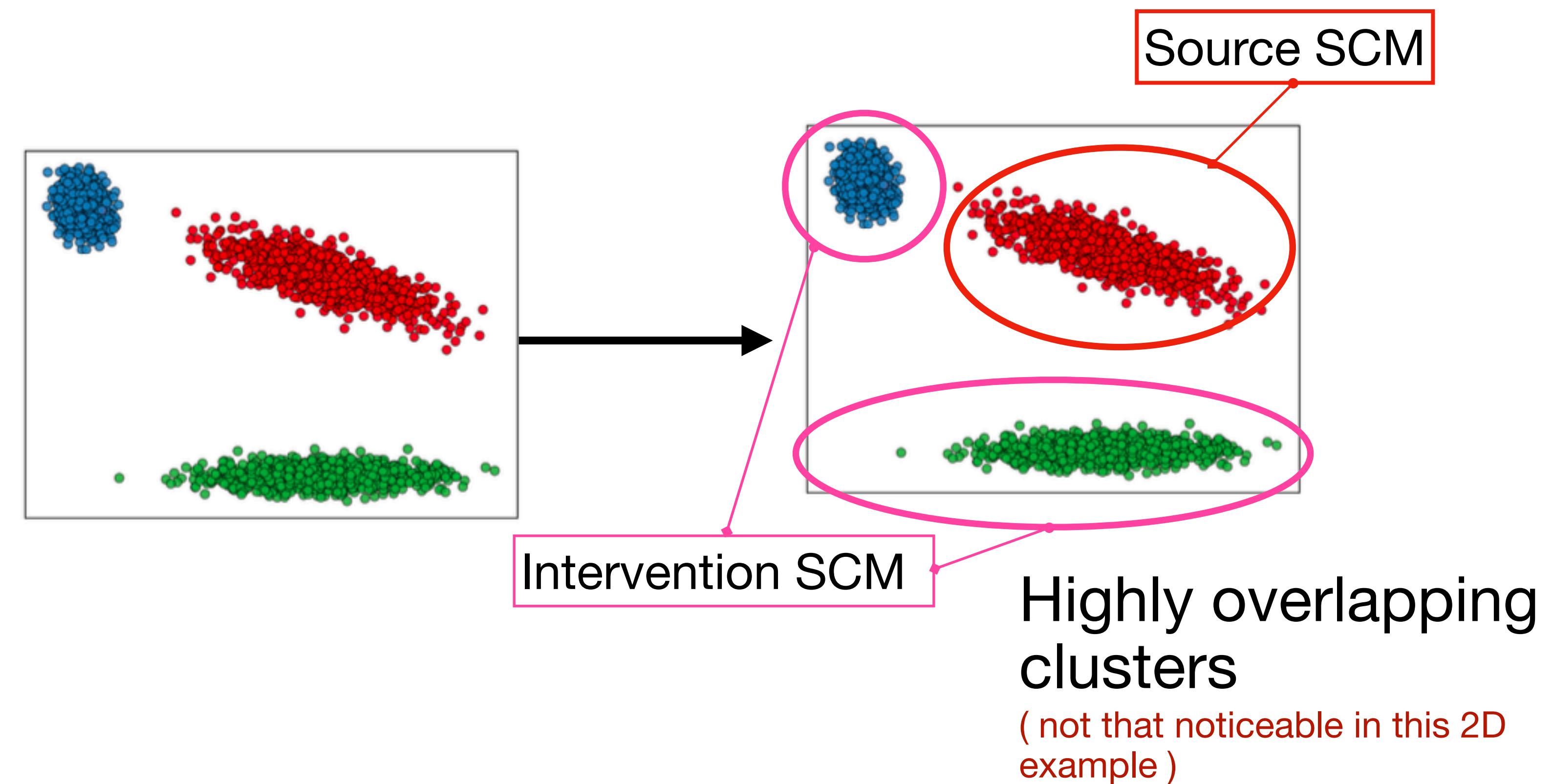


# Approximate posterior inference

- Variational inference

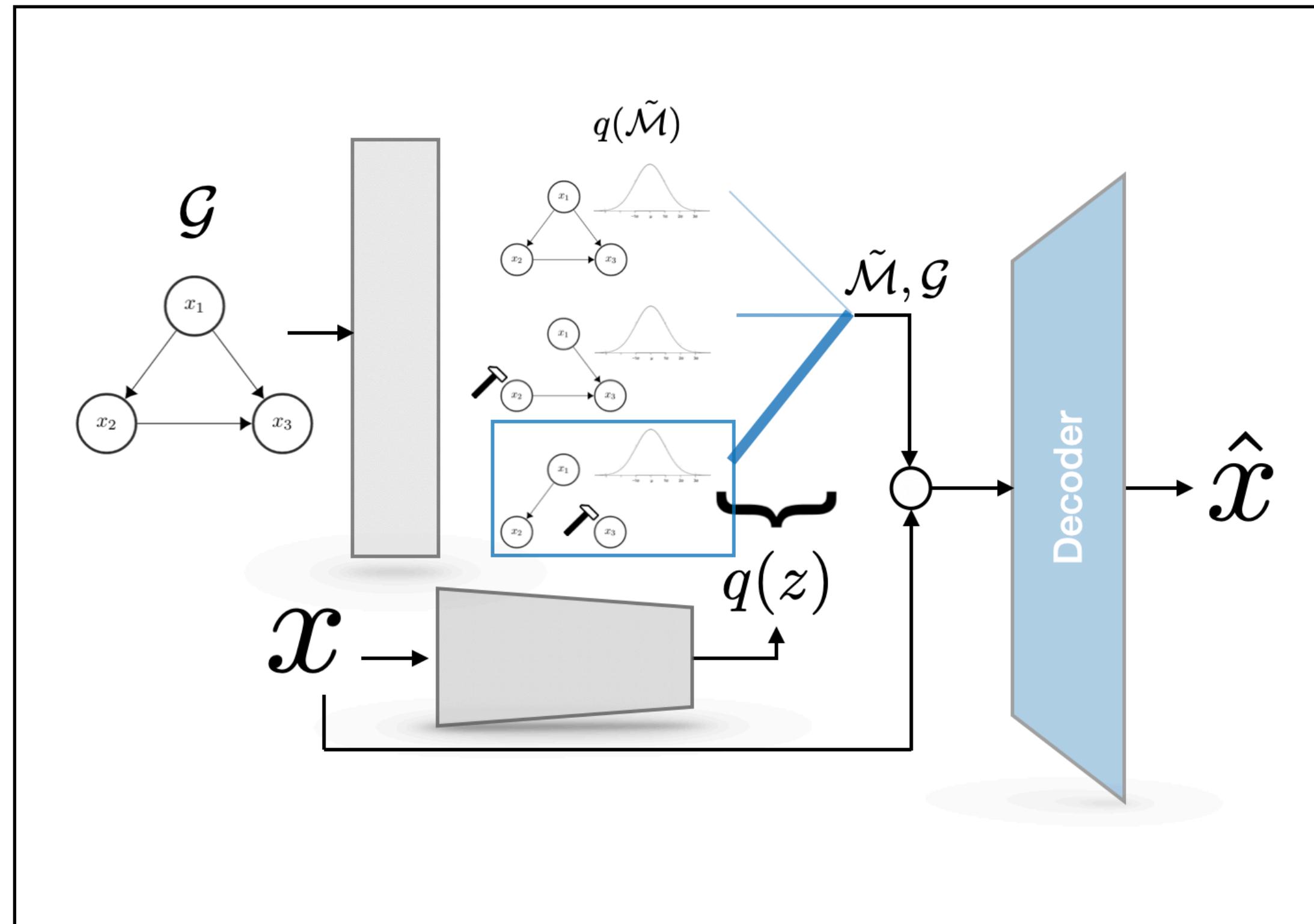
Similar to a clustering problem

( each cluster as few degrees of freedom )



# Approximate posterior inference

$$\max_{\theta, \phi} \sum_{i=1}^N \text{ELBO}_{q(z^{(i)}, \tilde{\mathcal{M}}; \phi)}(x^{(i)}, \mathcal{G}; \theta)$$



# Experiments with single node interventions

Model Type	<b>d = 10</b>		<b>d = 20</b>	
	ER1	ER4	ER1	ER4
<i>Stochastic Interventions:</i>				
Linear Gaussian	$0.9816 \pm 0.019$	$0.9492 \pm 0.007$	$0.9815 \pm 0.005$	$0.9785 \pm 0.006$
Non-Linear Gaussian	$0.9643 \pm 0.015$	$0.9242 \pm 0.006$	$0.9655 \pm 0.005$	$0.9431 \pm 0.003$
Non-Linear Non-Gaussian	$0.8314 \pm 0.062$	$0.7812 \pm 0.1074$	$0.8929 \pm 0.027$	$0.8412 \pm 0.087$
<i>Imperfect Interventions:</i>				
Linear Gaussian	$0.9262 \pm 0.030$	$0.8802 \pm 0.040$	$0.8072 \pm 0.098$	$0.7819 \pm 0.054$
Non-Linear Gaussian	$0.9541 \pm 0.012$	$0.9205 \pm 0.005$	$0.9670 \pm 0.007$	$0.9434 \pm 0.004$
Non-Linear Non-Gaussian	$0.9103 \pm 0.004$	$0.9090 \pm 0.004$	$0.9345 \pm 0.005$	$0.9196 \pm 0.005$
<i>Atomic Interventions:</i>				
Linear Gaussian	$0.9141 \pm 0.0018$	$0.9148 \pm 0.0021$	$0.9507 \pm 0.0059$	$0.9461 \pm 0.0030$
Non-Linear Gaussian	$0.9138 \pm 0.0053$	$0.9195 \pm 0.0041$	$0.9506 \pm 0.0056$	$0.9435 \pm 0.0032$
Non-Linear Non-Gaussian	$0.9100 \pm 0.0056$	$0.9034 \pm 0.0115$	$0.9314 \pm 0.0104$	$0.9306 \pm 0.0041$

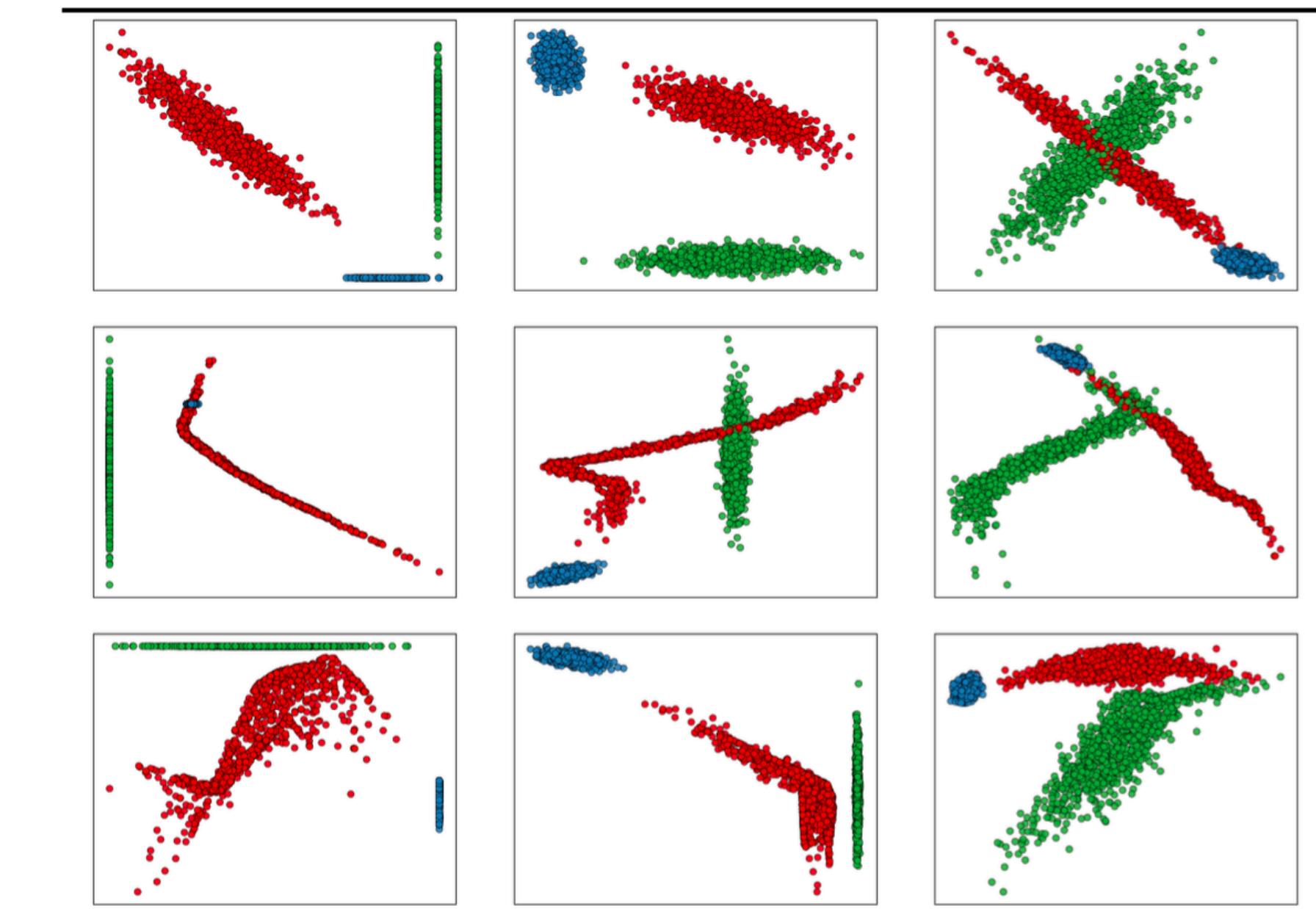


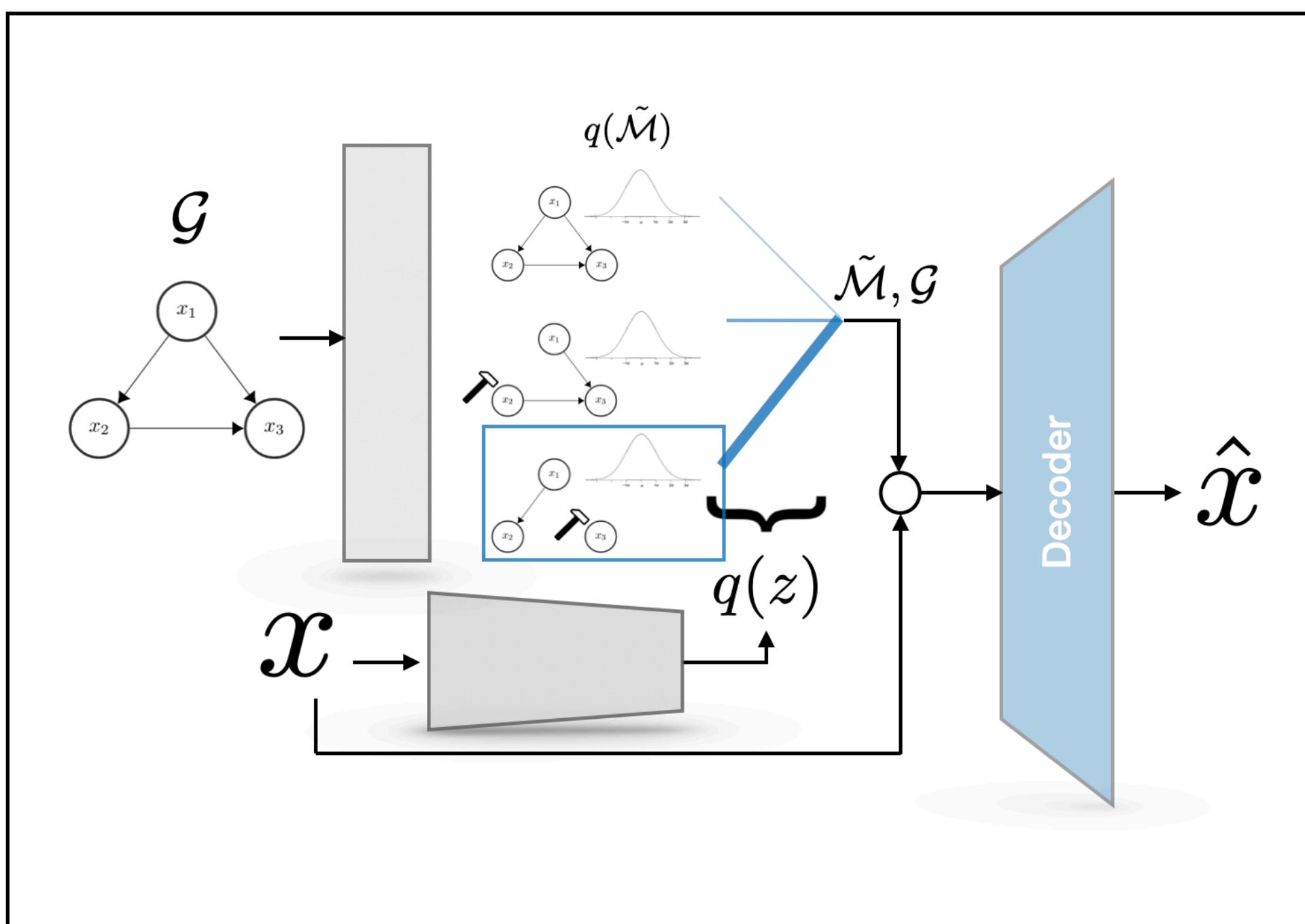
Table 3.2: Results for the single-node intervention experiment. The metric used is the average rand index.

# Can we recover the causal graph?

# Score based method

$$\Lambda^* = \arg \max_{\Lambda} S(\Lambda)$$

$$S(\Lambda) = \max_{\theta, \phi} \mathbb{E}_{\mathcal{G} \sim \text{DAG}(\mathcal{G}; \Lambda)} \left[ \sum_{i=1}^N \text{ELBO}_{q(z^{(i)}, \tilde{\mathcal{M}}; \phi)}(x^{(i)}, \mathcal{G}; \theta) + \log p(\mathcal{G}) \right]$$



# Model variants

- “latent” (new)

**For each sample:**

- do not know correspondence to intervention regime;

**For each intervention:**

- do not know experimental conditions

# Model variants

- “latent” (new)
- “unknown”

**For each sample:**

- ~~do~~ not know correspondence to intervention regime;

**For each intervention:**

- do not know experimental conditions

# Model variants

- “latent” (new)
- “unknown”
- **“known”**

**For each sample:**

- ~~do~~ not know correspondence to intervention regime;

**For each intervention:**

- ~~do~~ not know experimental conditions

# Model variants

- “latent” (new) assume there is no intervention data.
- “unknown”
- “known”
- **“observational”**

# Model variants

- “latent” (new)
- “unknown”
- “known”
- “observational”
- **“semi-supervised”** (new)

# Model variants

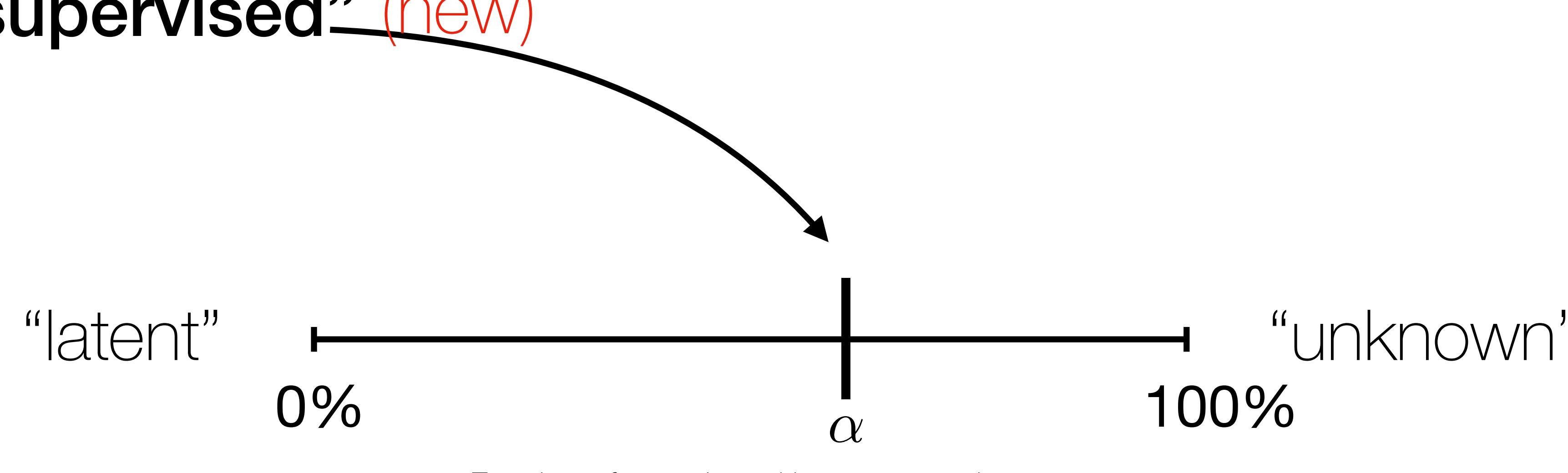
- “latent” (new)
- “unknown”
- “known”
- “observational”
- “semi-supervised” (new)

**For each sample:**

- (fraction) know correspondence to intervention regime;

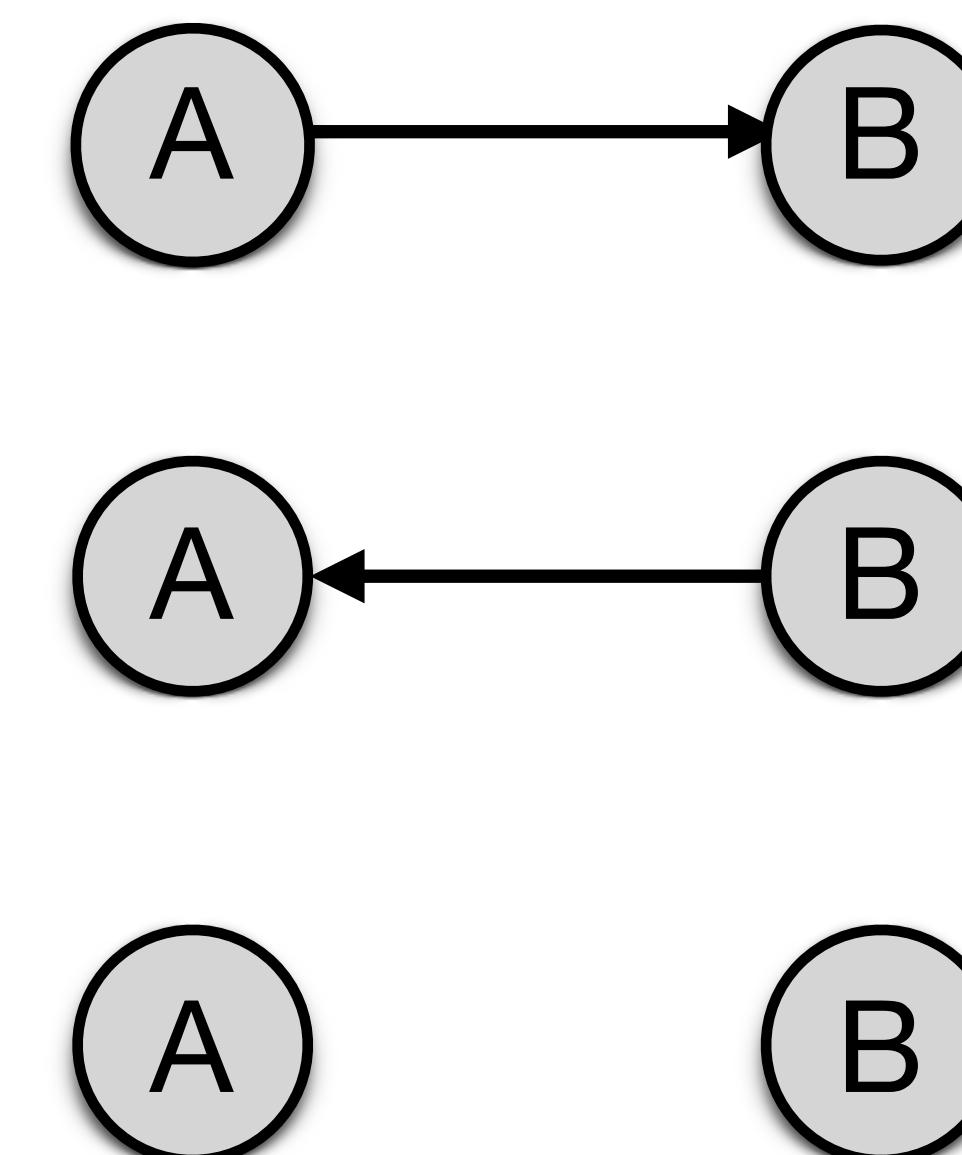
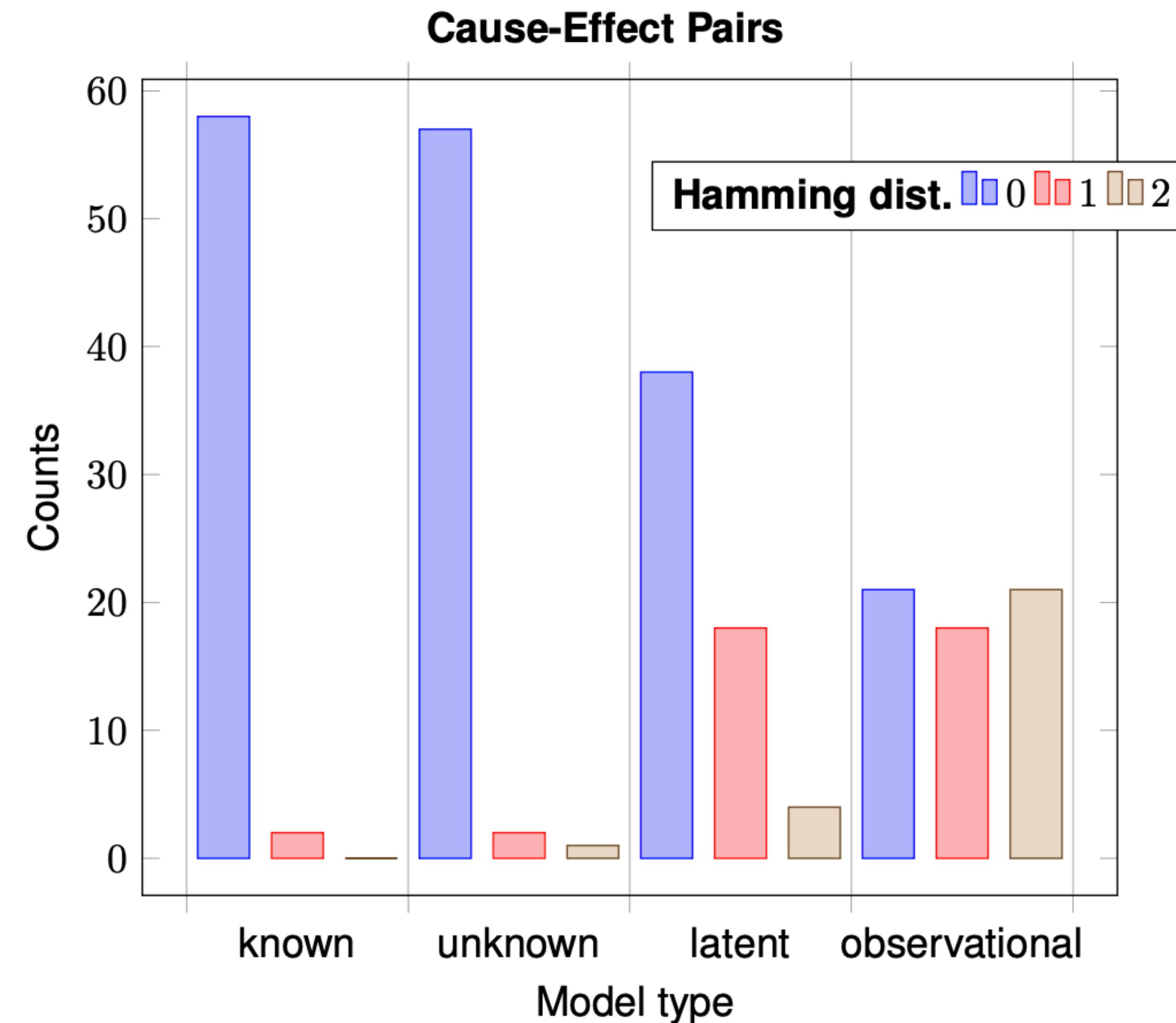
**For each intervention:**

- do not know experimental conditions



# Experiments

# Experiments with cause effect pairs



# Experiments on synthetic data

Model Type	$e$	latent	unknown	known	observational
<i>Stochastic Interventions:</i>					
Linear Gaussian		$5.9 \pm 6.2$	$3.4 \pm 3.2$	$0.5 \pm 1.3$	$10.3 \pm 7.8$
Non-Linear Gaussian	1	$12.2 \pm 3.9$	$10.3 \pm 2.5$	$7.0 \pm 3.6$	$13.7 \pm 3.8$
Non-Linear Non-Gaussian		$8.7 \pm 6.6$	$8.0 \pm 2.7$	$6.6 \pm 2.2$	$11.3 \pm 5.0$
Linear Gaussian		$27.2 \pm 6.2$	$24.1 \pm 5.8$	$15.6 \pm 6.0$	$39.6 \pm 5.0$
Non-Linear Gaussian	4	$35.8 \pm 3.8$	$30.3 \pm 5.3$	$27.7 \pm 4.3$	$37.5 \pm 5.2$
Non-Linear Non-Gaussian		$36.1 \pm 4.4$	$35.5 \pm 8.1$	$31.5 \pm 5.6$	$40.2 \pm 6.9$
<i>Imperfect Interventions:</i>					
Linear Gaussian		$5.8 \pm 4.2$	$6.2 \pm 3.06$	$4.7 \pm 3.6$	$10.4 \pm 2.9$
Non-Linear Gaussian	1	$9.3 \pm 2.4$	$8.9 \pm 2.5$	$7.8 \pm 3.9$	$10.5 \pm 2.8$
Non-Linear Non-Gaussian		$8.8 \pm 3.0$	$9.1 \pm 3.5$	$7.9 \pm 1.4$	$11.5 \pm 5.4$
Linear Gaussian		$35.9 \pm 8.3$	$29.7 \pm 5.6$	$17.7 \pm 7.9$	$39.1 \pm 9.1$
Non-Linear Gaussian	4	$32.1 \pm 6.0$	$32.6 \pm 5.8$	$32.8 \pm 5.4$	$39.8 \pm 9.3$
Non-Linear Non-Gaussian		$30.4 \pm 12.2$	$30.2 \pm 11.2$	$25.8 \pm 3.9$	$36.7 \pm 9.8$

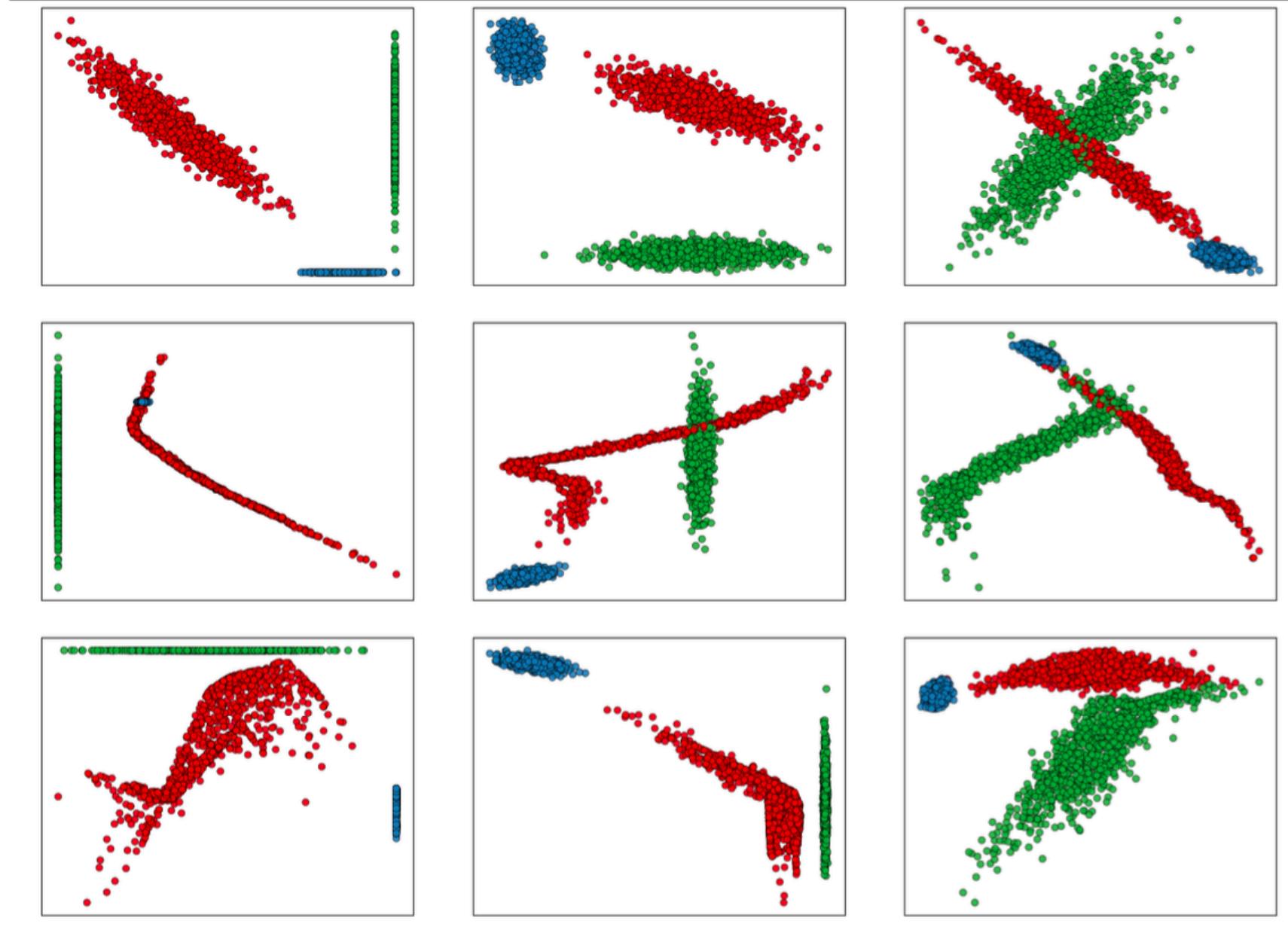


Table 4.1: Hamming distances on synthetic 10 variable SCMs.

# Experiments on synthetic data

Model Type	$e$	latent	unknown	known	observational
<i>Stochastic Interventions:</i>					
Linear Gaussian		$5.9 \pm 6.2$	$3.4 \pm 3.2$	$0.5 \pm 1.3$	$10.3 \pm 7.8$
Non-Linear Gaussian	1	$12.2 \pm 3.9$	$10.3 \pm 2.5$	$7.0 \pm 3.6$	$13.7 \pm 3.8$
Non-Linear Non-Gaussian		$8.7 \pm 6.6$	$8.0 \pm 2.7$	$6.6 \pm 2.2$	$11.3 \pm 5.0$
Linear Gaussian		$27.2 \pm 6.2$	$24.1 \pm 5.8$	$15.6 \pm 6.0$	$39.6 \pm 5.0$
Non-Linear Gaussian	4	$35.8 \pm 3.8$	$30.3 \pm 5.3$	$27.7 \pm 4.3$	$37.5 \pm 5.2$
Non-Linear Non-Gaussian		$36.1 \pm 4.4$	$35.5 \pm 8.1$	$31.5 \pm 5.6$	$40.2 \pm 6.9$
<i>Imperfect Interventions:</i>					
Linear Gaussian		$5.8 \pm 4.2$	$6.2 \pm 3.06$	$4.7 \pm 3.6$	$10.4 \pm 2.9$
Non-Linear Gaussian	1	$9.3 \pm 2.4$	$8.9 \pm 2.5$	$7.8 \pm 3.9$	$10.5 \pm 2.8$
Non-Linear Non-Gaussian		$8.8 \pm 3.0$	$9.1 \pm 3.5$	$7.9 \pm 1.4$	$11.5 \pm 5.4$
Linear Gaussian		$35.9 \pm 8.3$	$29.7 \pm 5.6$	$17.7 \pm 7.9$	$39.1 \pm 9.1$
Non-Linear Gaussian	4	$32.1 \pm 6.0$	$32.6 \pm 5.8$	$32.8 \pm 5.4$	$39.8 \pm 9.3$
Non-Linear Non-Gaussian		$30.4 \pm 12.2$	$30.2 \pm 11.2$	$25.8 \pm 3.9$	$36.7 \pm 9.8$

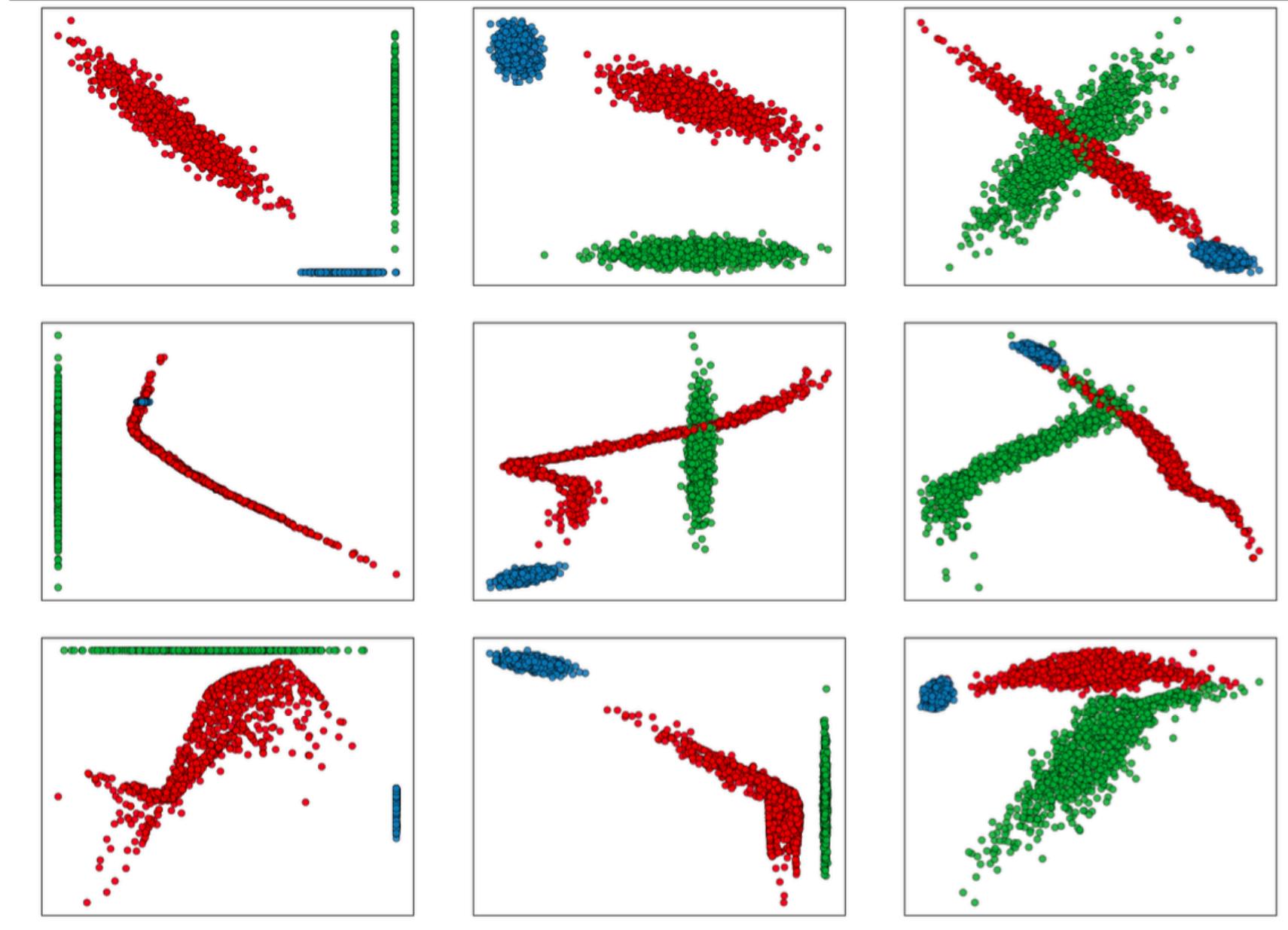


Table 4.1: Hamming distances on synthetic 10 variable SCMs.

# Experiments on synthetic data

Model Type	$e$	latent	unknown	known	observational
<i>Stochastic Interventions:</i>					
Linear Gaussian		$5.9 \pm 6.2$	$3.4 \pm 3.2$	$0.5 \pm 1.3$	$10.3 \pm 7.8$
Non-Linear Gaussian	1	$12.2 \pm 3.9$	$10.3 \pm 2.5$	$7.0 \pm 3.6$	$13.7 \pm 3.8$
Non-Linear Non-Gaussian		$8.7 \pm 6.6$	$8.0 \pm 2.7$	$6.6 \pm 2.2$	$11.3 \pm 5.0$
Linear Gaussian		$27.2 \pm 6.2$	$24.1 \pm 5.8$	$15.6 \pm 6.0$	$39.6 \pm 5.0$
Non-Linear Gaussian	4	$35.8 \pm 3.8$	$30.3 \pm 5.3$	$27.7 \pm 4.3$	$37.5 \pm 5.2$
Non-Linear Non-Gaussian		$36.1 \pm 4.4$	$35.5 \pm 8.1$	$31.5 \pm 5.6$	$40.2 \pm 6.9$
<i>Imperfect Interventions:</i>					
Linear Gaussian		$5.8 \pm 4.2$	$6.2 \pm 3.06$	$4.7 \pm 3.6$	$10.4 \pm 2.9$
Non-Linear Gaussian	1	$9.3 \pm 2.4$	$8.9 \pm 2.5$	$7.8 \pm 3.9$	$10.5 \pm 2.8$
Non-Linear Non-Gaussian		$8.8 \pm 3.0$	$9.1 \pm 3.5$	$7.9 \pm 1.4$	$11.5 \pm 5.4$
Linear Gaussian		$35.9 \pm 8.3$	$29.7 \pm 5.6$	$17.7 \pm 7.9$	$39.1 \pm 9.1$
Non-Linear Gaussian	4	$32.1 \pm 6.0$	$32.6 \pm 5.8$	$32.8 \pm 5.4$	$39.8 \pm 9.3$
Non-Linear Non-Gaussian		$30.4 \pm 12.2$	$30.2 \pm 11.2$	$25.8 \pm 3.9$	$36.7 \pm 9.8$

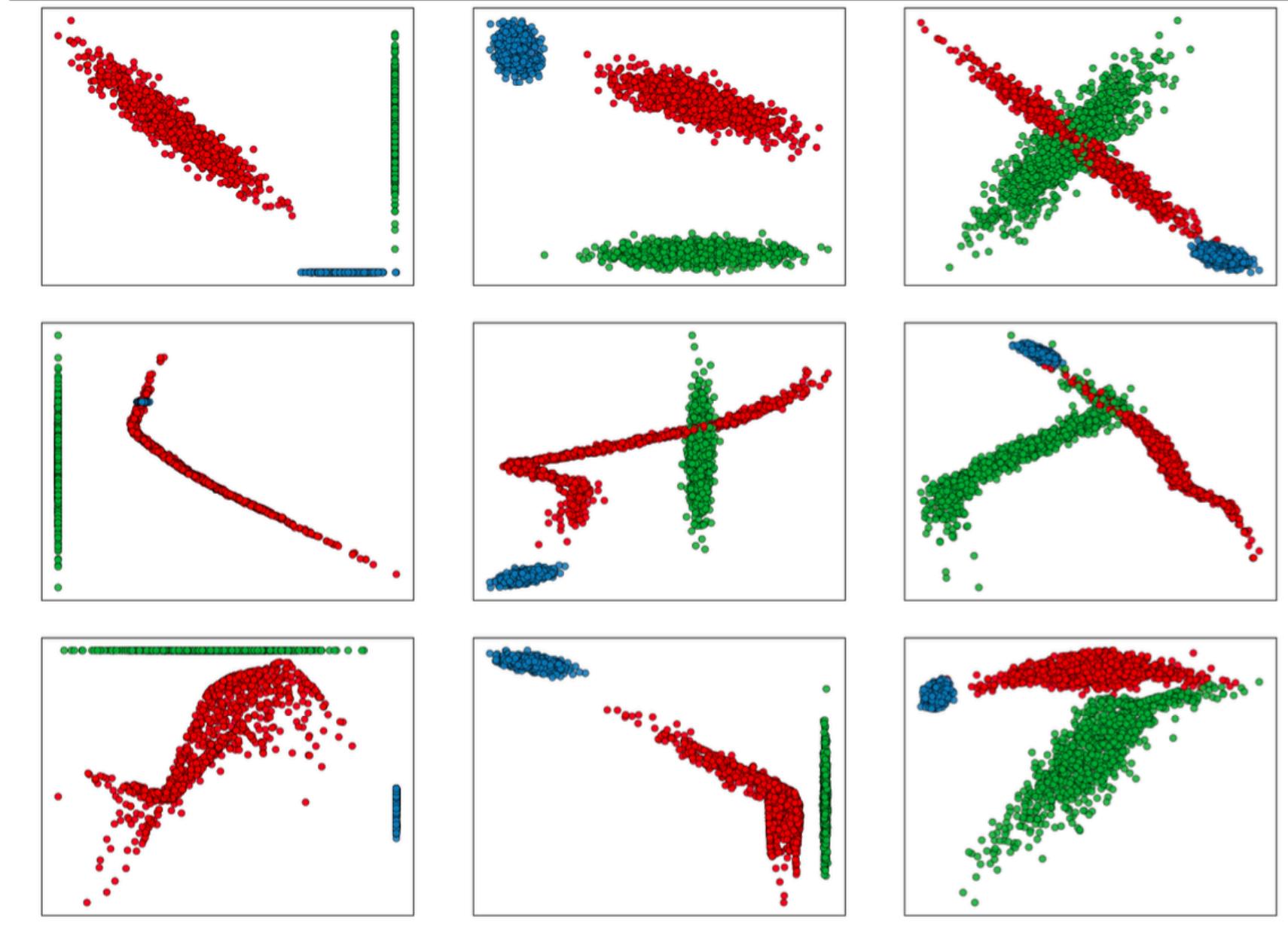
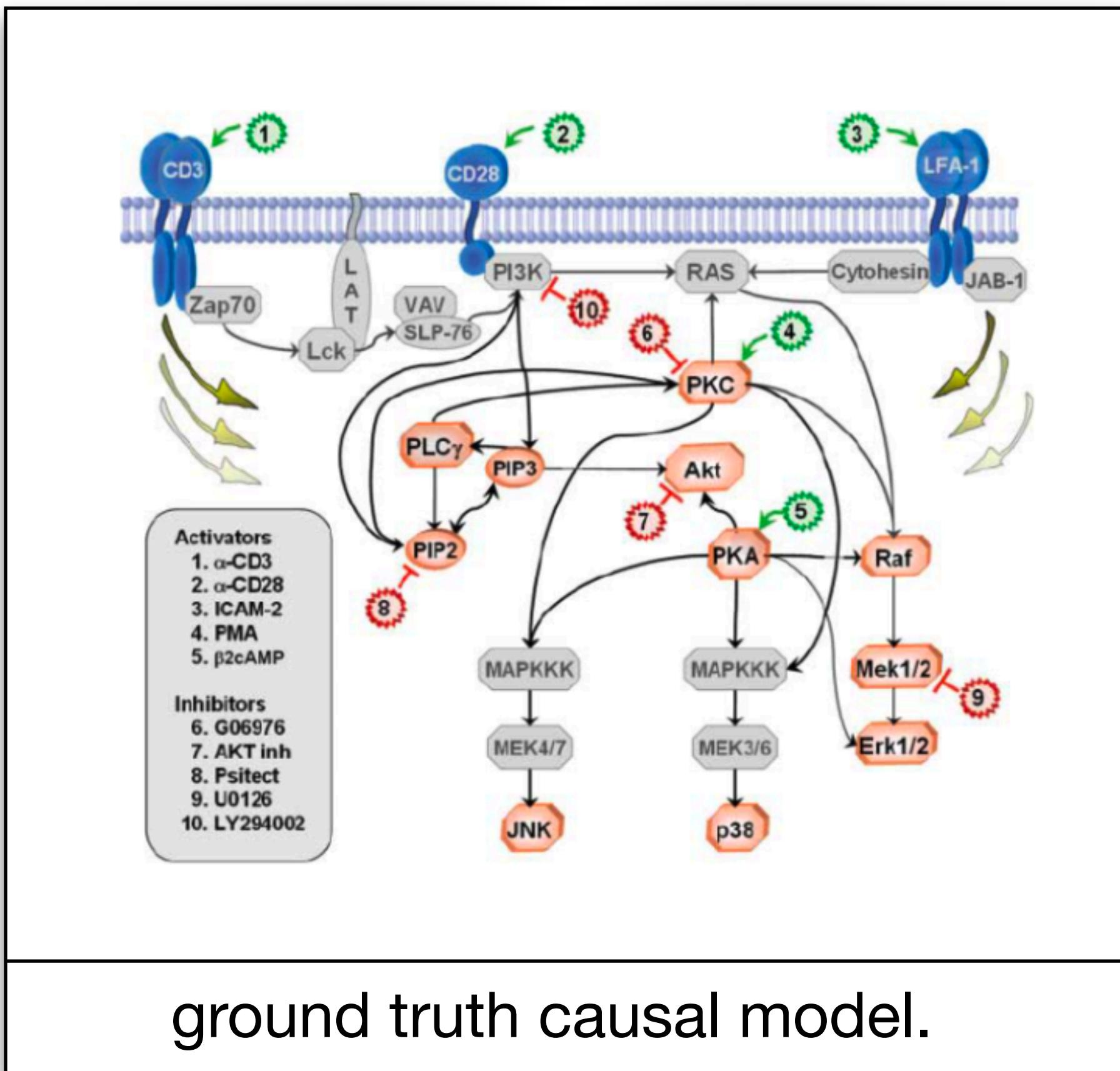


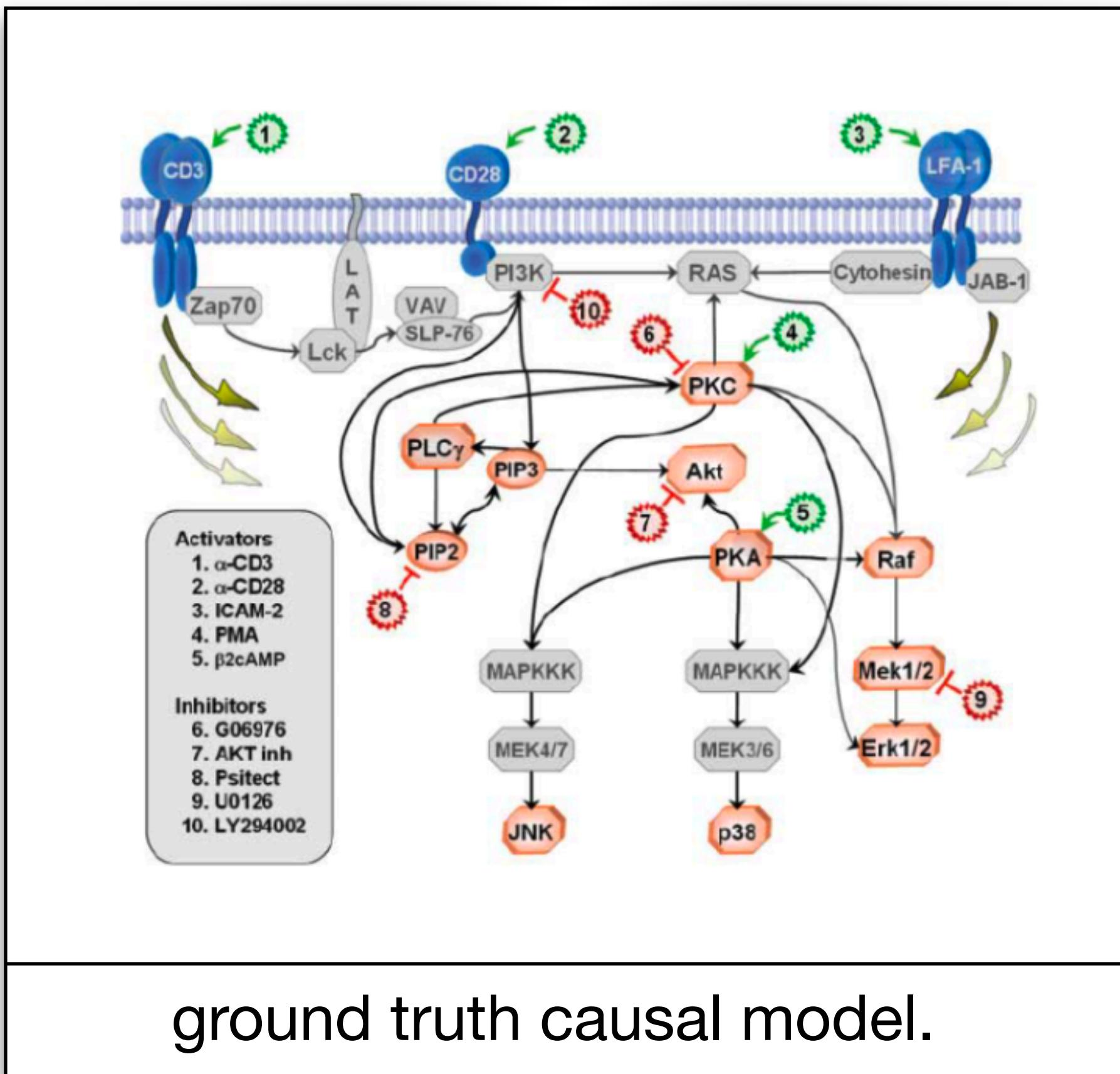
Table 4.1: Hamming distances on synthetic 10 variable SCMs.

# Experiments on real data



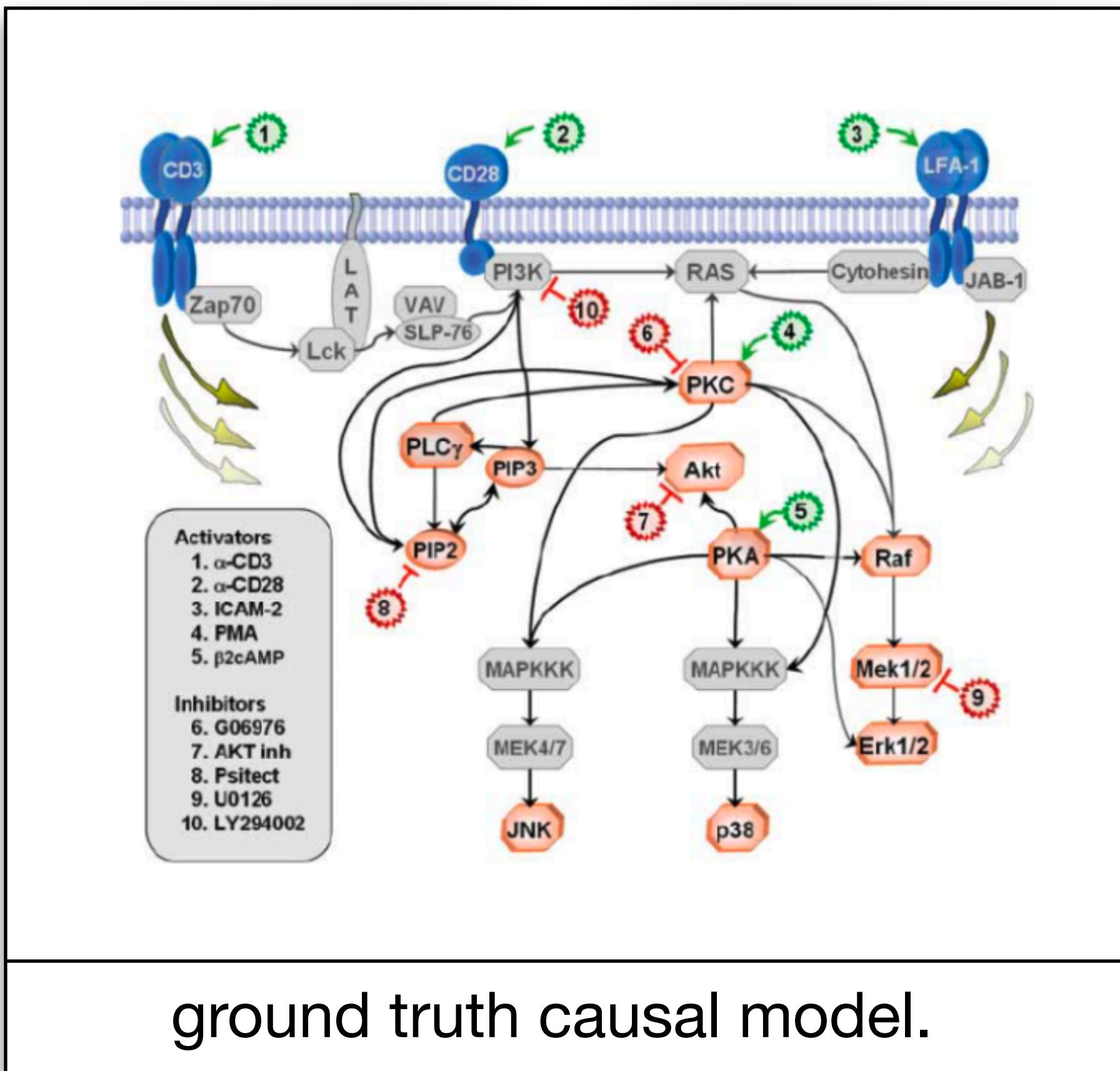
	HD	tp	fn	fp	rev	$F_1$ score
GIES (Hauser and Bühlmann, 2012)	38	10	0	41	7	0.33
CAM (Bühlmann et al., 2014)	35	12	1	30	4	0.51
IGSP (Wang et al., 2017)	18	4	6	5	7	0.42
DCDI-G (Brouillard et al., 2020)	36	6	2	25	9	0.31
DCDI-DSF (Brouillard et al., 2020)	33	6	2	22	9	0.33
FCI (Spirtes et al., 1993)	35	4	12	21	5	0.22
Imperfect Linear Gaussian (ours)	33	7	11	22	3	0.30
Imperfect Non-Linear Gaussian (ours)	19	7	11	8	0	0.42
Imperfect Normalizing Flow (ours)	30	9	9	21	1	0.38
Perfect Linear Gaussian (ours)	23	8	10	13	3	0.41
Perfect Non-Linear Gaussian (ours)	24	11	7	17	1	0.48
Perfect Normalizing Flow (ours)	23	7	11	12	2	0.38

# Experiments on real data



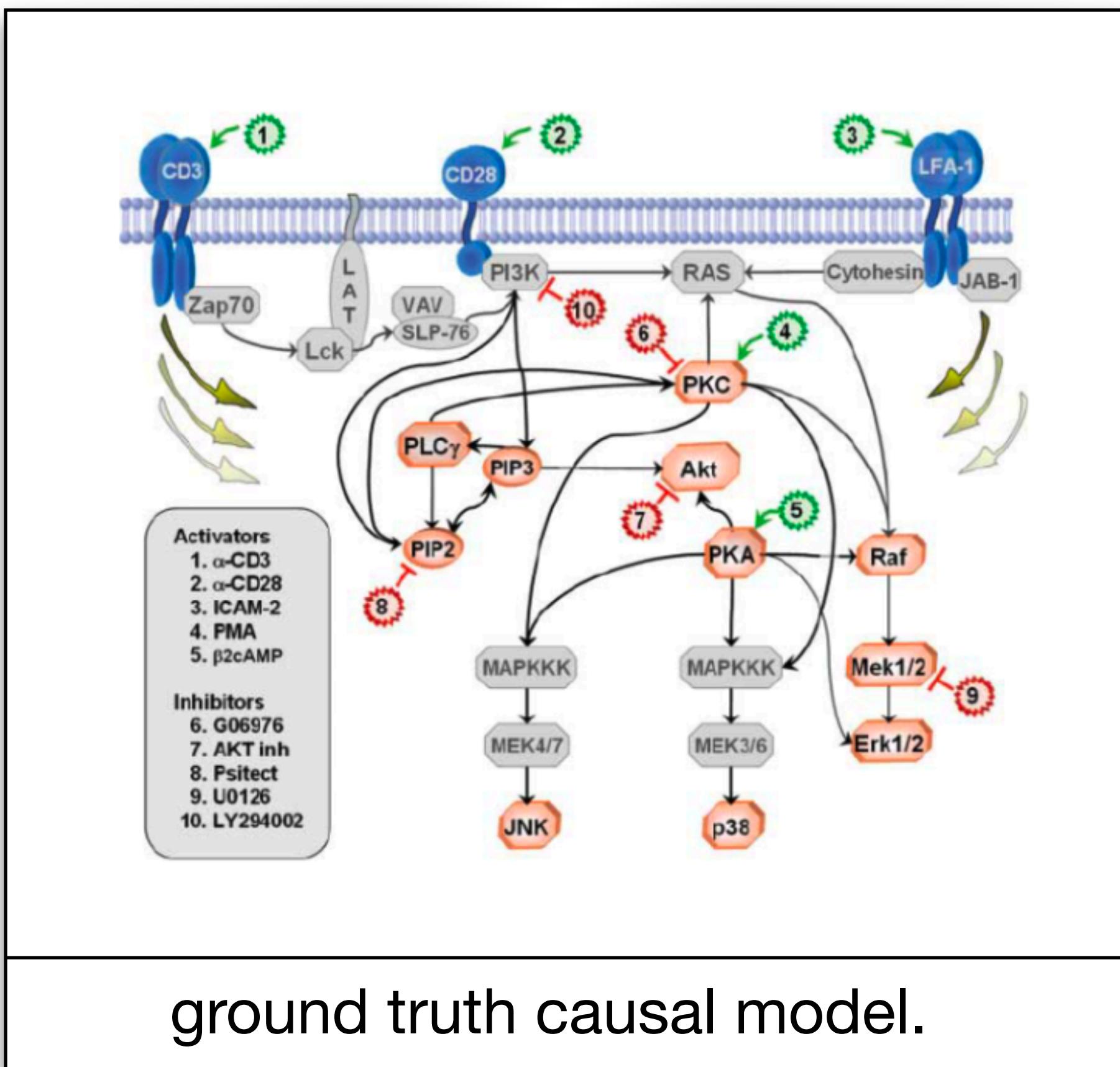
	HD	tp	fn	fp	rev	$F_1$ score
GIES (Hauser and Bühlmann, 2012)	38	10	0	41	7	0.33
CAM (Bühlmann et al., 2014)	35	12	1	30	4	0.51
IGSP (Wang et al., 2017)	18	4	6	5	7	0.42
DCDI-G (Brouillard et al., 2020)	36	6	2	25	9	0.31
DCDI-DSF (Brouillard et al., 2020)	33	6	2	22	9	0.33
FCI (Spirtes et al., 1993)	35	4	12	21	5	0.22
Imperfect Linear Gaussian (ours)	33	7	11	22	3	0.30
Imperfect Non-Linear Gaussian (ours)	19	7	11	8	0	0.42
Imperfect Normalizing Flow (ours)	30	9	9	21	1	0.38
Perfect Linear Gaussian (ours)	23	8	10	13	3	0.41
Perfect Non-Linear Gaussian (ours)	24	11	7	17	1	0.48
Perfect Normalizing Flow (ours)	23	7	11	12	2	0.38

# Experiments on real data



	HD	tp	fn	fp	rev	$F_1$ score
GIES (Hauser and Bühlmann, 2012)	38	10	0	41	7	0.33
CAM (Bühlmann et al., 2014)	35	12	1	30	4	0.51
IGSP (Wang et al., 2017)	18	4	6	5	7	0.42
DCDI-G (Brouillard et al., 2020)	36	6	2	25	9	0.31
DCDI-DSF (Brouillard et al., 2020)	33	6	2	22	9	0.33
FCI (Spirtes et al., 1993)	35	4	12	21	5	0.22
Imperfect Linear Gaussian (ours)	33	7	11	22	3	0.30
Imperfect Non-Linear Gaussian (ours)	19	7	11	8	0	0.42
Imperfect Normalizing Flow (ours)	30	9	9	21	1	0.38
Perfect Linear Gaussian (ours)	23	8	10	13	3	0.41
Perfect Non-Linear Gaussian (ours)	24	11	7	17	1	0.48
Perfect Normalizing Flow (ours)	23	7	11	12	2	0.38

# Experiments on real data



Edges in the wrong direction

	HD	tp	fn	fp	rev	$F_1$ score
GIES (Hauser and Bühlmann, 2012)	38	10	0	41	7	0.33
CAM (Bühlmann et al., 2014)	35	12	1	30	4	0.51
IGSP (Wang et al., 2017)	18	4	6	5	7	0.42
DCDI-G (Brouillard et al., 2020)	36	6	2	25	9	0.31
DCDI-DSF (Brouillard et al., 2020)	33	6	2	22	9	0.33
FCI (Spirtes et al., 1993)	35	4	12	21	5	0.22
Imperfect Linear Gaussian (ours)	33	7	11	22	3	0.30
Imperfect Non-Linear Gaussian (ours)	19	7	11	8	0	0.42
Imperfect Normalizing Flow (ours)	30	9	9	21	1	0.38
Perfect Linear Gaussian (ours)	23	8	10	13	3	0.41
Perfect Non-Linear Gaussian (ours)	24	11	7	17	1	0.48
Perfect Normalizing Flow (ours)	23	7	11	12	2	0.38

# Conclusion

- Make heterogeneity or non-stationarity your friend.
- Causal Discovery under Latent Interventions

# Causal discovery and Deep Learning

Faria et. al. 2022; Pearl et. al. 2018; Neal et. al. 2022; Schölkopf et. al. 2019; Bengio et. al. 2020; Peters et. al. 2017, Ke et. al. 2020, Brouillard et. al. 2020

Gonçalo R. A. Faria  
@goncalorafaria



Deep Learning Sessions Portugal