

# Causal Effect Inference for Structured Treatments



Jean Kaddour

Yuchen Zhu

Qi Liu

Matt Kusner

Ricardo Silva

# Roadmap

- Motivation
- Generalized Robinson Decomposition
- Quasi-Oracle Convergence Rate
- Structured Intervention Networks
- Experiments
- Summary

# Roadmap

- **Motivation**
- Generalized Robinson Decomposition
- Quasi-Oracle Convergence Rate
- Structured Intervention Networks
- Experiments
- Summary

# Causal Effects of Structured Treatments

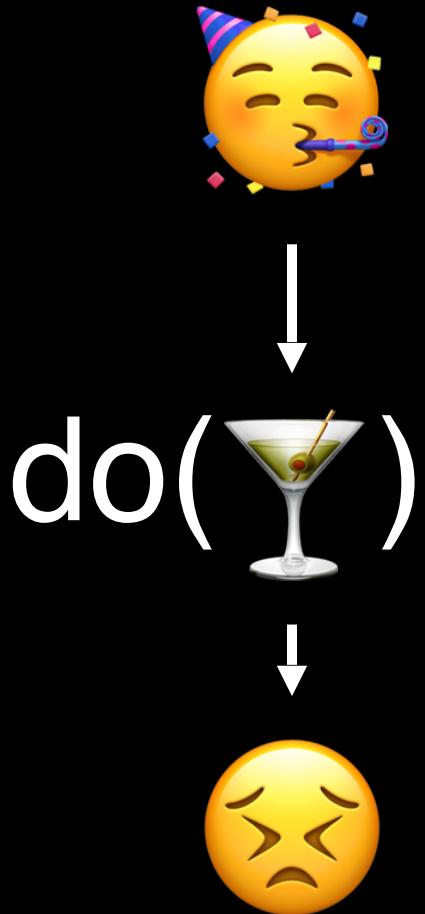


# Causal Effects of Structured Treatments

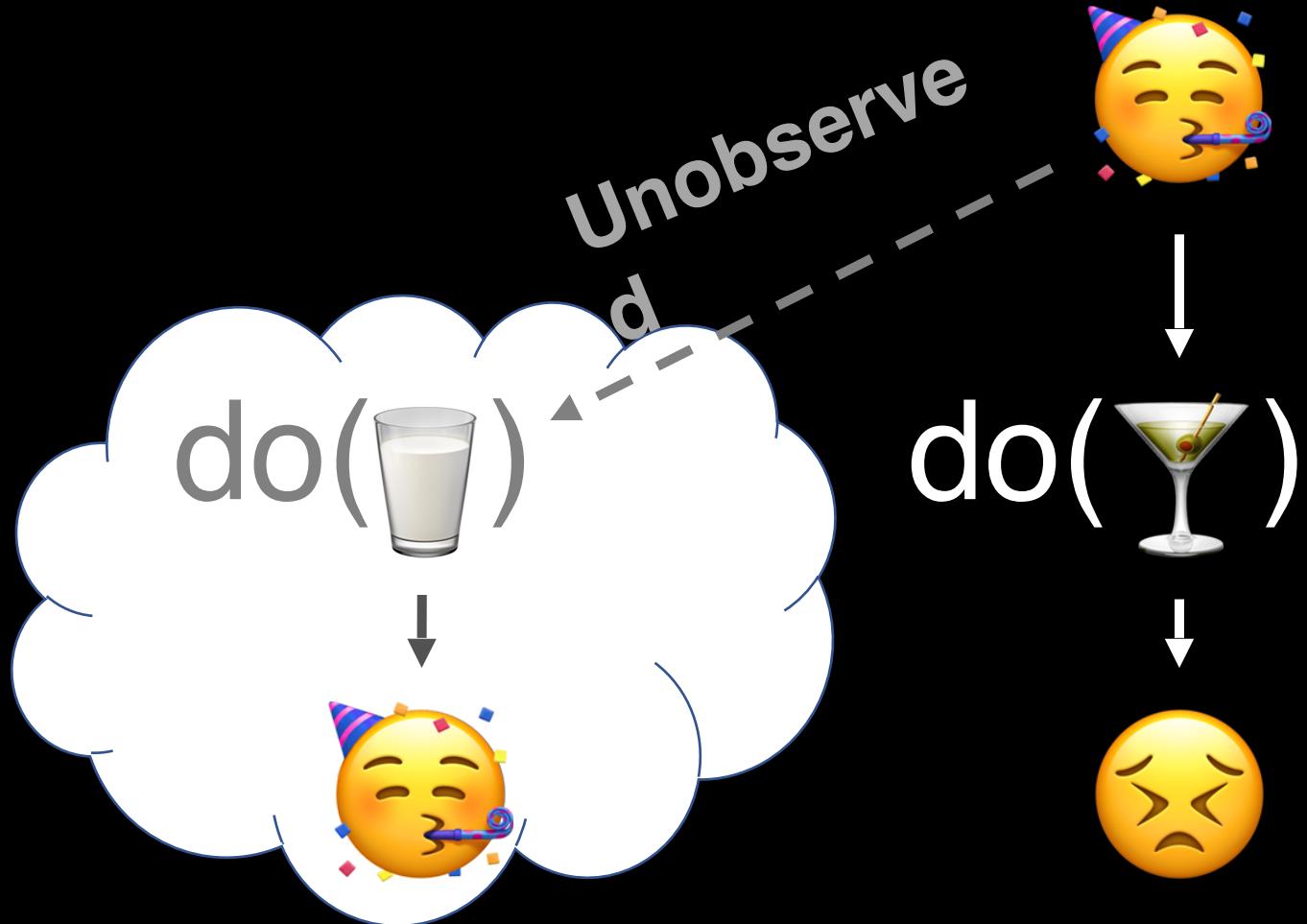


do(  
    )

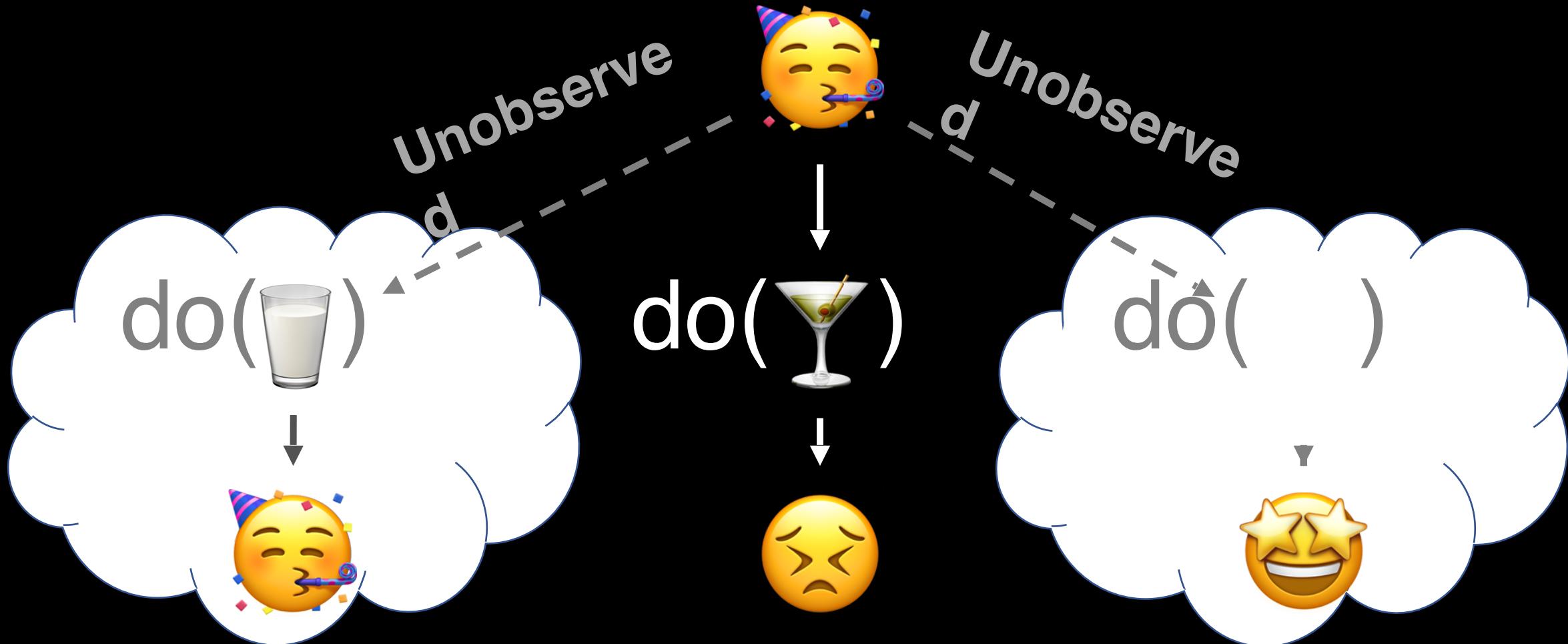
# Causal Effects of Structured Treatments



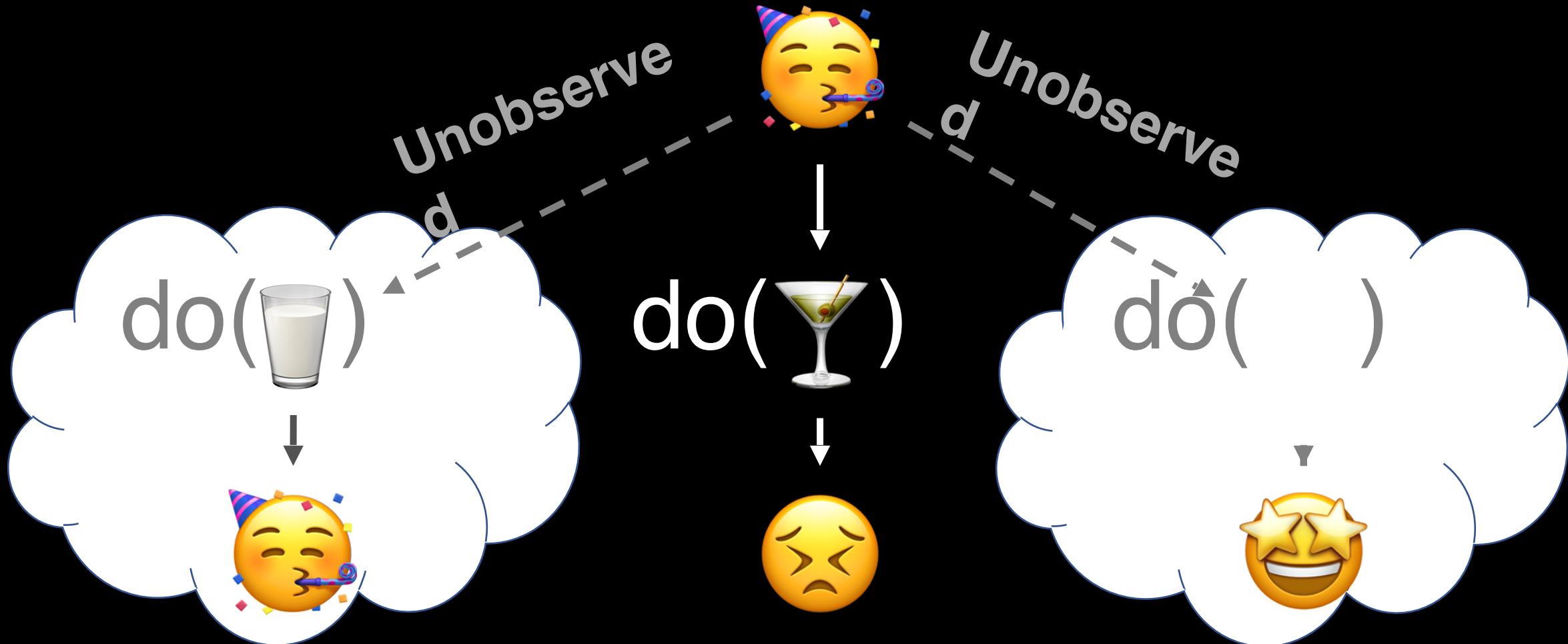
# Causal Effects of Structured Treatments



# Causal Effects of Structured Treatments

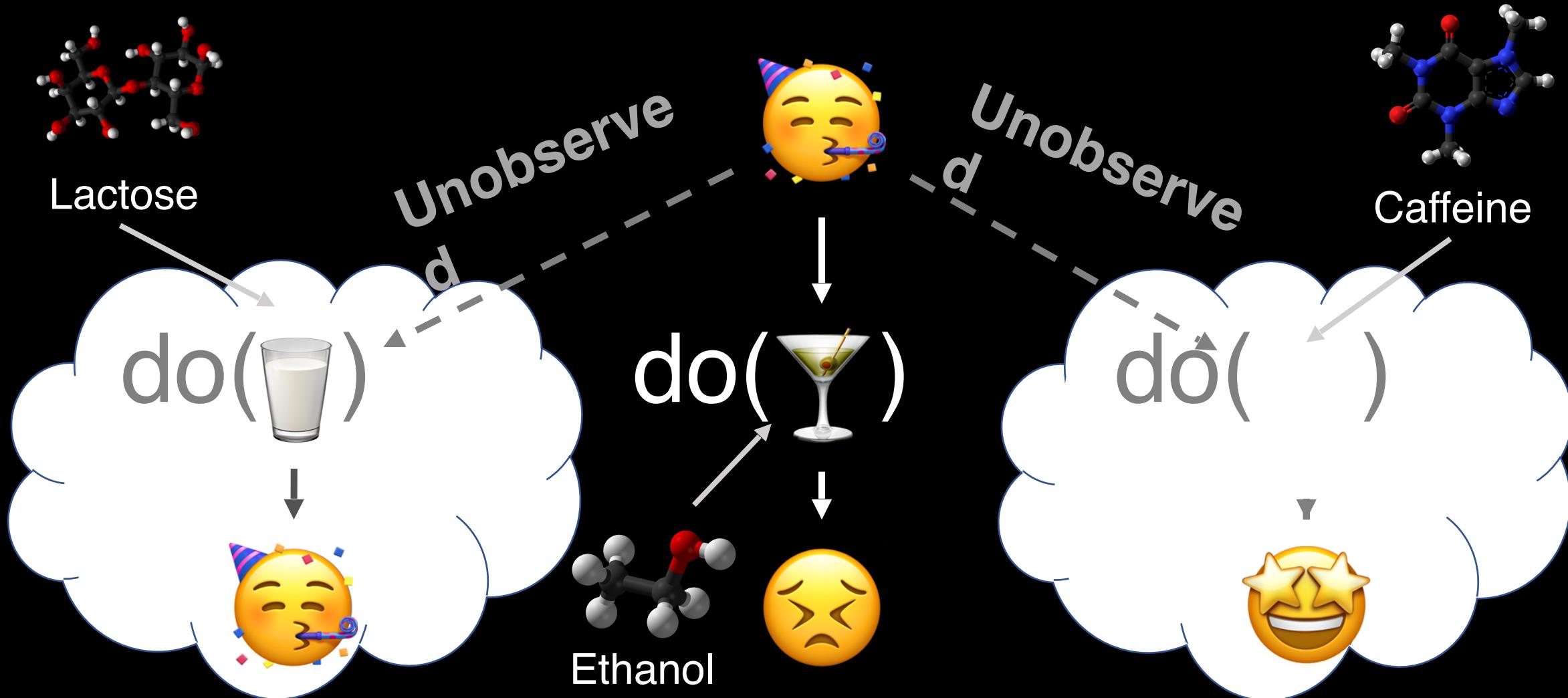


# Causal Effects of Structured Treatments



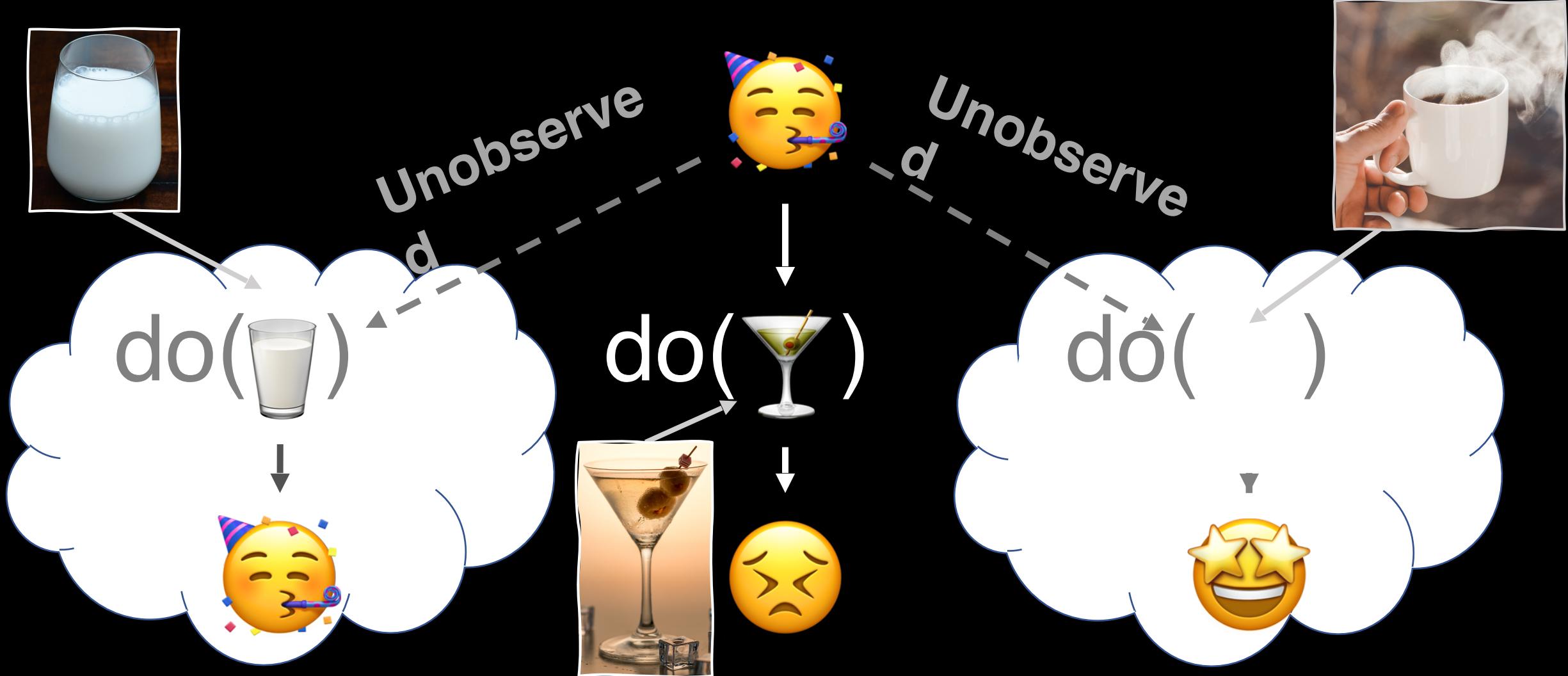
$$\text{CATE: } \tau(\text{, , }) = \mathbb{E}[Y | \text{, do}(\text{ })] - \mathbb{E}[Y | \text{, do}(\text{ })]$$

# Causal Effects of Structured Treatments



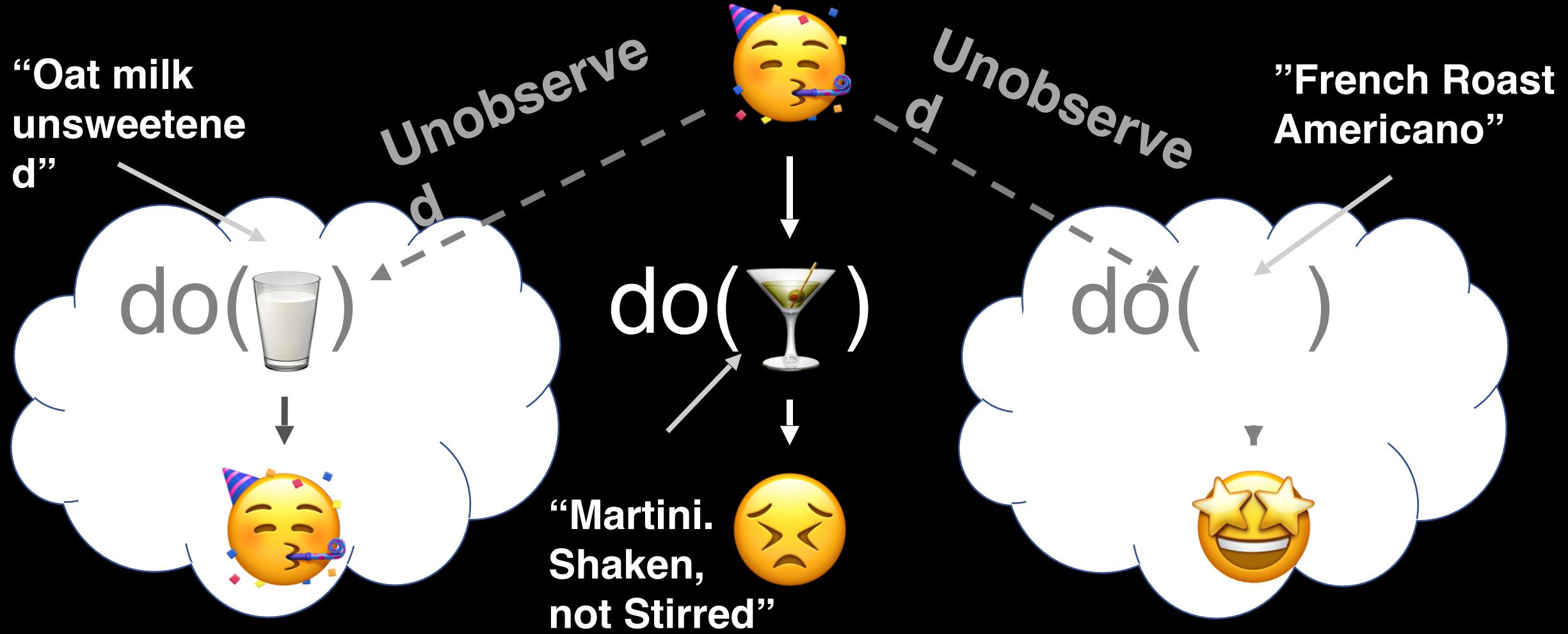
$$\text{CATE: } \tau(\ , \ , \ ) = \mathbb{E}[Y | \ , \ \text{do}(\ )] - \mathbb{E}[Y | \ , \ \text{do}(\ )]$$

# Causal Effects of Structured Treatments



$$\text{CATE: } \tau(\text{, , }) = \mathbb{E}[Y | \text{, } do(\text{ })] - \mathbb{E}[Y | \text{, } do(\text{ })]$$

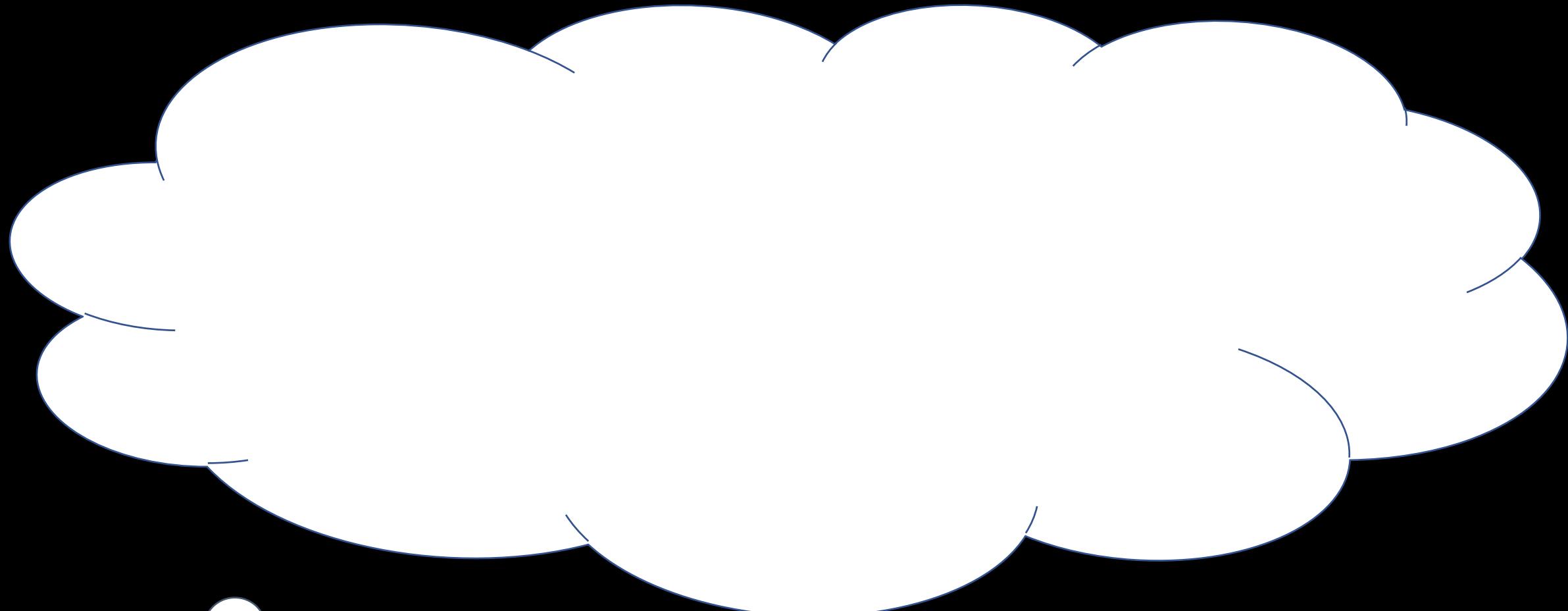
# Causal Effects of Structured Treatments



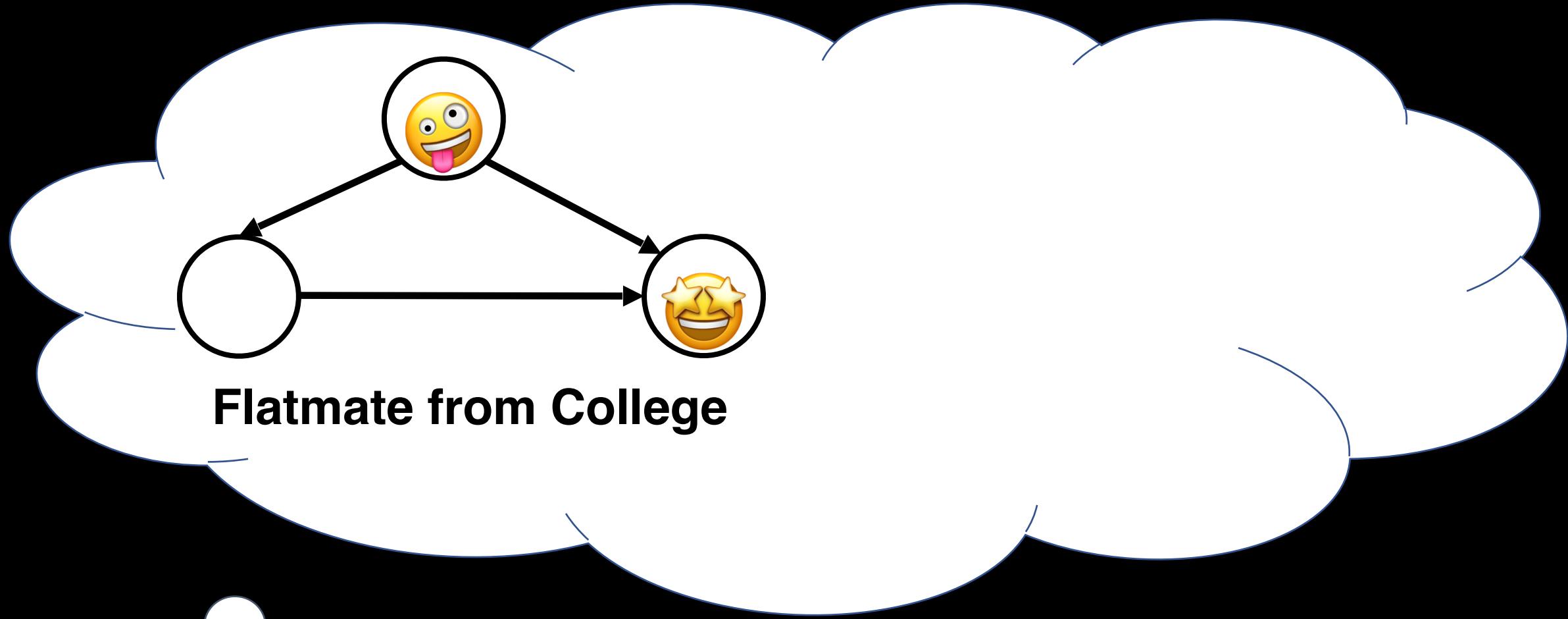
$$\text{CATE: } \tau(\ , \ , \ ) = \mathbb{E}[Y | \ , \text{do}(\ )] - \mathbb{E}[Y | \ , \text{do}(\ )]$$

# You only live once!

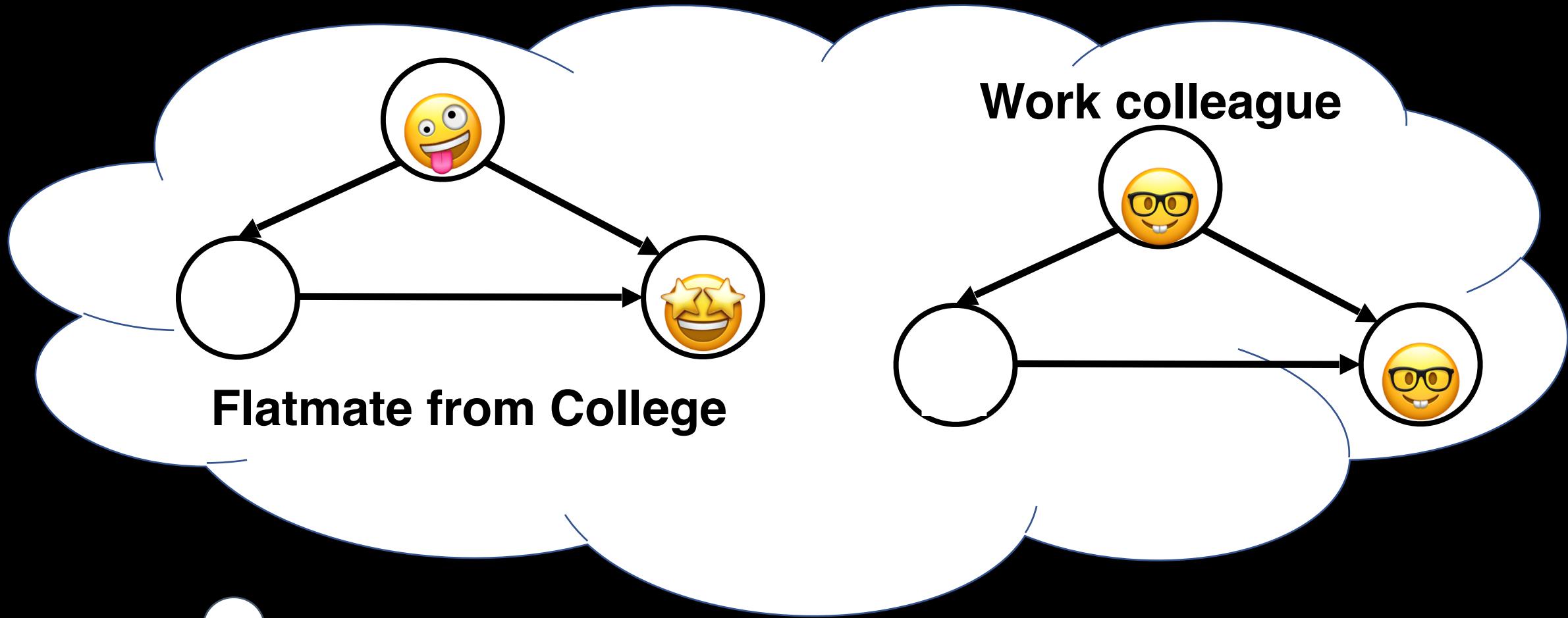
# You only live once!



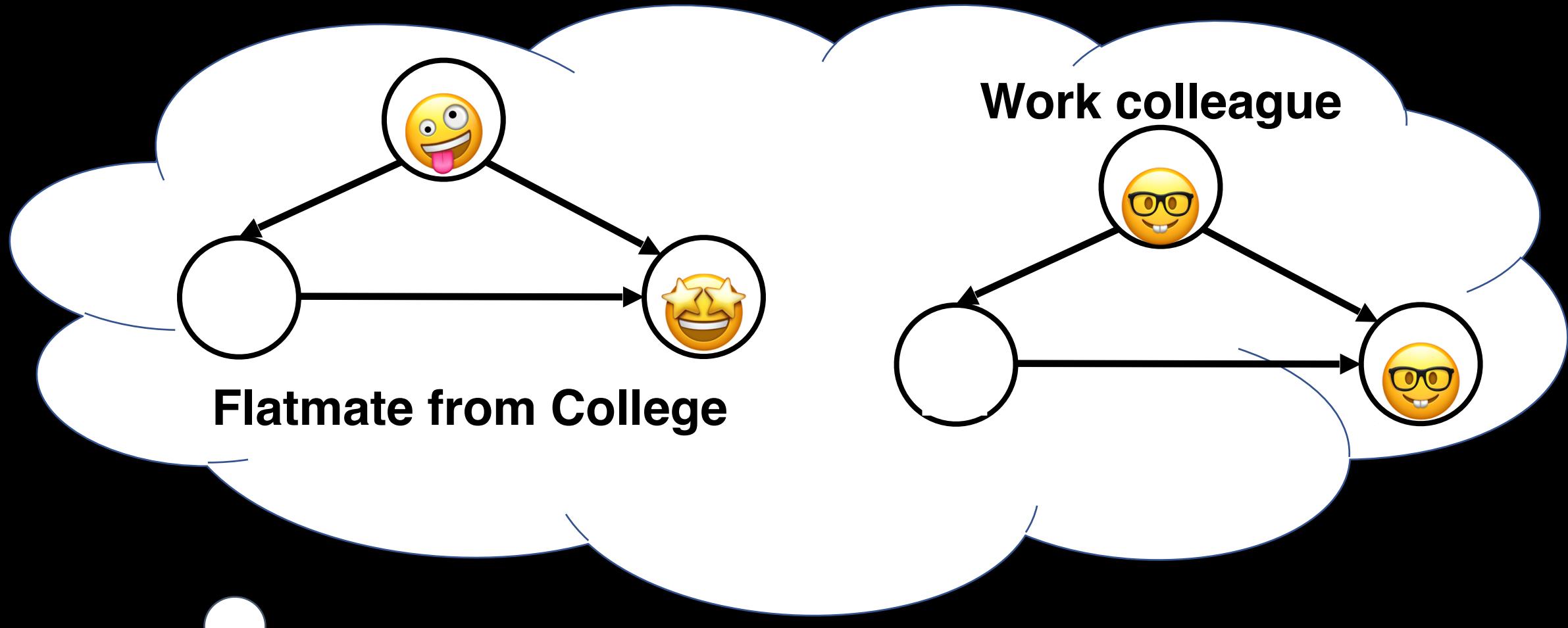
# You only live once!



# You only live once!

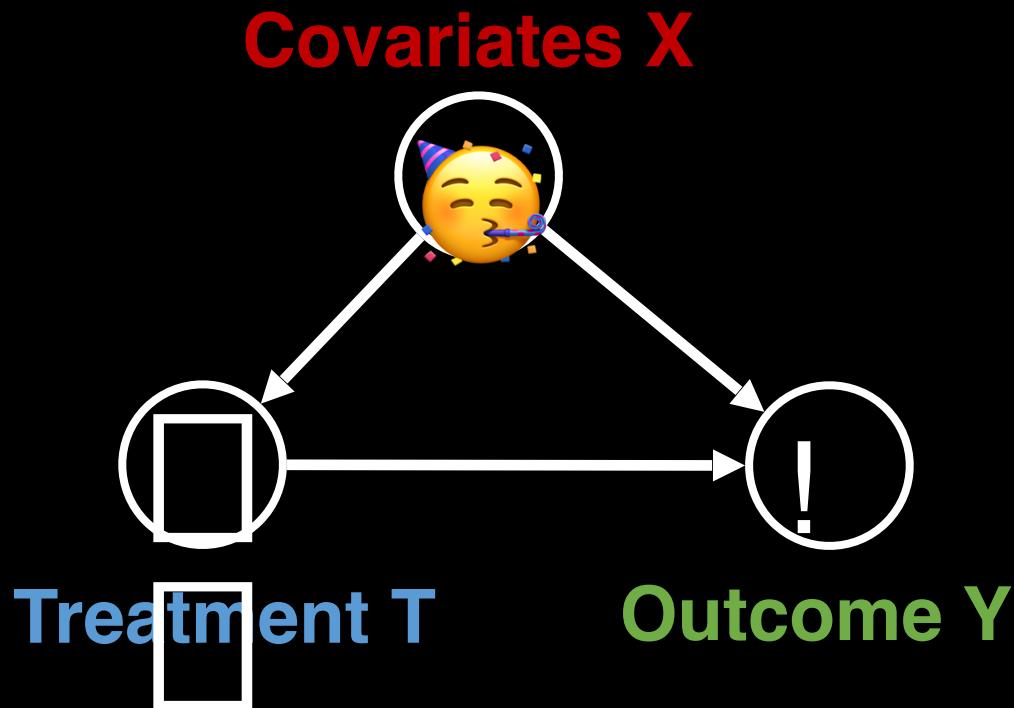


# You only live once!

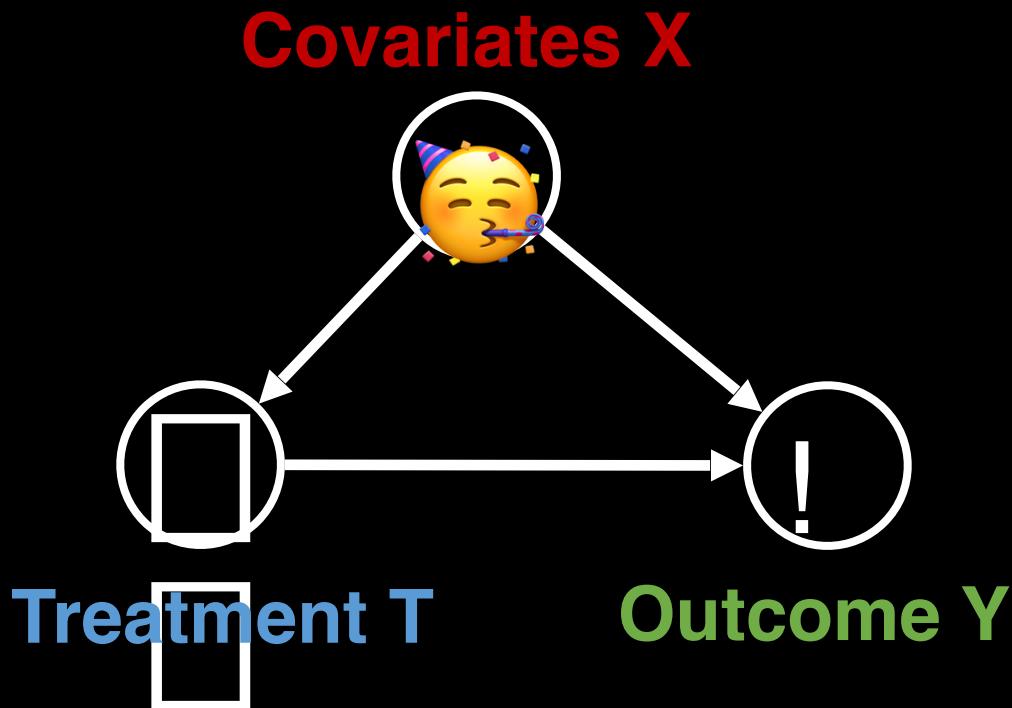


Can we learn from others?

# Goal: Estimate CATE from observational data



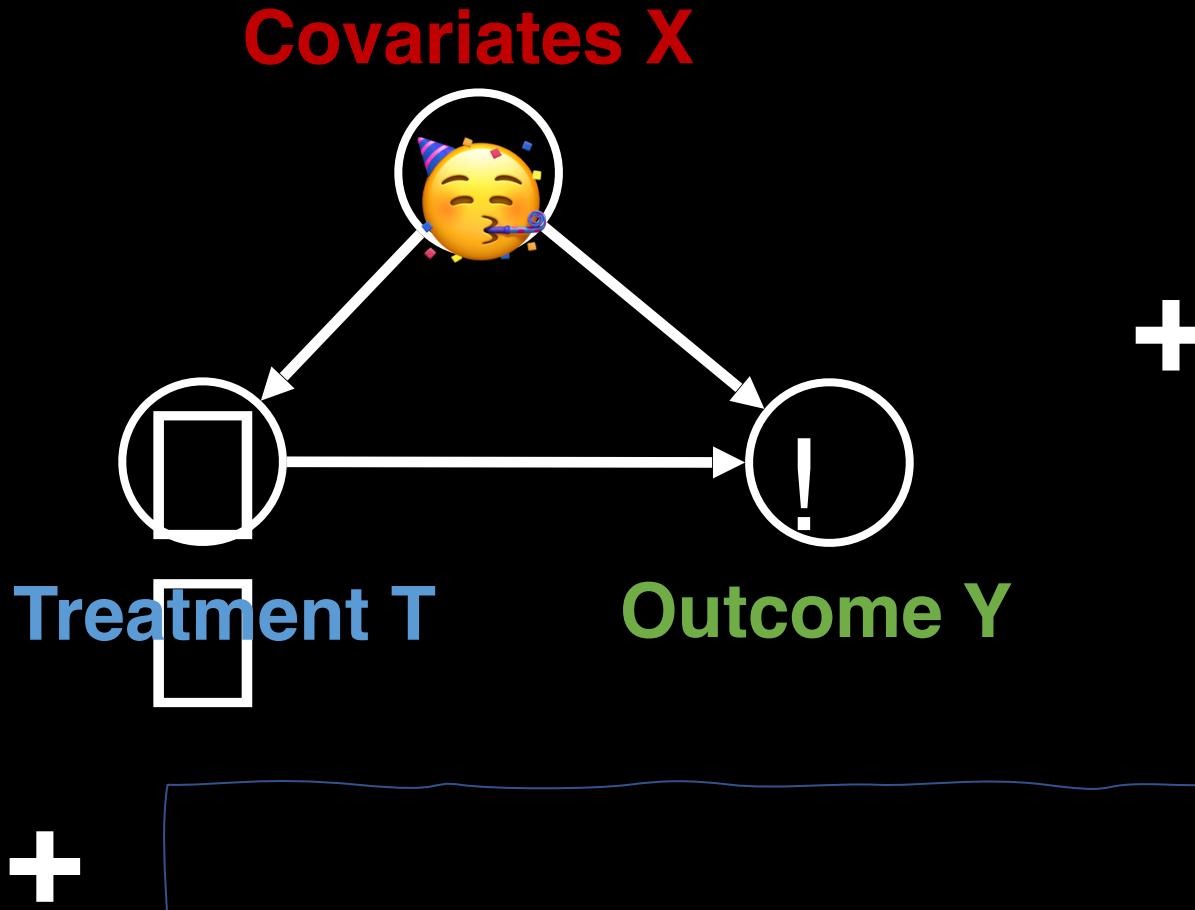
# Goal: Estimate CATE from observational data



+

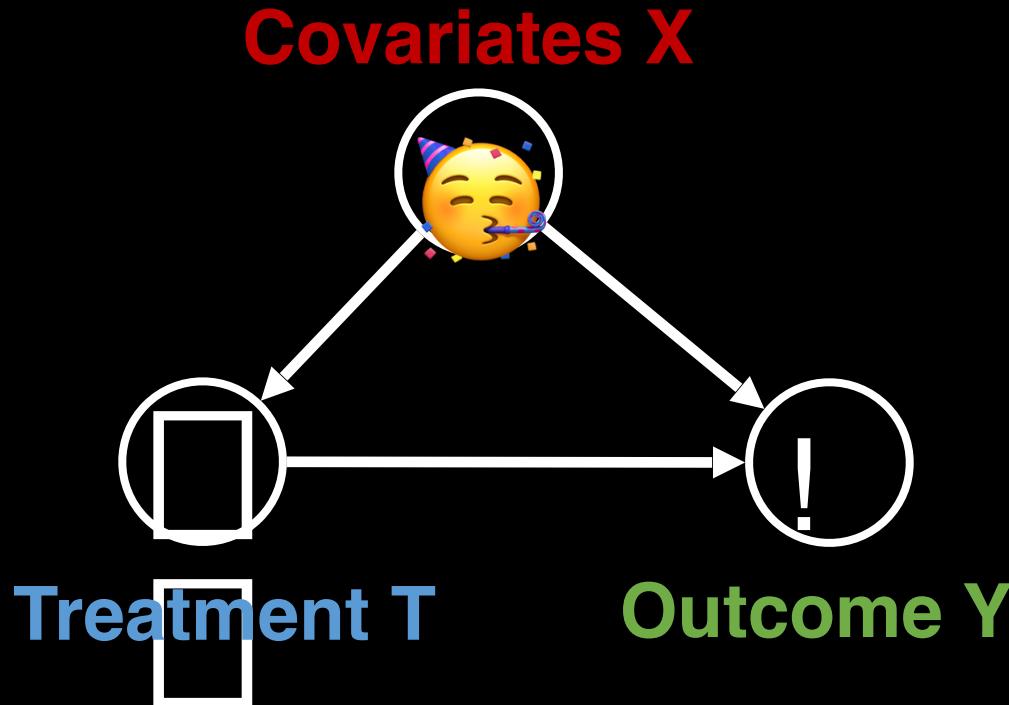
| X | T | Y |
|---|---|---|
| 🎉 | 🍸 | 😢 |
| 😜 | # | 🤩 |
| 🤓 | 🥛 | 🤓 |

# Goal: Estimate CATE from observational data



| X | T | Y |
|---|---|---|
| 🎉 | 🍸 | 😢 |
| 😜 | # | 🤩 |
| 🤓 | 🥛 | 🤓 |

# Goal: Estimate CATE from observational data



+

| X | T | Y |
|---|---|---|
| 🎉 | 🍸 | 😢 |
| 😜 | # | 🤩 |
| 🤓 | 🥛 | 🤓 |

+

# □ □

=

□ #

□ □



Wait, can't I simply learn a model

$$\hat{f} \approx \mathbb{E}[Y \mid \mathbf{X}, \mathbf{T}]$$

Wait, can't I simply learn a model

$$\hat{f} \approx \mathbb{E}[Y \mid \mathbf{X}, \mathbf{T}]$$

and then subtract predictions

$$\hat{\tau}(t', t, \mathbf{x}) = \hat{f}(\mathbf{x}, t') - \hat{f}(\mathbf{x}, t)?$$

Technically, yes! But...

# Technically, yes! But...



# Technically, yes! But...

*“When solving a problem of interest,  
do not solve a more general problem  
as an intermediate step.*



# Technically, yes! But...

*“When solving a problem of interest,  
do not solve a more general problem  
as an intermediate step.*

*Try to get the answer that you really  
need but not a more general one.”*

**Vladimir Vapnik, 2006.**



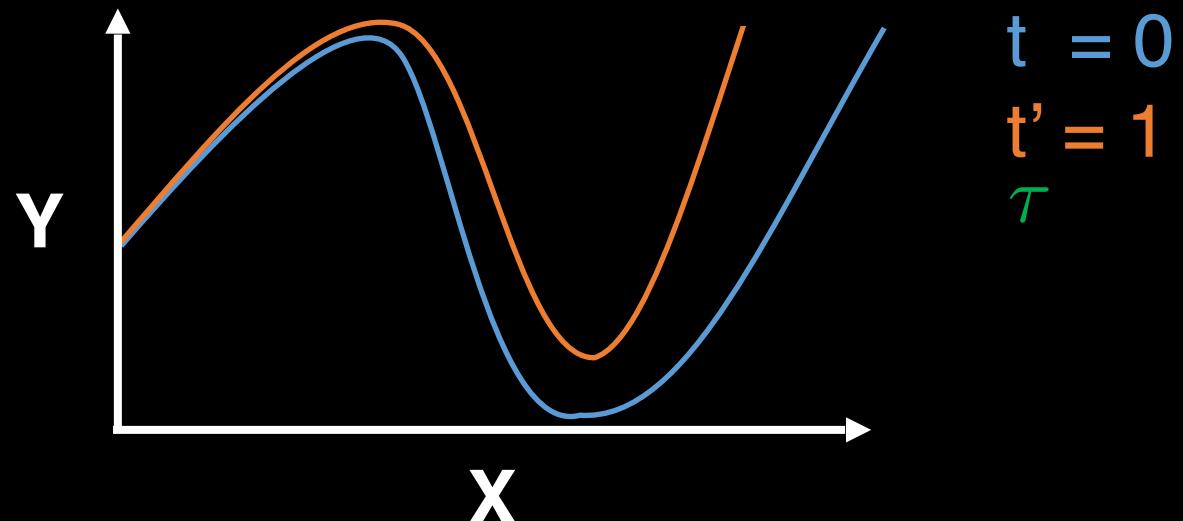
# Why targeting the effect directly helps

# Why targeting the effect directly helps

- **Reason 1:** the effect often exhibits a simpler structure

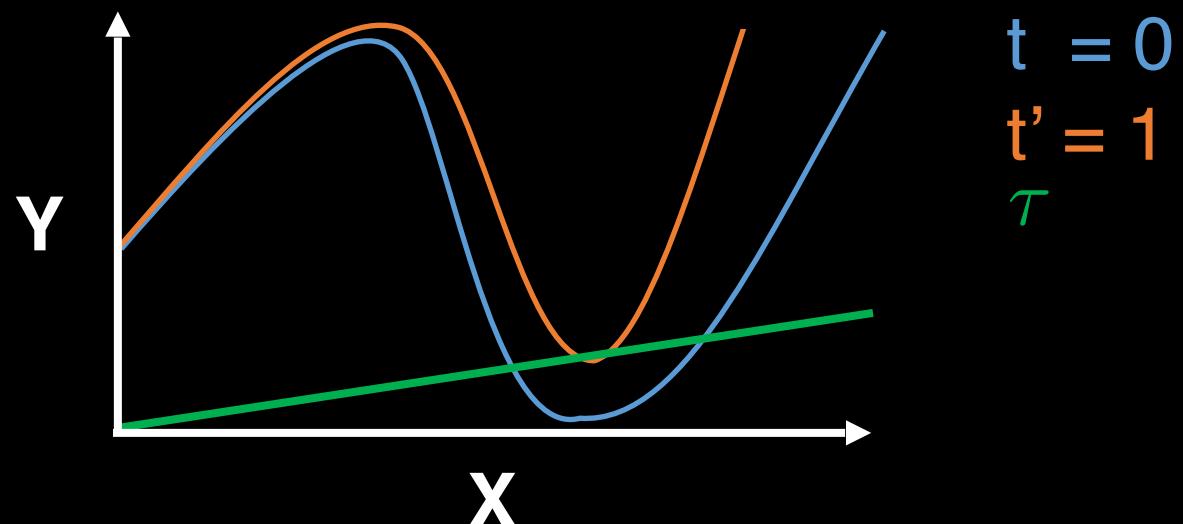
# Why targeting the effect directly helps

- Reason 1: the effect often exhibits a simpler structure
- For example,  $\hat{f}(\mathbf{X}, \mathbf{T})$  can be very non-smooth for rarely treated  $\mathbf{X}$



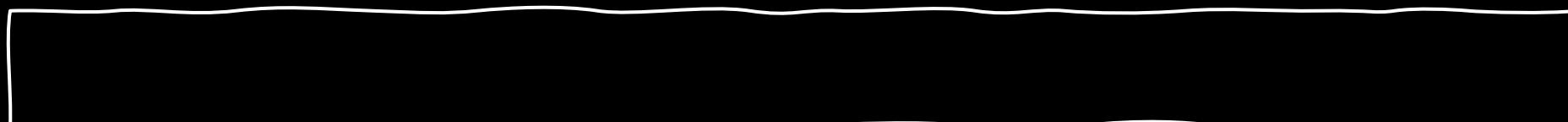
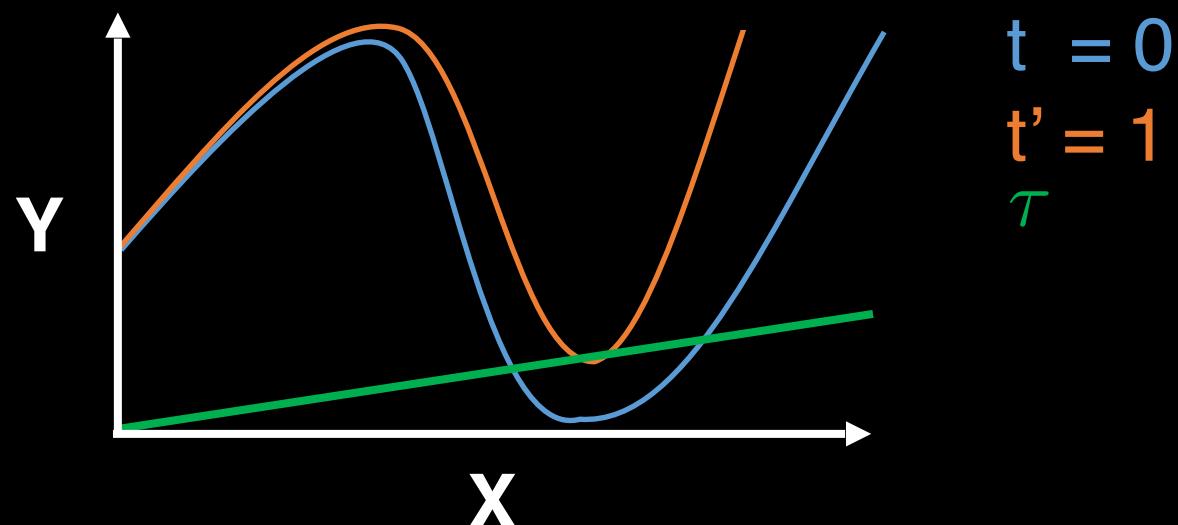
# Why targeting the effect directly helps

- Reason 1: the effect often exhibits a simpler structure
  - For example,  $\hat{f}(\mathbf{X}, \mathbf{T})$  can be very non-smooth for rarely treated  $\mathbf{X}$
  - However,  $f(\mathbf{X}, t') - f(\mathbf{X}, t)$  might be (almost) linear across  $\mathbf{X}$



# Why targeting the effect directly helps

- Reason 1: the effect often exhibits a simpler structure
- For example,  $\hat{f}(\mathbf{X}, \mathbf{T})$  can be very non-smooth for rarely treated  $\mathbf{X}$
- However,  $f(\mathbf{X}, t') - f(\mathbf{X}, t)$  might be (almost) linear across  $\mathbf{X}$



# Why targeting the effect directly helps

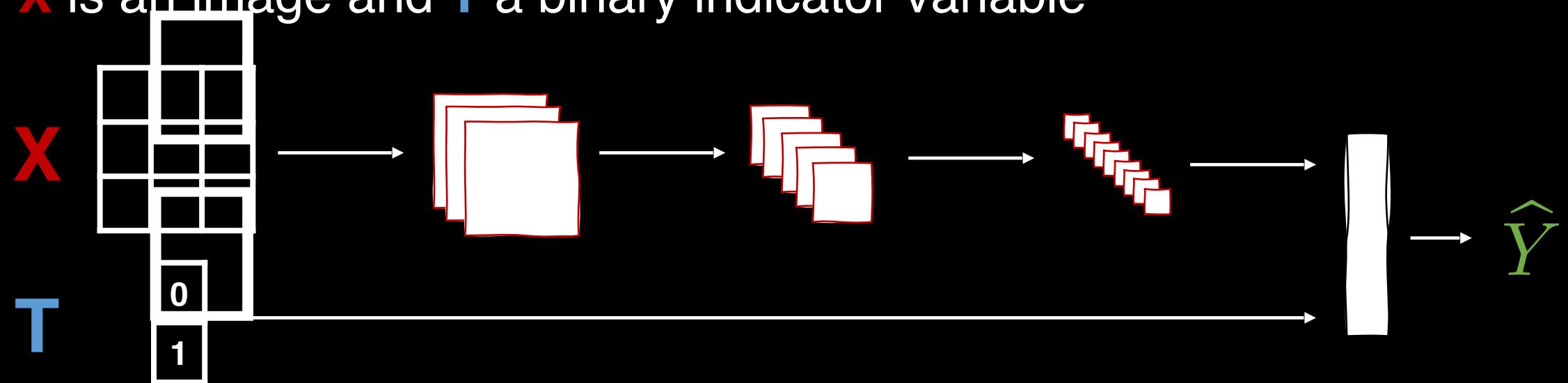
- **Reason 2:** it reduces regularization bias

# Why targeting the effect directly helps

- **Reason 2:** it reduces regularization bias
- e.g., consider **X** and **T** with vastly different dimensionalities

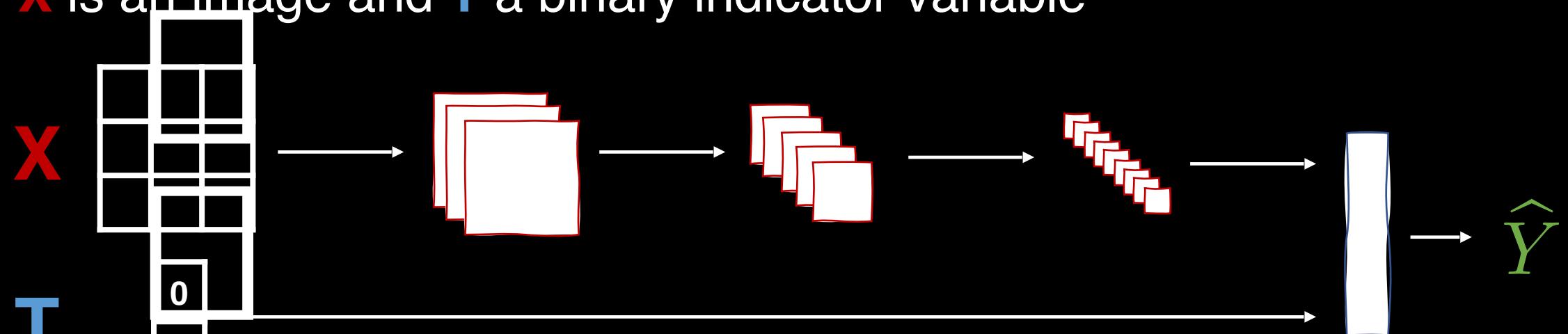
# Why targeting the effect directly helps

- **Reason 2:** it reduces regularization bias
  - e.g., consider  $\mathbf{X}$  and  $\mathbf{T}$  with vastly different dimensionalities
  - $\mathbf{X}$  is an image and  $\mathbf{T}$  a binary indicator variable



# Why targeting the effect directly helps

- **Reason 2:** it reduces regularization bias
  - e.g., consider  $\mathbf{X}$  and  $\mathbf{T}$  with vastly different dimensionalities
  - $\mathbf{X}$  is an image and  $\mathbf{T}$  a binary indicator variable



- Regularizing the model in favour of predicting the outcome can make it ignore the lower-dimensional variable (towards zero-effect)

# Previous work in CATE estimation

- **S-Learner** (Hill, 2011)
- **T-Learner** (Athey & Imbens, 2016)
- **R-Learner** (Nie & Wager, 2017)
- **CFRnet/TARnet** (Shalit et al., 2017)
- **X-Learner** (Künzel et al., 2018)
- Perfect Match (Schwab et al. 2018)
- **Multitask-Learner** (Alaa & van der Schaar, 2018)
- **Bayesian Causal Forest** (Hahn et al., 2020)
- **VCnet** (Nie et al., 2021)
- **FlexTENet** (Curth & van der Schaar, 2021)
- ...

# Previous work in CATE estimation

- **S-Learner** (Hill, 2011)
- **T-Learner** (Athey & Imbens, 2016)
- **R-Learner** (Nie & Wager, 2017)
- **CFRnet/TARnet** (Shalit et al., 2017)
- **X-Learner** (Künzel et al., 2018)
- Perfect Match (Schwab et al. 2018)
- **Multitask-Learner** (Alaa & van der Schaar, 2018)
- **Bayesian Causal Forest** (Hahn et al., 2020)
- **VCnet** (Nie et al., 2021)
- **FlexTENet** (Curth & van der Schaar, 2021)
- ...

Most of these deal with binary or scalar-continuous treatments.

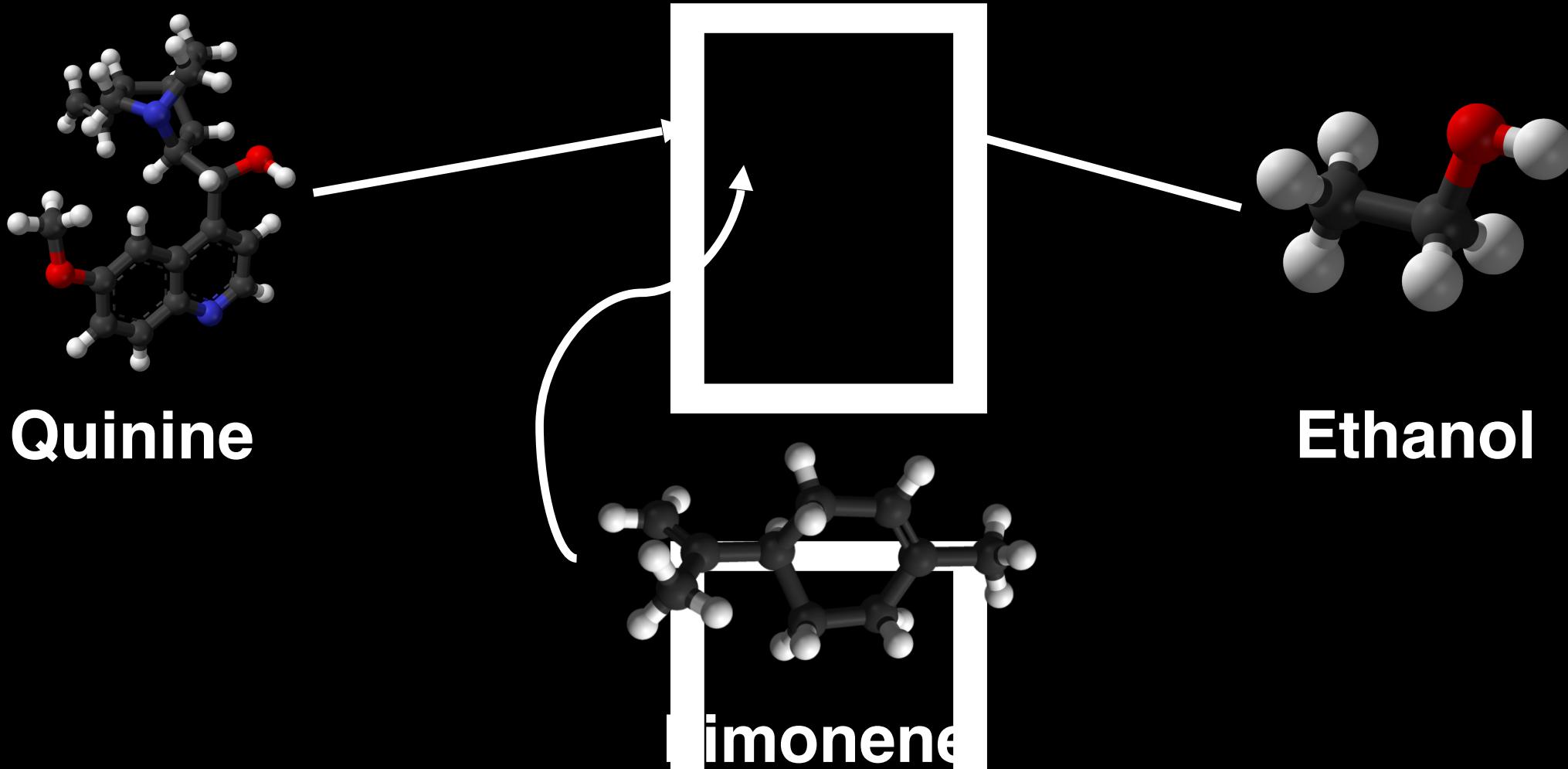
# Why structured treatments?

- 1. Data-Efficiency**
- 2. (Infinitely)-Many-Treatments-Setup**
- 3. Generalization to unseen treatments**

# 1. Data-Efficiency

# 1. Data-Efficiency

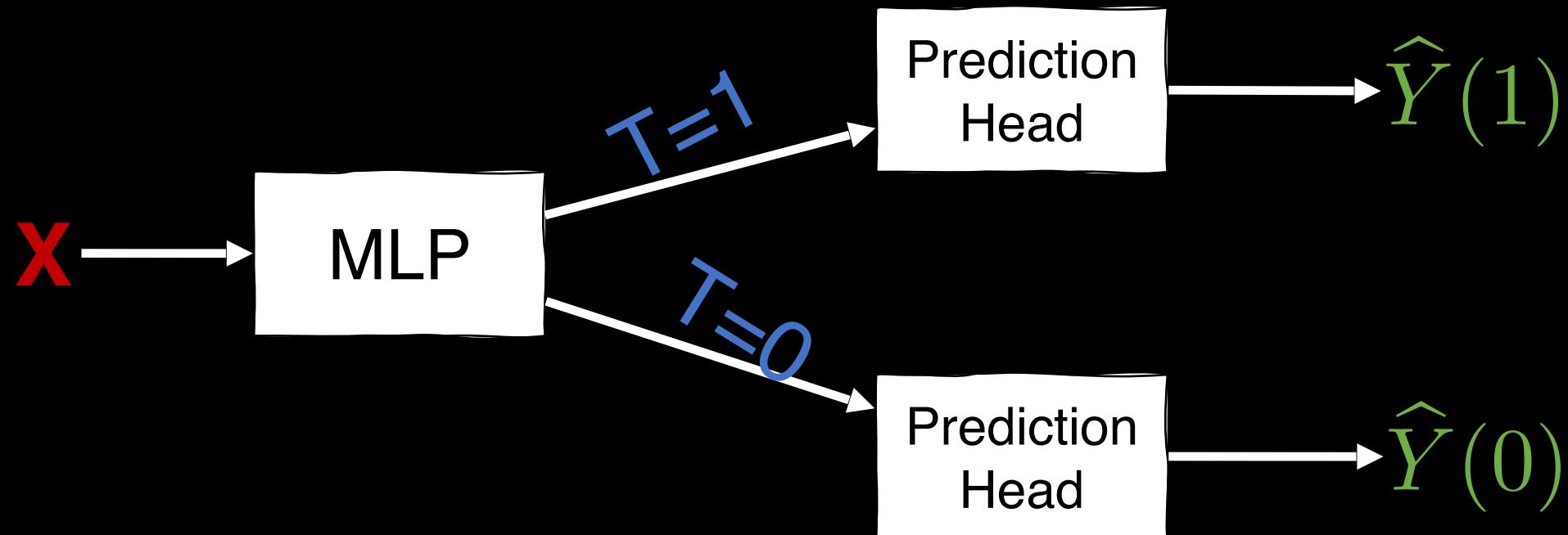
**Example:** the molecular graph of a drink



## 2. (Infinitely-)Many-Treatments settings

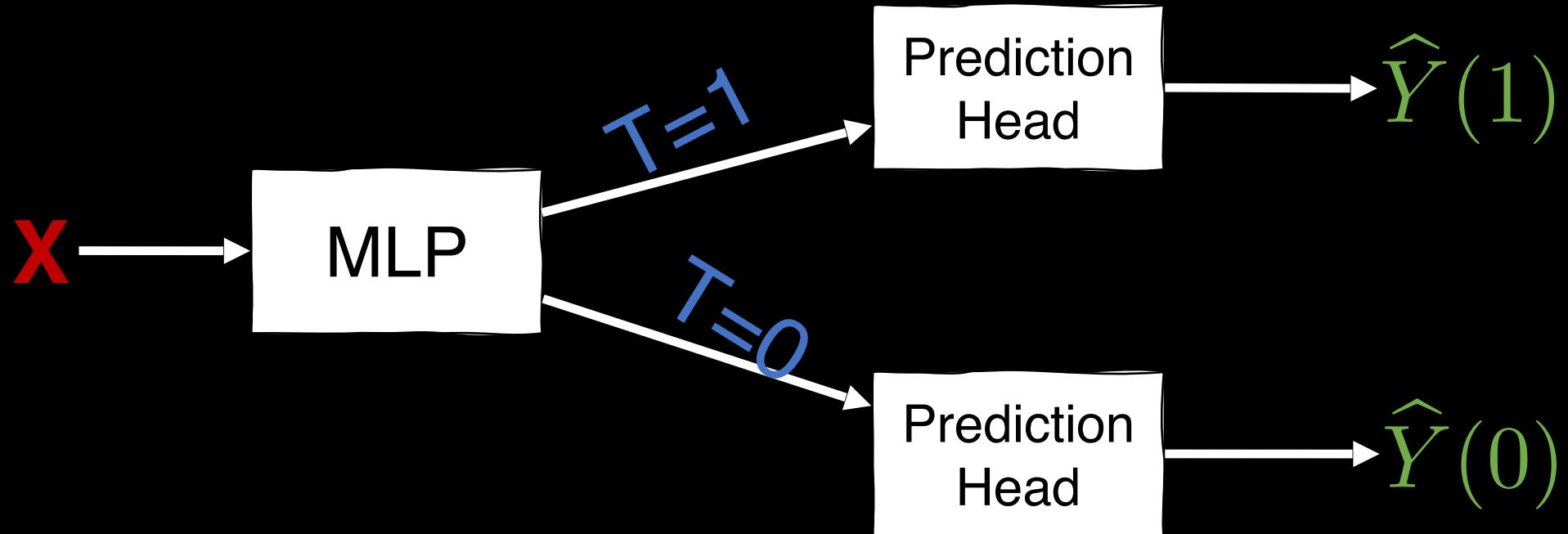
## 2. (Infinitely-)Many-Treatments settings

Typical NN architectures for CATE estimation



## 2. (Infinitely-)Many-Treatments settings

Typical NN architectures for CATE estimation



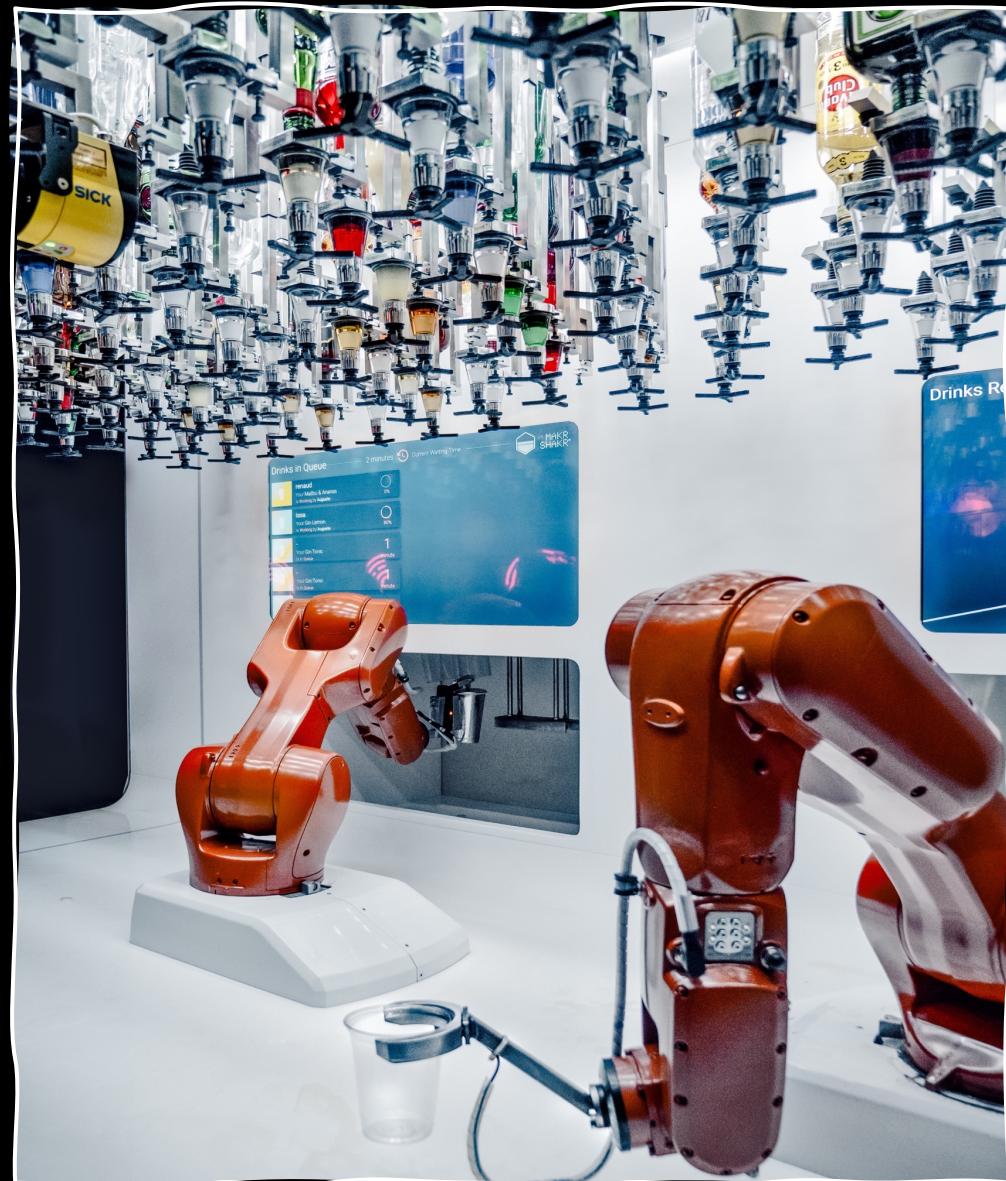
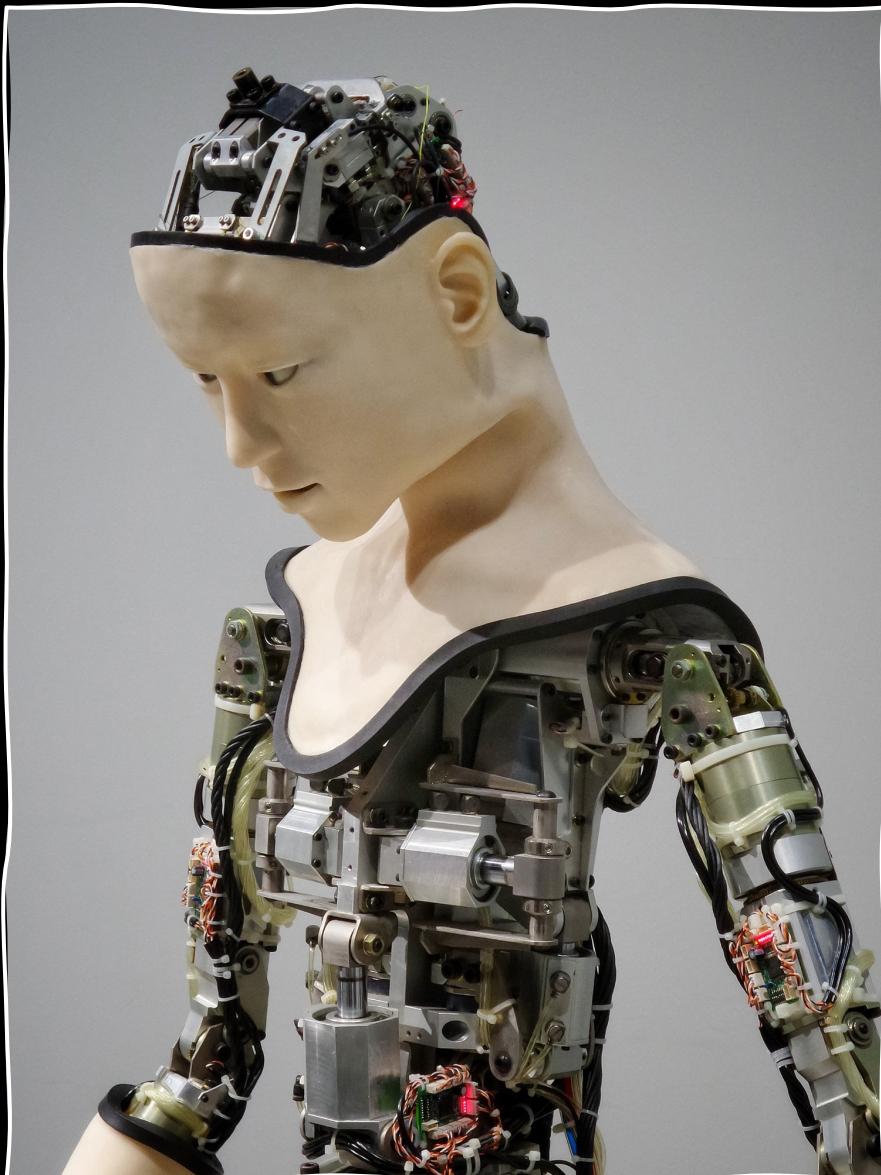
**This does not scale well in the treatment options!**

## 2. (Infinitely-)Many-Treatments settings



# 3. Generalization to unseen treatments

### 3. Generalization to unseen treatments



# How to deal with structured treatments?

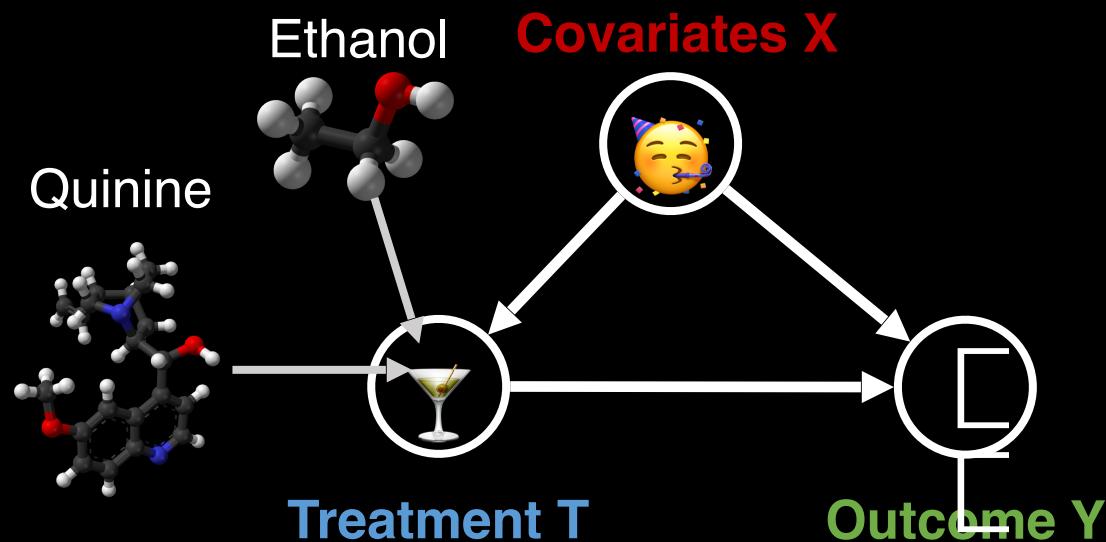
# How to deal with structured treatments?

- **Issue 1:** How can we represent the CATE with struct. treatments?

# How to deal with structured treatments?

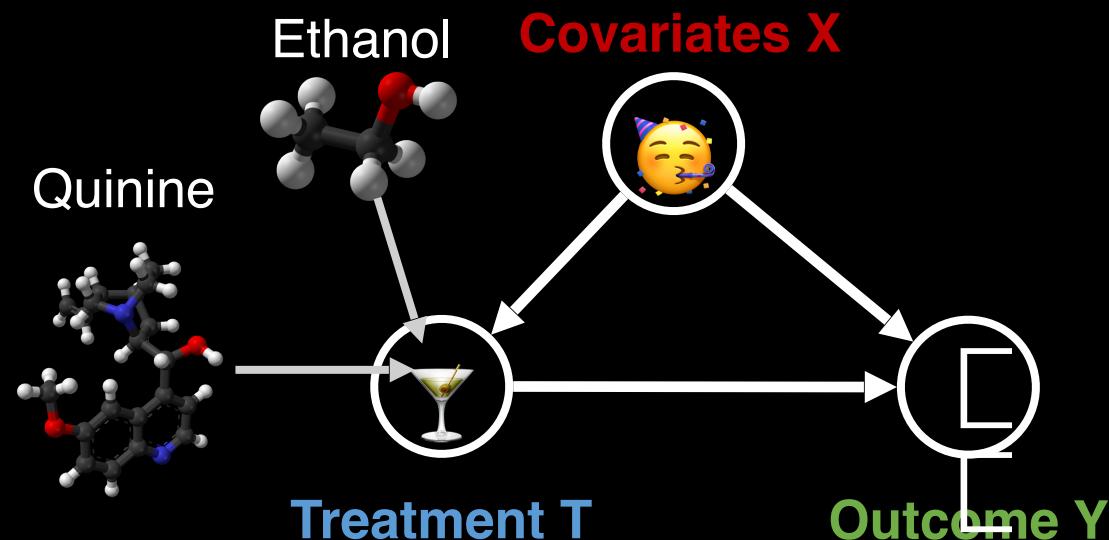
- Issue 1: How can we represent the CATE with struct. treatments?

➤ Goal: Taking covariates AND treatment features into account



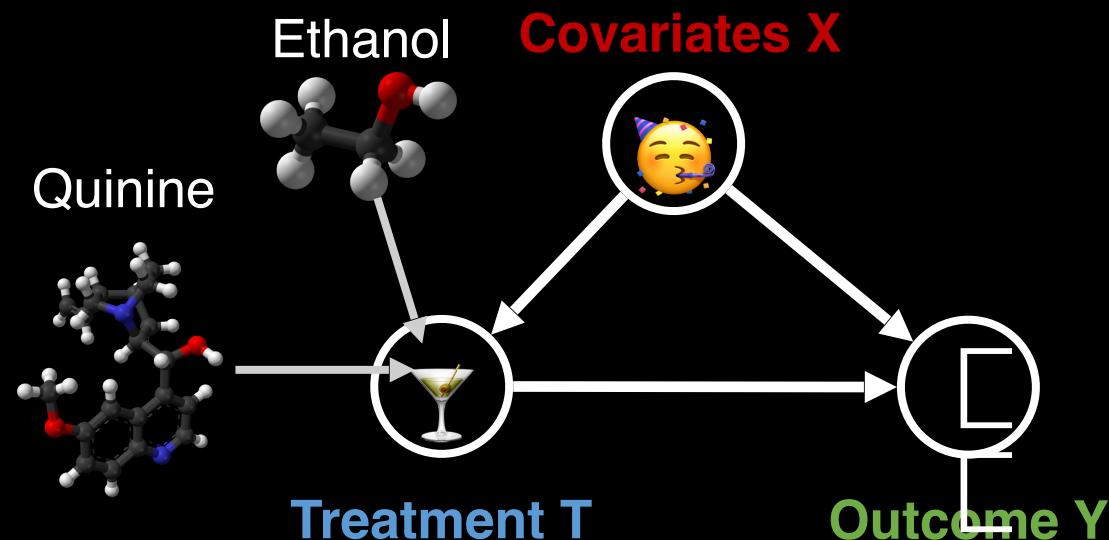
# How to deal with structured treatments?

- **Issue 1:** How can we represent the CATE with struct. treatments?
  - **Goal:** Taking **covariates** AND **treatment features** into account
- **Issue 2:** How can we learn the relevant feature maps of the effect?



# How to deal with structured treatments?

- **Issue 1:** How can we represent the CATE with struct. treatments?
  - **Goal:** Taking **covariates** AND **treatment features** into account
- **Issue 2:** How can we learn the relevant feature maps of the effect?
  - **Goal:** Deriving a trainable objective that targets the effect



# Roadmap

- Motivation
- **Generalized Robinson Decomposition**
- Quasi-Oracle Convergence Rate
- Structured Intervention Networks
- Experiments
- Summary

# Robinson Decomposition/R-Learner<sup>1</sup>

# Robinson Decomposition/R-Learner<sup>1</sup>

- Allows us to construct a learnable objective of the binary CATE

# Robinson Decomposition/R-Learner<sup>1</sup>

- Allows us to construct a learnable objective of the binary CATE
- Define the *propensity score*  $e(\mathbf{x}) \triangleq p(T = 1 \mid \mathbf{x})$

# Robinson Decomposition/R-Learner<sup>1</sup>

- Allows us to construct a learnable objective of the binary CATE
- Define the *propensity score*  $e(\mathbf{x}) \triangleq p(T = 1 \mid \mathbf{x})$
- Define the *conditional mean outcome*  $m(\mathbf{x}) \triangleq \mathbb{E}[Y \mid \mathbf{x}]$

# Robinson Decomposition/R-Learner<sup>1</sup>

- Allows us to construct a learnable objective of the binary CATE
- Define the *propensity score*  $e(\mathbf{x}) \triangleq p(T = 1 \mid \mathbf{x})$
- Define the *conditional mean outcome*  $m(\mathbf{x}) \triangleq \mathbb{E}[Y \mid \mathbf{x}]$
- Define  $\tilde{y}_i \triangleq y_i - \hat{m}(\mathbf{x}_i)$  and  $\tilde{t}_i \triangleq t_i - \hat{e}(\mathbf{x}_i)$  yield the objective

$$\hat{\tau}_b(\cdot) = \arg \min_{\tau_b} \left\{ \frac{1}{n} \sum_{i=1}^n (\tilde{y}_i - \tilde{t}_i \times \tau_b(\mathbf{x}_i))^2 + \Lambda(\tau_b(\cdot)) \right\}$$

- We call  $\hat{m}(\mathbf{x})$  and  $\hat{e}(\mathbf{x})$  the (estimated) **nuisance components**

# Generalizing RD to structured treatments I

# Generalizing RD to structured treatments I

- **Product Effect Assumption:** Re-parameterize the outcome surface as

$$Y = g(\mathbf{X})^\top h(\mathbf{T}) + \varepsilon$$

where  $g : \mathcal{X} \rightarrow \mathbb{R}^d$ ,  $h : \mathcal{T} \rightarrow \mathbb{R}^d$  are feature maps

# Generalizing RD to structured treatments I

- **Product Effect Assumption:** Re-parameterize the outcome surface as

$$Y = g(\mathbf{X})^\top h(\mathbf{T}) + \varepsilon$$

where  $g : \mathcal{X} \rightarrow \mathbb{R}^d$ ,  $h : \mathcal{T} \rightarrow \mathbb{R}^d$  are feature maps

- **Universality property:** As we let the dimensionality of  $g(\cdot)$  and  $h(\cdot)$  grow, we may approximate any bounded function in  $\mathcal{C}(\mathcal{X} \times \mathcal{T})$

# Generalizing RD to structured treatments I

- **Product Effect Assumption:** Re-parameterize the outcome surface as

$$Y = g(\mathbf{X})^\top h(\mathbf{T}) + \varepsilon$$

where  $g : \mathcal{X} \rightarrow \mathbb{R}^d$ ,  $h : \mathcal{T} \rightarrow \mathbb{R}^d$  are feature maps

- **Universality property:** As we let the dimensionality of  $g(\cdot)$  and  $h(\cdot)$  grow, we may approximate any bounded function in  $\mathcal{C}(\mathcal{X} \times \mathcal{T})$

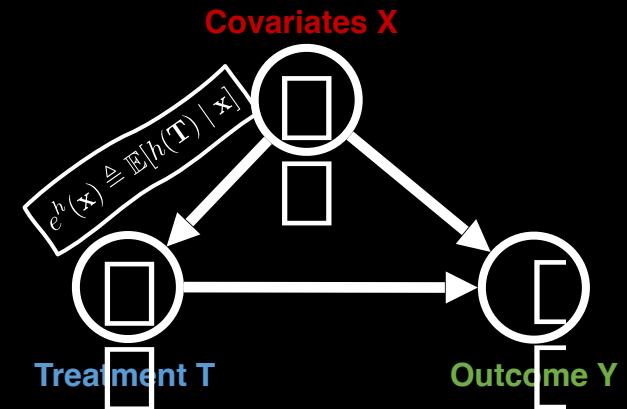
~~Issue 1: How can we represent the CATE effect?~~

$$\tau(\mathbf{t}', \mathbf{t}, \mathbf{x}) = g(\mathbf{x})^\top (h(\mathbf{t}') - h(\mathbf{t}))$$

# Generalizing RD to structured treatments II

- Define *propensity features*

$$e^h(\mathbf{x}) \triangleq \mathbb{E}[h(\mathbf{T}) \mid \mathbf{x}]$$

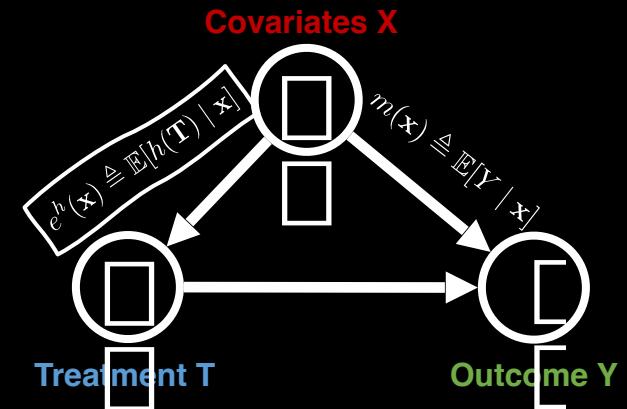


# Generalizing RD to structured treatments II

- Define *propensity features*

$$e^h(\mathbf{x}) \triangleq \mathbb{E}[h(\mathbf{T}) \mid \mathbf{x}]$$

- Recall  $m(\mathbf{x}) \triangleq \mathbb{E}[Y \mid \mathbf{x}] = g(\mathbf{x})^\top e^h(\mathbf{x})$



# Generalizing RD to structured treatments II

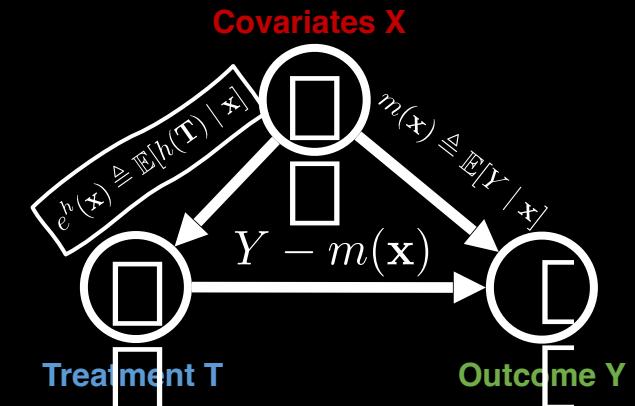
- Define *propensity features*

$$e^h(\mathbf{x}) \triangleq \mathbb{E}[h(\mathbf{T}) \mid \mathbf{x}]$$

- Recall  $m(\mathbf{x}) \triangleq \mathbb{E}[Y \mid \mathbf{x}] = g(\mathbf{x})^\top e^h(\mathbf{x})$

- Following the same steps as for the binary treatment case, we yield

$$Y - m(\mathbf{X}) = g(\mathbf{X})^\top (h(\mathbf{T}) - e^h(\mathbf{X})) + \varepsilon$$



# Generalizing RD to structured treatments II

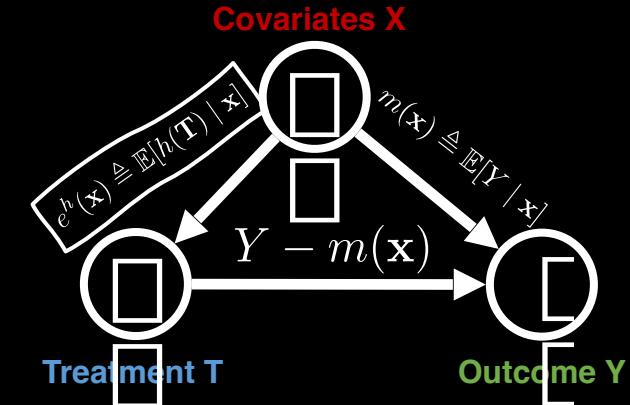
- Define *propensity features*

$$e^h(\mathbf{x}) \triangleq \mathbb{E}[h(\mathbf{T}) \mid \mathbf{x}]$$

- Recall  $m(\mathbf{x}) \triangleq \mathbb{E}[Y \mid \mathbf{x}] = g(\mathbf{x})^\top e^h(\mathbf{x})$

- Following the same steps as for the binary treatment case, we yield

$$Y - m(\mathbf{X}) = g(\mathbf{X})^\top (h(\mathbf{T}) - e^h(\mathbf{X})) + \varepsilon$$



Issue 2: How can we learn the relevant feature maps of the effect?

Solution: For a fixed  $h(\cdot)$ , a generalization to structured treatments is

$$\hat{g}(\cdot) = \arg \min_g \left\{ \frac{1}{n} \sum_{i=1}^n \left( Y_i - \hat{m}(\mathbf{X}_i) - g(\mathbf{X}_i)^\top (h(\mathbf{T}_i) - \hat{e}^h(\mathbf{X}_i)) \right)^2 \right\}$$

# Roadmap

- Motivation
- **Generalized Robinson Decomposition**
- Quasi-Oracle Convergence Rate
- Structured Intervention Networks
- Experiments
- Summary

# Why is our decomposition useful\*?

\* *Main statement in Theorem 2 of paper.*

# Why is our decomposition useful\*?

$$\hat{f}(\mathbf{x}, \mathbf{t}) := \Psi(\mathbf{x})^T \Theta \Phi(\mathbf{t})$$

\* *Main statement in Theorem 2 of paper.*

# Why is our decomposition useful\*?

$$\hat{f}(\mathbf{x}, \mathbf{t}) := \Psi(\mathbf{x})^T \Theta \Phi(\mathbf{t})$$

$$f^*(\mathbf{x}, \mathbf{t}) := \mathbb{E}[Y | \mathbf{x}, \mathbf{t}]$$

\* *Main statement in Theorem 2 of paper.*

# Why is our decomposition useful\*?

$$\hat{f}(\cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}) := \Psi(\cdot_{\mathbf{x}})^T \Theta \Phi(\cdot_{\mathbf{t}})$$



$$f^*(\cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}) := \mathbb{E}[Y \mid \cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}]$$

\* *Main statement in Theorem 2 of paper.*

# Why is our decomposition useful\*?

$$\boxed{\begin{aligned}\hat{f}(\cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}) &:= \Psi(\cdot_{\mathbf{x}})^T \Theta \Phi(\cdot_{\mathbf{t}}) \\ &\downarrow \tilde{O}(n^{-\frac{1}{2(1+p)}}) \\ f^*(\cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}) &:= \mathbb{E}[Y \mid \cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}]\end{aligned}}$$

\* Main statement in Theorem 2 of paper.

# Why is our decomposition useful\*?

$$\hat{f}(\cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}) := \Psi(\cdot_{\mathbf{x}})^T \Theta \Phi(\cdot_{\mathbf{t}})$$

$$\downarrow \tilde{O}(n^{-\frac{1}{2(1+p)}})$$

$$f^*(\cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}) := \mathbb{E}[Y \mid \cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}]$$

$$\hat{m}(\cdot_{\mathbf{x}}) \rightarrow m(\cdot_{\mathbf{x}})$$

$$\hat{e}^h(\cdot_{\mathbf{x}}) \rightarrow e^h(\cdot_{\mathbf{x}})$$

\* Main statement in Theorem 2 of paper.

# Why is our decomposition useful\*?

$$\hat{f}(\cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}) := \Psi(\cdot_{\mathbf{x}})^T \Theta \Phi(\cdot_{\mathbf{t}})$$

$$\downarrow \tilde{O}(n^{-\frac{1}{2(1+p)}})$$

$$f^*(\cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}) := \mathbb{E}[Y \mid \cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}]$$

$$\hat{m}(\cdot_{\mathbf{x}}) \xrightarrow{O(n^{-\frac{1}{4}})} m(\cdot_{\mathbf{x}})$$

$$\hat{e}^h(\cdot_{\mathbf{x}}) \xrightarrow{O(n^{-\frac{1}{4}})} e^h(\cdot_{\mathbf{x}})$$

\* Main statement in Theorem 2 of paper.

# Why is our decomposition useful\*?

$$\hat{f}(\cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}) := \Psi(\cdot_{\mathbf{x}})^T \Theta \Phi(\cdot_{\mathbf{t}})$$

$$\downarrow \tilde{O}(n^{-\frac{1}{2(1+p)}})$$

$$f^*(\cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}) := \mathbb{E}[Y \mid \cdot_{\mathbf{x}}, \cdot_{\mathbf{t}}]$$

$$\hat{m}(\cdot_{\mathbf{x}}) \xrightarrow{O(n^{-\frac{1}{4}})} m(\cdot_{\mathbf{x}})$$

$$\hat{e}^h(\cdot_{\mathbf{x}}) \xrightarrow{O(n^{-\frac{1}{4}})} e^h(\cdot_{\mathbf{x}})$$

$$\text{Overlap: } \mathcal{P}_{\Psi(X) \times \Phi(T)} > 0$$

\* Main statement in Theorem 2 of paper.

# What does this mean?

- The target or nuisance functions cannot converge faster than  $O(n^{-1/2})$

# What does this mean?

- The target or nuisance functions cannot converge faster than  $O(n^{-1/2})$
- Usually this rate caps the rate of the target function - see the discussion in e.g. *Chernozhukov et al., 2018* (Double Machine Learning)

# What does this mean?

- The target or nuisance functions cannot converge faster than  $O(n^{-1/2})$
- Usually this rate caps the rate of the target function - see the discussion in e.g. *Chernozhukov et al., 2018* (Double Machine Learning)

We show that in the fixed features setting, the target function converges at almost  $n^{-\frac{1}{2(1+p)}}$  rate as long as the nuisance functions converge at  $n^{-1/4}$  rate.

# Proof outline - Notation

The regret quantities:

$$R(\Theta) = L(\Theta) - L(\Theta^*)$$

$$\tilde{R}_n(\Theta) = \tilde{L}_n(\Theta) - \tilde{L}_n(\Theta^*)$$

$$\hat{R}_n(\Theta) = \hat{L}_n(\Theta) - \hat{L}_n(\Theta^*)$$

$$L(f_\Theta) = L(\Theta) = \mathbb{E} \left[ \left\{ (Y - m^*(\mathbf{X})) - \boldsymbol{\alpha}(\mathbf{X})^T \boldsymbol{\Theta} (\boldsymbol{\beta}(\mathbf{T}) - e^P(\mathbf{X})) \right\}^2 \right] \quad (39)$$

$$\tilde{L}_n(f_\Theta) = \tilde{L}_n(\Theta) = \sum_{l=1}^n \left[ \left\{ (Y - m^*(\mathbf{X}_l)) - \boldsymbol{\alpha}(\mathbf{X}_l)^T \boldsymbol{\Theta} (\boldsymbol{\beta}(\mathbf{T}_l) - e^P(\mathbf{X}_l)) \right\}^2 \right] \quad (40)$$

$$\hat{L}_n(f_\Theta) = \hat{L}_n(\Theta) = \sum_{l=1}^n \left[ \left\{ (Y - \hat{m}(\mathbf{X}_l)) - \boldsymbol{\alpha}(\mathbf{X}_l)^T \boldsymbol{\Theta} (\boldsymbol{\beta}(\mathbf{T}_l) - \hat{e}^P(\mathbf{X}_l)) \right\}^2 \right] \quad (41)$$

# Proof outline - oracle rate

The regret quantities:

$$\begin{aligned} R(\Theta) &= L(\Theta) - L(\Theta^*) \\ \tilde{R}_n(\Theta) &= \tilde{L}_n(\Theta) - \tilde{L}_n(\Theta^*) \\ \hat{R}_n(\Theta) &= \hat{L}_n(\Theta) - \hat{L}_n(\Theta^*) \end{aligned}$$

**Lemma 4.** Let  $\check{L}(f_\Theta \in \mathcal{H}_c)$  be a loss function, and  $\check{R}(f_\Theta; c) = \check{L}(f_\Theta) - \check{L}(f_{\Theta_c})$  be the associated  $c$ -regret. Suppose  $\rho(r)$  is a positive, continuous, increasing function. If,  $\forall 1 \leq c \leq C$  and some  $k > 1$ , the following inequality holds for all  $f_\Theta \in \mathcal{H}_c$ :

$$\frac{1}{k} \check{R}(f_\Theta; c) - \rho(c) \leq R(f_\Theta; c) \leq k \check{R}(f_\Theta; c) + \rho(c) \quad (45)$$

Then, writing  $\kappa_1 = 2k + \frac{1}{k}$  and  $\kappa_2 = 2k^2 + 3$ , any solution to the regularized minimization problem with  $\Lambda(c) \geq \rho(c)$ ,

$$f_{\check{\Theta}} \in \arg \min_{f_\Theta \in \mathcal{H}_C} \{\check{L}(f_\Theta) + \kappa_1 \Lambda(f_\Theta)_\mathcal{H}\} \quad (46)$$

also satisfied the following risk bound:

$$L(f_{\check{\Theta}}) \leq \inf_{f_\Theta \in \mathcal{H}_C} \{L(f_\Theta) + \kappa_2 \Lambda(f_\Theta)_\mathcal{H}\} \quad (47)$$

# Proof outline - oracle rate

Mendelson and Neeman (2010) for  $\tilde{R}$ :

$$\rho_n(c) = U(\epsilon) \left\{ 1 + \log(n) + \log(\log(c + e)) \right\} \left( \frac{(c + 1)^p \log(n)}{\sqrt{n}} \right)^{2/(1+p)} \quad (53)$$

With 53, Lemma 4 immediately implies that penalized regression over  $\mathcal{H}_C$  with the oracle loss function  $\tilde{L}_n(\cdot)$  and regularizer  $\kappa_1 \rho_n(c)$  satisfies the bound below with high probability:

$$R(\tilde{\Theta}_n) = L(\tilde{\Theta}_n) - L(\Theta^*) \leq \inf_{\Theta \in \mathcal{H}_C} \{L((\Theta) + \kappa_2 \rho_n(\|\Theta\|_{\mathcal{H}})\} - L(\Theta^*) \quad (54)$$

# Proof outline - oracle rate

Furthermore, Corollary 2.7 in [36] gives that for any  $1 < c < C$ ,

$$\inf_{\Theta \in \mathcal{H}_C} \{L(\Theta) + \kappa_2 \rho_n(\|\Theta\|_{\mathcal{H}})\} \leq L(\Theta^*) + \{L(\Theta_c^*) - L(\Theta^*)\} + \kappa_2 \rho_n(c) \quad (55)$$

Finally, note that for large enough  $c$ ,

$$\{L(\Theta_c^*) - L(\Theta^*)\} = 0, \quad (56)$$

so the error is dominated by  $\rho_n(c)$ , at

$$R(\tilde{\Theta}_n) = \mathcal{O}\left((\log(n))^{\frac{3+p}{1+p}} n^{-\frac{1}{1+p}}\right) = \tilde{\mathcal{O}}(n^{-\frac{1}{1+p}}), \quad (57)$$

where  $\tilde{\mathcal{O}}$  notation ignores the logarithmic factors.

# Proof outline - Bridging $\hat{R}$ and $\tilde{R}$

Overlap + Boundedness of  $Y$ :

$$\left| \hat{R}_n(\Theta; c) - \tilde{R}_n(\Theta; c) \right| \leq 0.125R(\Theta; c) + o(\rho_n(c)) \quad (119)$$

# Roadmap

- Motivation
- **Generalized Robinson Decomposition**
- Quasi-Oracle Convergence Rate
- Structured Intervention Networks
- Experiments
- Summary

# Recap: Generalized Robinson Decomposition

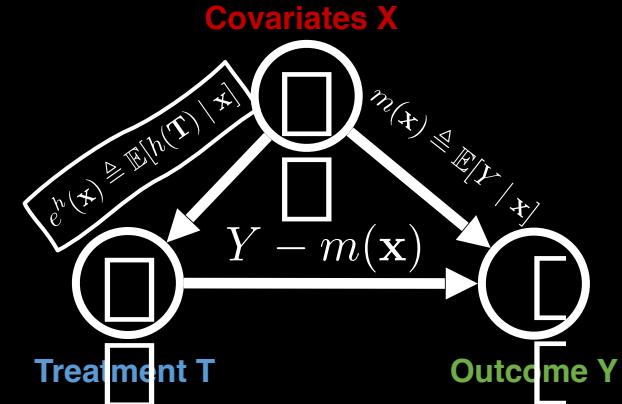
- Define *propensity features*

$$e^h(\mathbf{x}) \triangleq \mathbb{E}[h(\mathbf{T}) \mid \mathbf{x}]$$

- Recall  $m(\mathbf{x}) \triangleq \mathbb{E}[Y \mid \mathbf{x}] = g(\mathbf{x})^\top e^h(\mathbf{x})$

- Following the same steps as for the binary treatment case, we yield

$$Y - m(\mathbf{X}) = g(\mathbf{X})^\top (h(\mathbf{T}) - e^h(\mathbf{X})) + \varepsilon$$



Issue 2: How can we learn the relevant feature maps of the effect?

Solution: For a fixed  $h(\cdot)$ , a generalization to structured treatments is

$$\hat{g}(\cdot) = \arg \min_g \left\{ \frac{1}{n} \sum_{i=1}^n \left( Y_i - \hat{m}(\mathbf{X}_i) - g(\mathbf{X}_i)^\top (h(\mathbf{T}_i) - \hat{e}^h(\mathbf{X}_i)) \right)^2 \right\}$$

# Recap: Generalized Robinson Decomposition

How can we turn this into a practical learning algorithm?

Issue 2: How can we learn the relevant feature maps of the effect?

Solution: For a fixed  $h(\cdot)$ , a generalization to structured treatments is

$$\hat{g}(\cdot) = \arg \min_g \left\{ \frac{1}{n} \sum_{i=1}^n \left( Y_i - \hat{m}(\mathbf{X}_i) - g(\mathbf{X}_i)^\top (h(\mathbf{T}_i) - \hat{e}^h(\mathbf{X}_i)) \right)^2 \right\}$$

# Two-Stage Training Procedure

# Two-Stage Training Procedure

- **Stage 1:** Learn parameters of  $\hat{m}_{\theta}(\mathbf{X})$  based on objective

$$J_m(\boldsymbol{\theta}) = \sum_{i=1}^m (y_i - \hat{m}_{\boldsymbol{\theta}}(\mathbf{x}_i))^2$$

# Two-Stage Training Procedure

- Stage 1: Learn parameters of  $\hat{m}_{\theta}(\mathbf{X})$  based on objective

$$J_m(\boldsymbol{\theta}) = \sum_{i=1}^m (y_i - \hat{m}_{\boldsymbol{\theta}}(\mathbf{x}_i))^2$$

- Stage 2: Alternate between optimizing  $\hat{g}_{\psi}(\mathbf{X})$ ,  $\hat{h}_{\phi}(\mathbf{T})$  and  $\hat{e}_{\eta}^h(\mathbf{X})$

# Two-Stage Training Procedure

- Stage 1: Learn parameters of  $\hat{m}_{\boldsymbol{\theta}}(\mathbf{X})$  based on objective

$$J_m(\boldsymbol{\theta}) = \sum_{i=1}^m (y_i - \hat{m}_{\boldsymbol{\theta}}(\mathbf{x}_i))^2$$

- Stage 2: Alternate between optimizing  $\hat{g}_{\psi}(\mathbf{X}), \hat{h}_{\phi}(\mathbf{T})$  and  $\hat{e}_{\eta}^h(\mathbf{X})$

- a: Freeze  $\hat{m}_{\boldsymbol{\theta}}(\mathbf{X})$  and  $\hat{e}_{\eta}^h(\mathbf{X})$  to optimize  $\hat{g}_{\psi}(\mathbf{X}), \hat{h}_{\phi}(\mathbf{T})$  based on

$$J_{g,h}(\boldsymbol{\phi}, \boldsymbol{\psi}) = \sum_{i=1}^n \left( y_i - \left\{ \hat{m}_{\boldsymbol{\theta}}(\mathbf{x}_i) + \hat{g}_{\psi}(\mathbf{x}_i)^{\top} \left( \hat{h}_{\phi}(\mathbf{t}_i) - \hat{e}_{\eta}^h(\mathbf{x}_i) \right) \right\} \right)^2$$

# Two-Stage Training Procedure

- Stage 1: Learn parameters of  $\hat{m}_{\boldsymbol{\theta}}(\mathbf{X})$  based on objective

$$J_m(\boldsymbol{\theta}) = \sum_{i=1}^m (y_i - \hat{m}_{\boldsymbol{\theta}}(\mathbf{x}_i))^2$$

- Stage 2: Alternate between optimizing  $\hat{g}_{\psi}(\mathbf{X}), \hat{h}_{\phi}(\mathbf{T})$  and  $\hat{e}_{\eta}^h(\mathbf{X})$

- a: Freeze  $\hat{m}_{\boldsymbol{\theta}}(\mathbf{X})$  and  $\hat{e}_{\eta}^h(\mathbf{X})$  to optimize  $\hat{g}_{\psi}(\mathbf{X}), \hat{h}_{\phi}(\mathbf{T})$  based on

$$J_{g,h}(\boldsymbol{\phi}, \boldsymbol{\psi}) = \sum_{i=1}^n \left( y_i - \left\{ \hat{m}_{\boldsymbol{\theta}}(\mathbf{x}_i) + \hat{g}_{\psi}(\mathbf{x}_i)^{\top} \left( \hat{h}_{\phi}(\mathbf{t}_i) - \hat{e}_{\eta}^h(\mathbf{x}_i) \right) \right\} \right)^2$$

- b: Freeze  $\hat{m}_{\boldsymbol{\theta}}(\mathbf{X})$  and  $\hat{g}_{\psi}(\mathbf{X}), \hat{h}_{\phi}(\mathbf{T})$  optimize  $\hat{e}_{\eta}^h(\mathbf{X})$  based on

$$J_{e^h}(\boldsymbol{\eta}) = \sum_{i=1}^n \sum_{j=1}^d \left( \hat{h}_{\phi}(\mathbf{t}_i)^{(j)} - \hat{e}_{\eta}^h(\mathbf{x}_i)^{(j)} \right)^2$$

# Algorithm

---

a SIN Training.

---

**Input:** Stage 1 data  $\mathcal{D}_1 := \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ , Stage 2 data  $\mathcal{D}_2 := \{(\mathbf{x}_i, \mathbf{t}_i, y_i)\}_{i=1}^n$ . Step sizes  $\lambda_\theta, \lambda_\eta, \lambda_\psi, \lambda_\phi$ . Number of update steps  $K$ . Mini-batch sizes  $B_1, B_2$ .

```
1: Initialize parameters:  $\theta, \eta, \psi, \phi$ 
2: while not converged do ▷ Stage 1
3:   Sample mini-batch  $\{(\mathbf{x}_b, y_b)\}_{b=1}^{m_{B_1}}$ 
4:   Evaluate  $J_m(\theta)$ 
5:   Update  $\theta \leftarrow \theta - \lambda_\theta \hat{\nabla}_\theta J(\theta)$ 
6: end while
7: while not converged do ▷ Stage 2
8:   Sample mini-batch  $\{(\mathbf{x}_b, \mathbf{t}_b, y_b)\}_{b=1}^{m_{B_2}}$ 
9:   Evaluate  $J_{g,h}(\psi, \phi), J_{e^h}(\eta)$ 
10:  for  $k = 1$  to  $K$  do
11:    Update  $\phi \leftarrow \phi - \lambda_\phi \hat{\nabla}_\phi J_{g,h}(\psi, \phi)$ 
12:    Update  $\psi \leftarrow \psi - \lambda_\psi \hat{\nabla}_\psi J_{g,h}(\psi, \phi)$ 
13:  end for
14:  Update  $\eta \leftarrow \eta - \lambda_\eta \hat{\nabla}_\eta J_{e^h}(\eta)$ 
15: end while
```

---

# 280 character PyTorch-like code

```
# Initialize submodels and optimizers
m, e, g, h = MLP(), MLP(), MLP(), GNN()
m_opt, e_opt, g_opt, h_opt = Adam(m.params(), m_lr), Adam(e.params(), e_lr), ...

# Stage 1: Train m(x)
for batch in train_loader:
    X, Y = batch.X, batch.Y
    m_opt.zero_grad()
    F.mse_loss(m(X), Y).backward()
    m_opt.step()

# Stage 2: Train g(x), h(t), e(x)
for batch in train_loader:
    X, T, Y = batch.X, batch.T, batch.Y
    for _ in range(num_update_steps):
        g_opt.zero_grad()
        h_opt.zero_grad()
        F.mse_loss((g(X)*(h(T) - e(X))).sum(-1), (Y-m(X))).backward()
        g_opt.step()
        h_opt.step()
    e_opt.zero_grad()
    F.mse_loss(e(X), h(T)).backward()
    e_opt.step()
```

# Roadmap

- Motivation
- **Generalized Robinson Decomposition**
- Quasi-Oracle Convergence Rate
- Structured Intervention Networks
- Experiments
- Summary

# Experimental Setup

- **Data:** Two semi-synthetic datasets involving graph-treatments

## Small-World (SW)

**X:** Samples from multivar. uniform dist.

**T:** Watts–Strogatz small-world graphs

## The Cancer Genomic Atlas (TCGA)<sup>1</sup>

**X:** Gene expression data of cancer patients

1 | Data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>.

2 | L. Ruddigkeit, et al., Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17, 2012.

3 | Harada & Kashima, GraphITE: Estimating Individual Effects of Graph-structured Treatments, 2020.

# Experimental Setup

- **Data:** Two semi-synthetic datasets involving graph-treatments

## Small-World (SW)

**X:** Samples from multivar. uniform dist.

**T:** Watts–Strogatz small-world graphs

## The Cancer Genomic Atlas (TCGA)<sup>1</sup>

**X:** Gene expression data of cancer patients

- **Tasks:** Predicting in-sample/out-sample CATEs

1 | Data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>.

2 | L. Ruddigkeit, et al., Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17, 2012.

3 | Harada & Kashima, GraphITE: Estimating Individual Effects of Graph-structured Treatments, 2020.

# Experimental Setup

- **Data:** Two semi-synthetic datasets involving graph-treatments

## Small-World (SW)

**X:** Samples from multivar. uniform dist.

**T:** Watts–Strogatz small-world graphs

## The Cancer Genomic Atlas (TCGA)<sup>1</sup>

**X:** Gene expression data of cancer patients

- **Tasks:** Predicting in-sample/out-sample CATEs
- **Baselines:** GraphITE<sup>3</sup>, Vanilla Regression (GNN/CAT), Zero

1 | Data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>.

2 | L. Ruddigkeit, et al., Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17, 2012.

3 | Harada & Kashima, GraphITE: Estimating Individual Effects of Graph-structured Treatments, 2020.

# Experimental Setup

- **Data:** Two semi-synthetic datasets involving graph-treatments

## Small-World (SW)

**X:** Samples from multivar. uniform dist.

**T:** Watts–Strogatz small-world graphs

## The Cancer Genomic Atlas (TCGA)<sup>1</sup>

**X:** Gene expression data of cancer patients

- **Tasks:** Predicting in-sample/out-sample CATEs
- **Baselines:** GraphITE<sup>3</sup>, Vanilla Regression (GNN/CAT), Zero
- **Metric:** (Un-)Weighted expected Precision in Estimation of Het. Effects

$$\epsilon_{\text{UPEHE}}(\text{WPEHE}) \triangleq \int_{\mathcal{X}} (\hat{\tau}(\mathbf{t}', \mathbf{t}, \mathbf{x}) - \tau(\mathbf{t}', \mathbf{t}, \mathbf{x}))^2 p(\mathbf{t} \mid \mathbf{x}) p(\mathbf{t}' \mid \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

1 | Data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>.

2 | L. Ruddigkeit, et al., Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17, 2012.

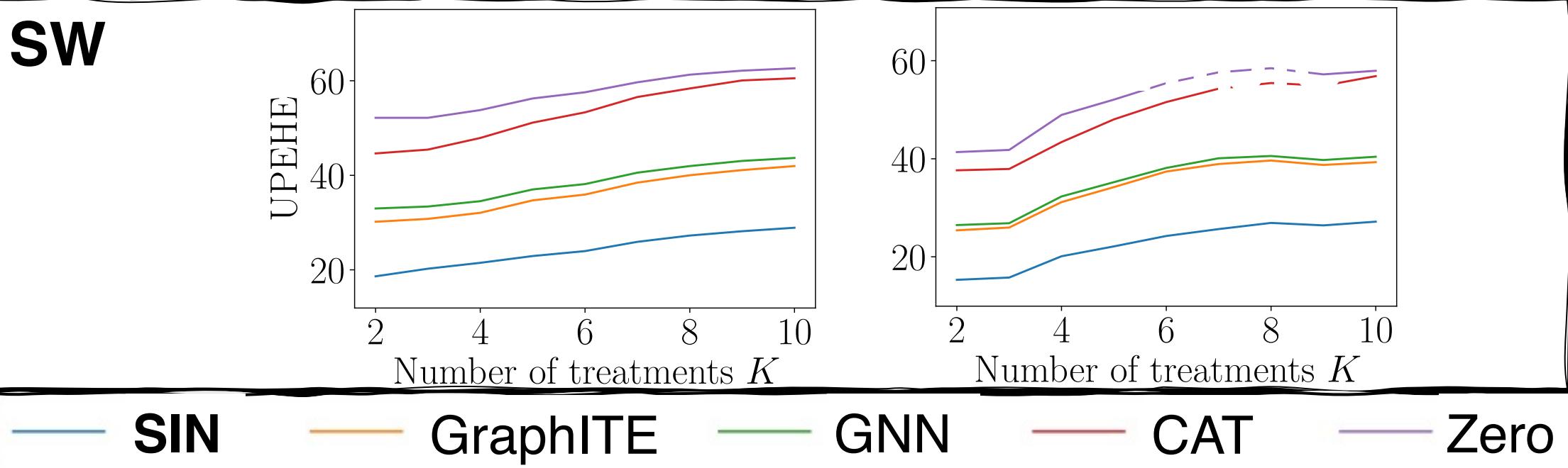
3 | Harada & Kashima, GraphITE: Estimating Individual Effects of Graph-structured Treatments, 2020.

# Results: In-Sample Out-Sample

# Results:

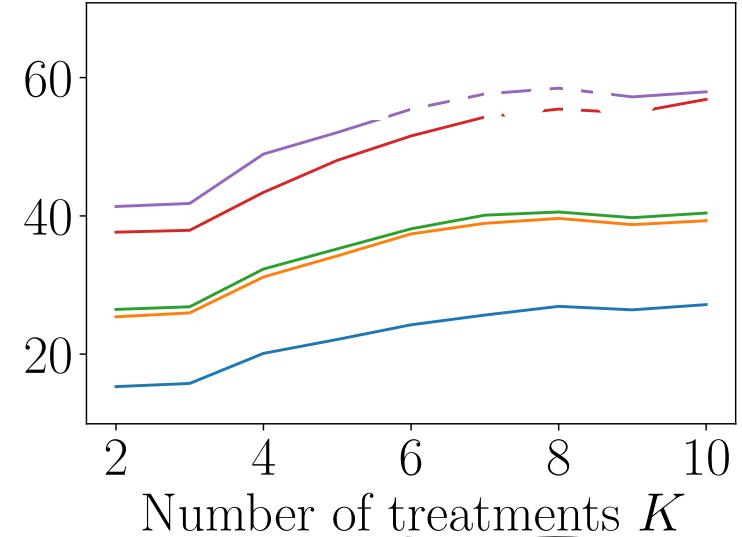
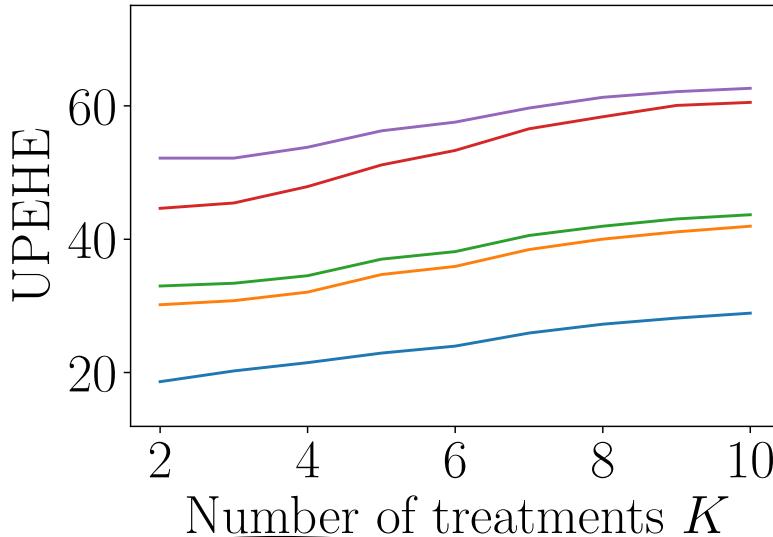
## In-Sample

## Out-

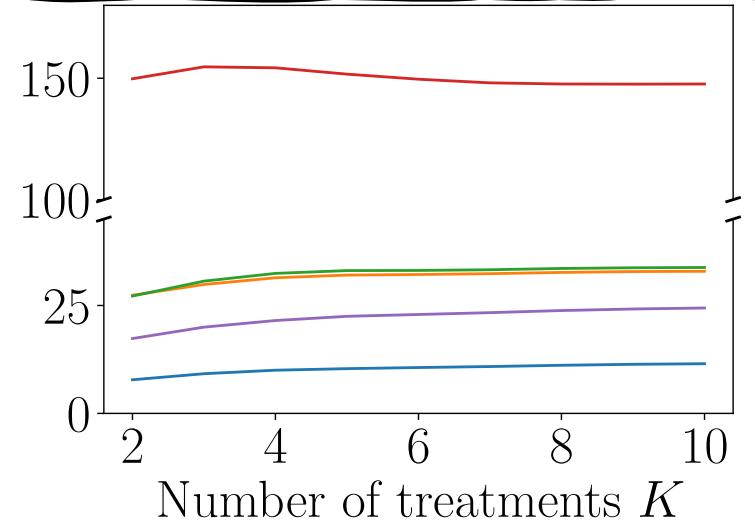
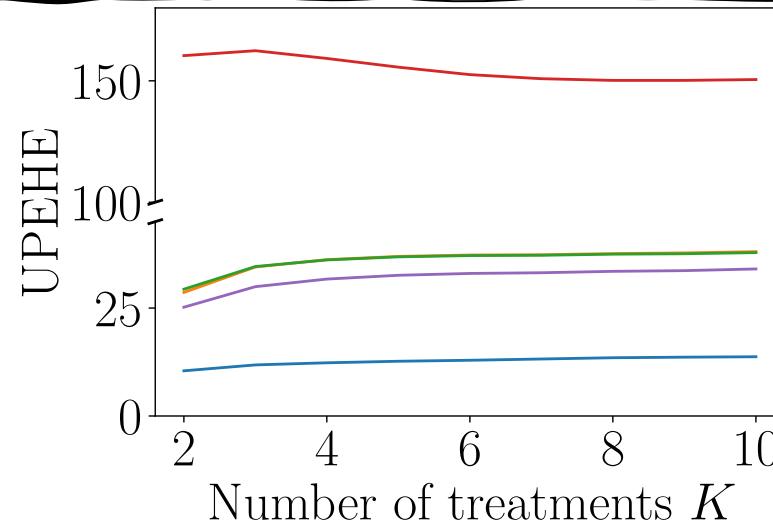


# Results: In-Sample Out-

**SW**



**TCGA**



**SIN**

**GraphITE**

**GNN**

**CAT**

**Zero**

# WPEHE for most likely K=6 treatments

| Method   | SW               |                  | TCGA               |                    |
|----------|------------------|------------------|--------------------|--------------------|
|          | In-sample        | Out-sample       | In-sample          | Out-sample         |
| Zero     | 56.26 $\pm$ 8.12 | 53.77 $\pm$ 8.93 | 26.63 $\pm$ 7.55   | 17.94 $\pm$ 4.86   |
| CAT      | 51.75 $\pm$ 8.85 | 49.76 $\pm$ 9.73 | 155.88 $\pm$ 52.82 | 146.62 $\pm$ 42.32 |
| GNN      | 37.10 $\pm$ 6.84 | 36.74 $\pm$ 7.42 | 30.67 $\pm$ 8.29   | 27.57 $\pm$ 7.95   |
| GraphITE | 34.81 $\pm$ 6.70 | 35.94 $\pm$ 8.07 | 30.31 $\pm$ 8.96   | 27.48 $\pm$ 8.95   |
| SIN      | 23.00 $\pm$ 4.56 | 23.19 $\pm$ 5.56 | 10.98 $\pm$ 3.45   | 8.15 $\pm$ 1.46    |

# Roadmap

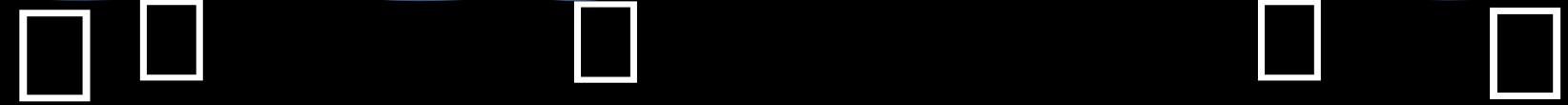
- Motivation
- **Generalized Robinson Decomposition**
- Quasi-Oracle Convergence Rate
- Structured Intervention Networks
- Experiments
- Summary

# Summary

## CATE Structured treatments (e.g. graphs)

- What?

$$\tau(\square, \square, \square) = \mathbb{E}[Y | \square, \text{do}(\square)] - \mathbb{E}[Y | \square, \text{do}(\square)]$$

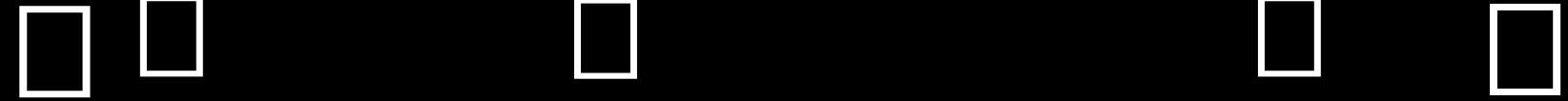


# Summary

## CATE Structured treatments (e.g. graphs)

- What?

$$\tau(\quad, \quad, \quad) = \mathbb{E}[Y | \quad, \text{do}(\quad)] - \mathbb{E}[Y | \quad, \text{do}(\quad)]$$



- Why?

1. Data-Efficiency
2. (Infinitely-)Many-Treatments-Settings
3. Generalization to unseen treatments

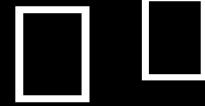
# Summary

## CATE Structured treatments (e.g. graphs)



- What?

$$\tau(\text{ }, \text{ }, \text{ }) = \mathbb{E}[Y | \text{ }, \text{do}(\text{ })] - \mathbb{E}[Y | \text{ }, \text{do}(\text{ })]$$



- Why?

1. Data-Efficiency

2. (Infinitely-)Many-Treatments-Settings

3. Generalization to unseen treatments

- How?

### Generalized Robinson Decomposition

$$Y - m(\mathbf{X}) = g(\mathbf{X})^\top (h(\mathbf{T}) - e^h(\mathbf{X})) + \varepsilon$$