



INSTITUTE OF ARTIFICIAL  
INTELLIGENCE (AI) IN MANAGEMENT



Munich Center for Machine Learning

# Bounds on Representation-Induced Confounding Bias for Treatment Effect Estimation

Valentyn Melnychuk, Dennis Frauen, Stefan Feuerriegel

LMU Munich & Munich Center for Machine Learning (MCML), Munich, Germany

ICLR 2024, Spotlight Presentation



# Introduction: Representation learning for CATE estimation

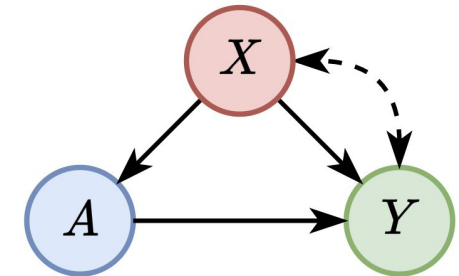
## Why this is important?

- State-of-the-art methods for conditional average treatment effect (CATE) estimation make widespread use of representation learning
- Low-dimensional (potentially constrained) representations reduce the variance, but, at the same time lose information about covariates, including information about confounders

## Problem formulation: representation-based CATE estimation

Given i.i.d. observational dataset  $\mathcal{D} = \{X_i, A_i, Y_i\}_{i=1}^n \sim \mathbb{P}(X, A, Y)$

- $X$  covariates
- $A$  binary treatments
- $Y$  continuous (factual) outcomes



Representation learning methods estimate the **conditional average treatment effect (CATE)**

$$\tau^x(x) = \mathbb{E}(Y[1] - Y[0] \mid X = x)$$

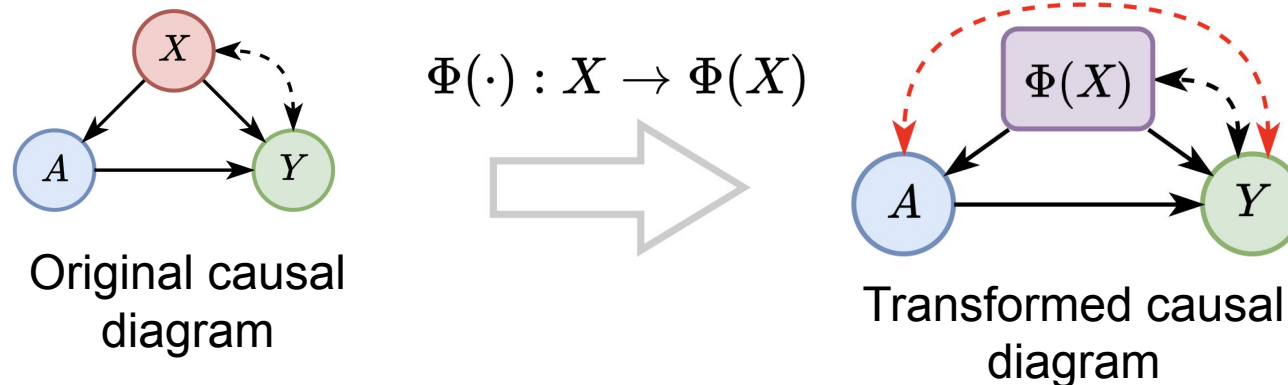
by (1) learning a low-dimensional (potentially constrained) representation  $\Phi(\cdot) : X \rightarrow \Phi(X)$  and by (2) estimating CATE wrt. representations

$$\tau^\phi(\phi) = \mathbb{E}(Y[1] - Y[0] \mid \Phi(X) = \phi) = \mu_1^\phi(\phi) - \mu_0^\phi(\phi) \quad \mu_a^\phi(\phi) = \mathbb{E}(Y \mid A = a, \Phi(X) = \phi)$$

# Introduction: Representation-induced confounding bias

- Constraints on the low-dimensional representations include:
  - treatment balancing with a probability metric:  $\text{dist} [\mathbb{P}(\Phi(X) \mid A = 0), \mathbb{P}(\Phi(X) \mid A = 1)] \approx 0$ .
  - invertibility:  $\Phi^{-1}(\Phi(X)) \approx X$ .
- Such low-dimensional representations can lead to a **representation-induced confounding bias (RICB)**, which we want to estimate / bound

**Problem formulation: representation-induced confounding bias**



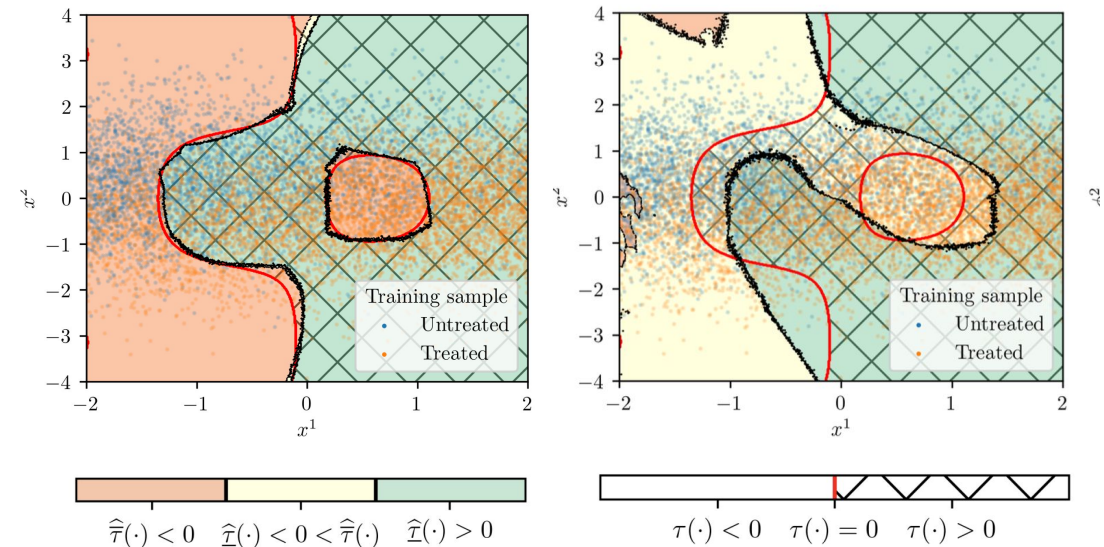
# Introduction: Research gap – Our contributions

## Research gap

- No work has studied the confounding bias (RICB) in low-dimensional (constrained) representations for CATE estimation

## Our contributions

- We formalize the representation-induced confounding bias (RICB)
- We propose a neural framework for estimating bounds based on the **Marginal Sensitivity Model**, which can be seen as a **refutation method** for representation learning CATE estimators
- We show that the estimated bounds are highly effective for the CATE-based decision-making



# Representation learning for CATE estimation: Assumptions

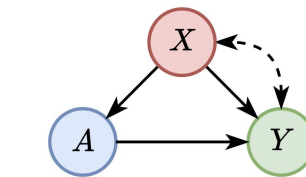
## Identifiability assumptions

- Potential outcomes framework (Neuman-Rubin):
  - **(i) Consistency.** If  $A = a$  is a treatment for some patient, then  $Y = Y[a]$
  - **(ii) Positivity (Overlap).** There is always a non-zero probability of receiving/not receiving any treatment, conditioning on the covariates:  $\epsilon > 0, \mathbb{P}(1 - \epsilon \geq \pi_a(X) \geq \epsilon) = 1$
  - **(iii) Exchangeability (Ignorability).** Current treatment is independent of the potential outcome, conditioning on the covariates  $A \perp\!\!\!\perp Y[a] \mid X$  for all  $a$ .
- Under assumptions (i)–(iii) CATE is identifiable

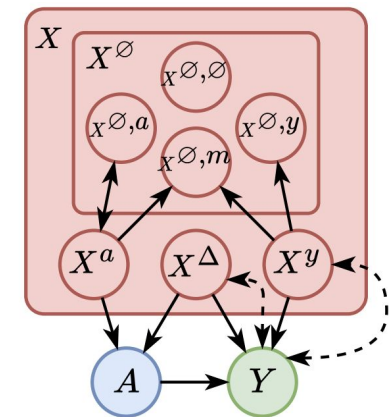
$$\tau^x(x) = \mu_1^x(x) - \mu_0^x(x) \qquad \mu_a^x(x) = \mathbb{E}(Y \mid A = a, X = x)$$

## Implicit partitioning assumption

- We assume an implicit partitioning (clustering) of  $X$  on  $\{X^\emptyset, X^a, X^y, X^\Delta\}$ 
  - (1) noise
  - (2) instruments
  - (3) outcome-predictive covariates
  - (4) confounders



Original causal diagram

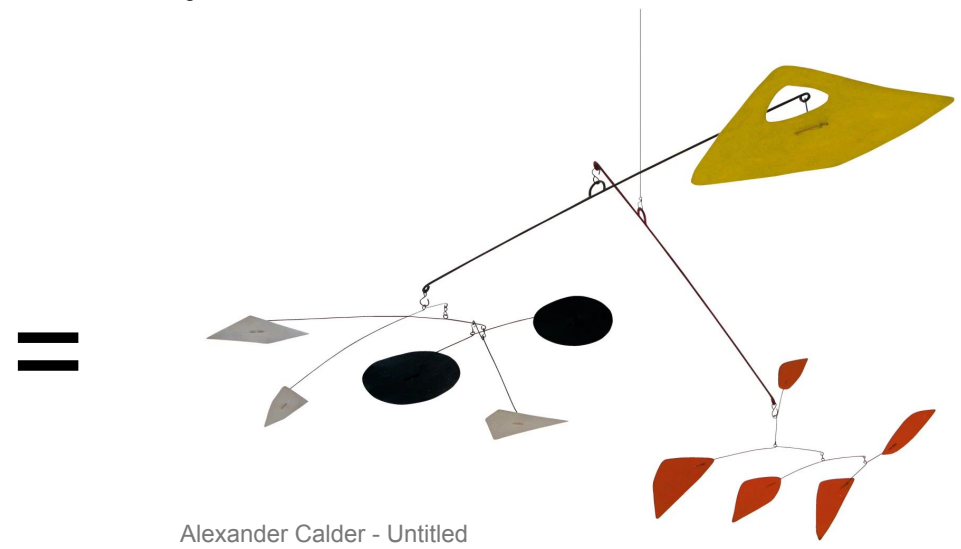
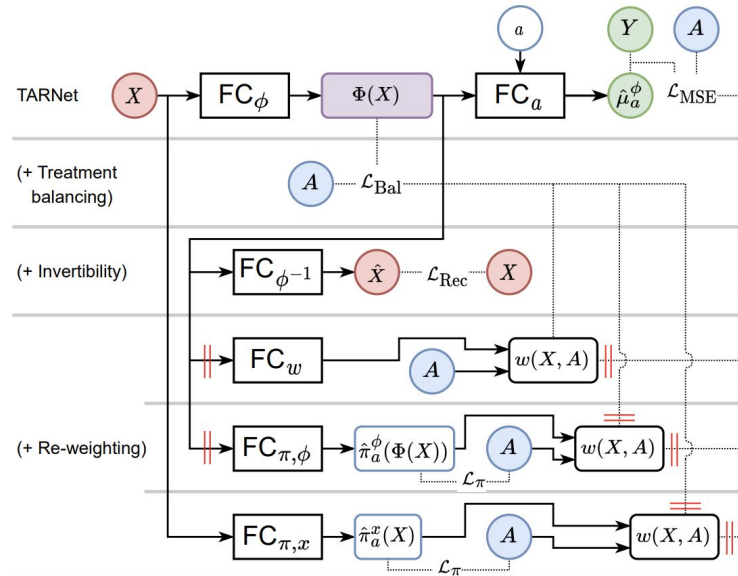


Clustered causal diagram

# Representation learning for CATE estimation: Methods

- Meta-learners (DR-learner, R-learner, etc.) can obtain the best **asymptotic performance and other properties** by fitting several models (nuisance functions and pseudo-outcome regression)
- Representation-based CATE estimators aim at **best-in-class estimation** with one model, but contain many trade-offs
- In low-sample regime, there is no universally best solution<sup>1</sup>

## Meta-learners vs. representation-based CATE estimators



<sup>1</sup> Alicia Curth and Mihaela van der Schaar. Nonparametric estimation of heterogeneous treatment effects: From theory to learning algorithms. In International Conference on Artificial Intelligence and Statistics, 2021.

# Types of representations: Valid representations

- We call a representation  $\Phi(\cdot)$  valid for CATE if it satisfies the following two equalities:

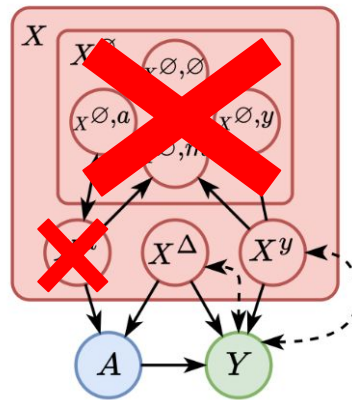
$$\tau^x(x) \stackrel{(i)}{=} \tau^\phi(\Phi(x)) \quad \text{and} \quad \tau^\phi(\phi) \stackrel{(ii)}{=} \mu_1^\phi(\phi) - \mu_0^\phi(\phi)$$

with  $\mu_a^\phi(\phi) = \mathbb{E}(Y | A = a, \Phi(X) = \phi)$

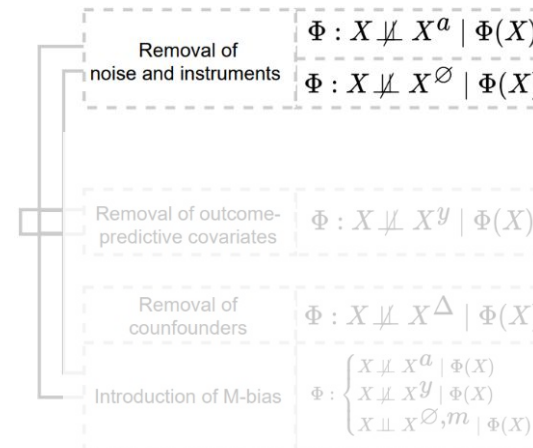
- Examples of valid representations:
  - Invertible representations (still help to reduce the variance when balanced)<sup>1</sup>
  - Removal of noise and instruments (achieved via balancing or lowering  $d_\phi$ )

## Valid representations

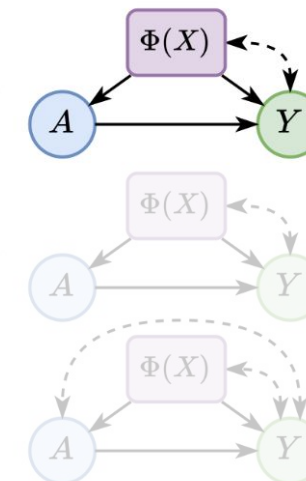
Original clustered causal diagram



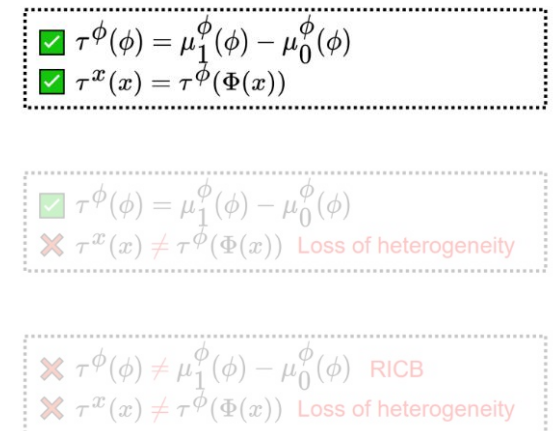
Constraints on  $\Phi(\cdot) : X \rightarrow \Phi(X)$



Transformed causal diagrams



Validity of  $\Phi(\cdot)$  for CATE



<sup>1</sup> Fredrik D. Johansson, Uri Shalit, Nathan Kallus, and David Sontag. Generalization bounds and representation learning for estimation of potential outcomes and causal effects. Journal of Machine Learning Research, 23:7489–7538, 2022.

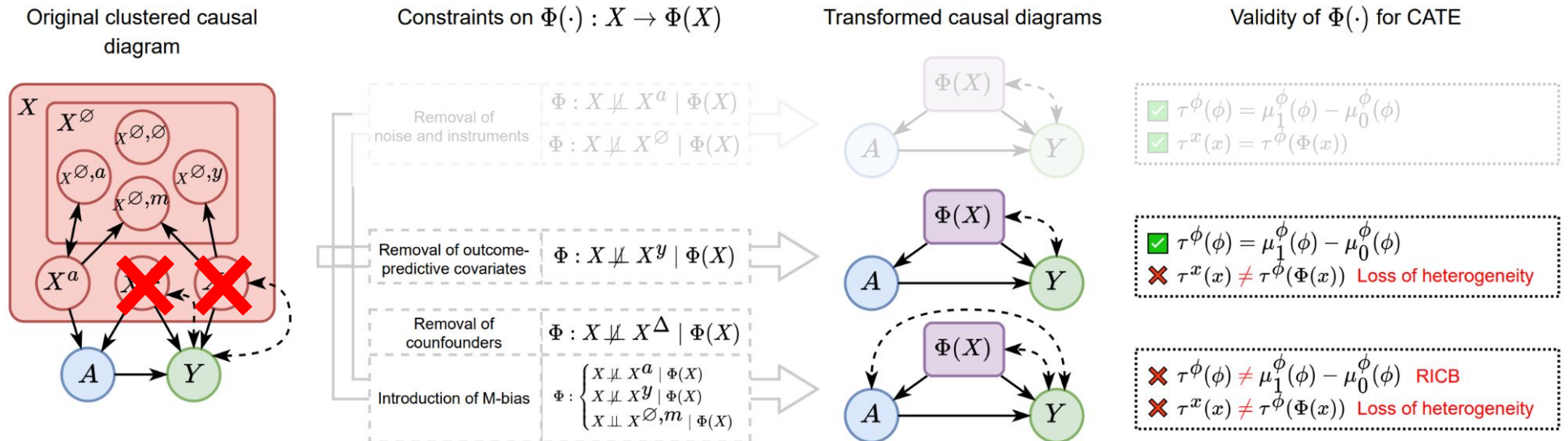
# Types of representations: Loss of heterogeneity

(i) **Loss of heterogeneity:** the treatment effect at the covariate (individual) level is different from the treatment effect at the representation (aggregated) level:

$$\tau^x(x) \neq \tau^\phi(\Phi(x))$$

- Happens whenever some information about  $X^\Delta$  or  $X^y$  is lost in the representation. E.g., propensity score is such a representation.
- Reasons: too low  $d_\phi$ , too large balancing

## Invalid representations





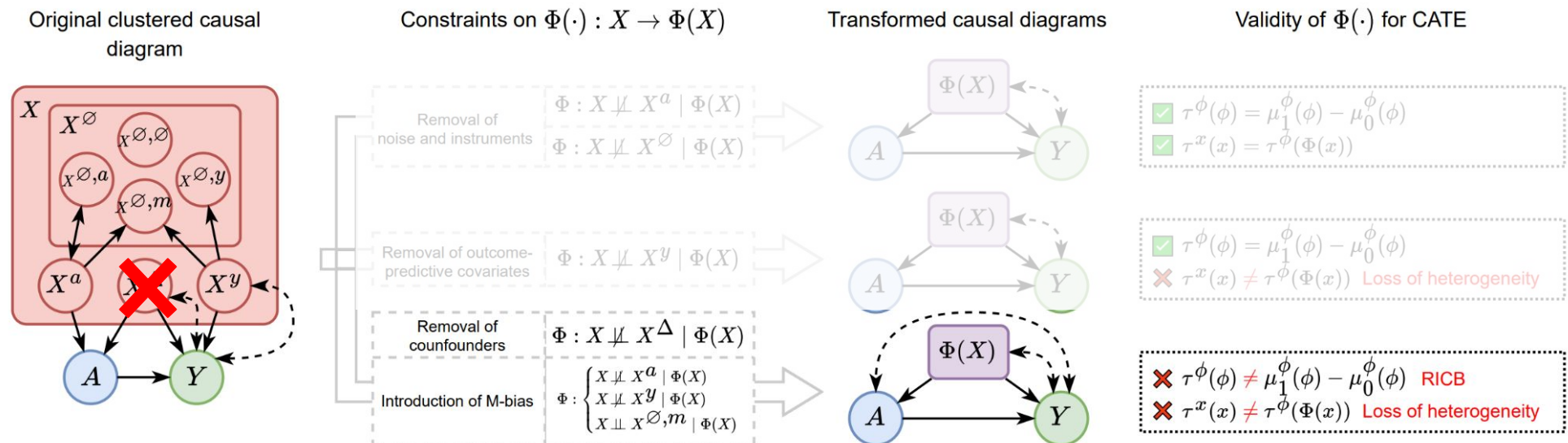
# Types of representations: RICB

(i) **Representation-induced confounding bias (RICB):** CATE wrt. representations is non-identifiable from observational data  $\mathbb{P}(\Phi(X), A, Y)$

$$\tau^\phi(\phi) \neq \mu_1^\phi(\phi) - \mu_0^\phi(\phi)$$

- Happens whenever some information about  $X^\Delta$  is lost in the representation or when M-bias is induced (this is rather a theoretic concept)
- Reasons: too low  $d_\phi$ , too large balancing

## Invalid representations



# Types of representations: Takeaways

## Takeaways

- The minimal sufficient and valid representation would aim to remove only the information about noise and instruments
- The loss of heterogeneity does not introduce bias but can only make CATE less individualized, namely, suitable only for subgroups
- The RICB automatically implies a loss of heterogeneity => We consider the RICB to be the main problem in representation learning methods for CATE
- RICB is an **infinite-sample confounding bias** (not a low-sample bias), present in the representations

# Partial identification of CATE under the RICB: Related work

Why is the direct inference of the RICB hard?

- Directly estimating RICB is (1) impractical and (2) intractable:

$$\tau^\phi(\Phi(x)) = \int_{\mathcal{X}_\Delta \times \mathcal{X}_Y} \tau^x(x) \mathbb{P}(X^\Delta = x^\Delta, X^y = x^y \mid \Phi(x)) dx^\Delta dx^y \neq \tau^x(x)$$

- The partitioning of  $X$  is unknown as well  $\{X^\emptyset, X^a, X^{\bar{y}}, X^\Delta\}$

Methods, affected by the RICB

Method	Invertibility	Balancing with	
		empirical probability metrics	loss re-weighting
TARNet (Shalit et al., 2017; Johansson et al., 2022)	–	–	–
BNN (Johansson et al., 2016); CFR (Shalit et al., 2017; Johansson et al., 2022); ESCFR (Wang et al., 2024)	–	IPM (MMD, WM)	–
RCFR (Johansson et al., 2018; 2022)	–	IPM (MMD, WM)	Learnable weights
DACPOL (Atan et al., 2018); CRN (Bica et al., 2020); ABCEI (Du et al., 2021); CT (Melnychuk et al., 2022); MitNet (Guo et al., 2023); BNCDE (Hess et al., 2024)	–	JSD (adversarial learning)	–
SITE (Yao et al., 2018)	Local similarity	Middle point distance	–
CFR-ISW (Hassanpour & Greiner, 2019a); DR-CFR (Hassanpour & Greiner, 2019b); DeR-CFR (Wu et al., 2022)	–	IPM (MMD, WM)	Representation propensity
DKLITE (Zhang et al., 2020)	Reconstruction loss	Counterfactual variance	–
BWCFR (Assaad et al., 2021)	–	IPM (MMD, WM)	Covariate propensity
PM (Schwab et al., 2018); StableCFR (Wu et al., 2023)	–	–	Upsampling via matching

IPM: integral probability metric; MMD: maximum mean discrepancy; WM: Wasserstein metric; JSD: Jensen-Shannon divergence

# Partial identification of CATE under the RICB: MSM

- Our idea is to employ a Marginal sensitivity model (MSM)<sup>1</sup> to perform the partial identification of the CATE (= find bounds on the RICB):

$$\Gamma(\phi)^{-1} \leq (\pi_0^\phi(\phi)/\pi_1^\phi(\phi)) (\pi_1^x(x)/\pi_0^x(x)) \leq \Gamma(\phi) \quad \text{for all } x \in \mathcal{X} \text{ s.t. } \Phi(x) = \phi,$$

where the sensitivity parameters can be **estimated from the combined data**  $\mathbb{P}(X, \Phi(X), A, Y)$

- Under the sensitivity constraint, the bounds on the RICB are given by

$$\underline{\tau}^\phi(\phi) = \underline{\mu}_1^\phi(\phi) - \overline{\mu}_0^\phi(\phi) \quad \text{and} \quad \overline{\tau}^\phi(\phi) = \overline{\mu}_1^\phi(\phi) - \underline{\mu}_0^\phi(\phi)$$

## Marginal sensitivity model

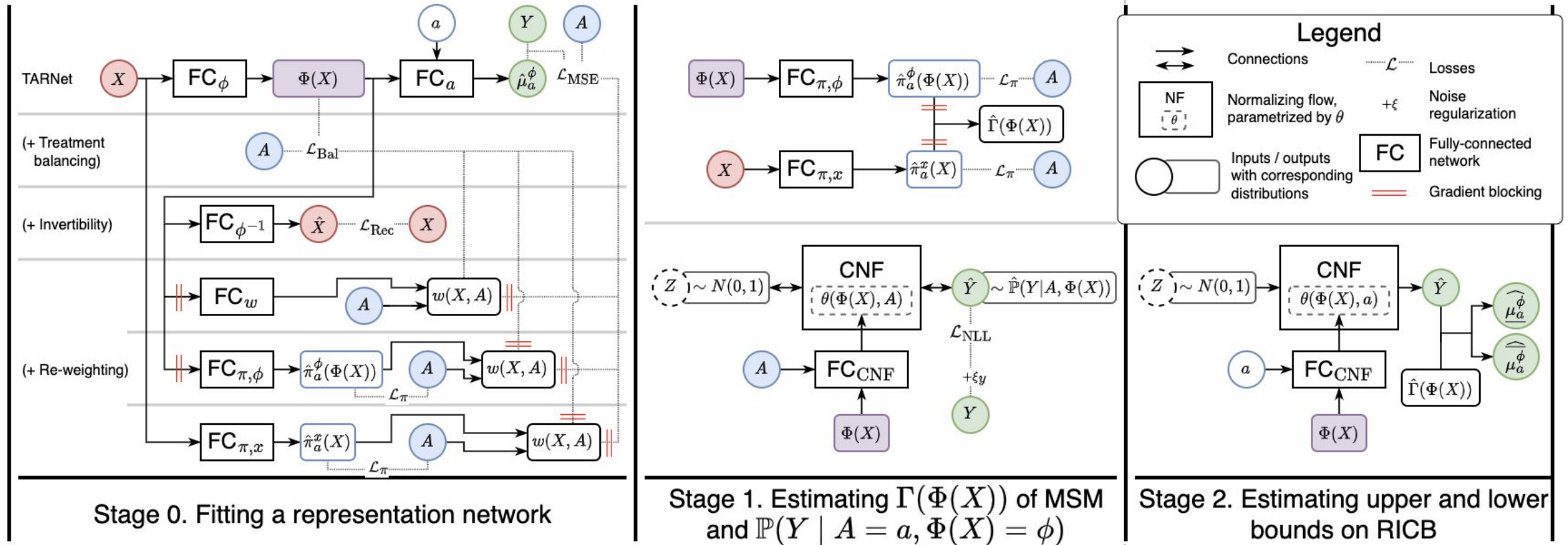
$$\underline{\mu}_a^\phi(\phi) = \frac{1}{s_-(a, \phi)} \int_{-\infty}^{\mathbb{F}^{-1}(c_-|a, \phi)} y \mathbb{P}(Y = y | a, \phi) dy + \frac{1}{s_+(a, \phi)} \int_{\mathbb{F}^{-1}(c_-|a, \phi)}^{+\infty} y \mathbb{P}(Y = y | a, \phi) dy,$$

$$\overline{\mu}_a^\phi(\phi) = \frac{1}{s_+(a, \phi)} \int_{-\infty}^{\mathbb{F}^{-1}(c_+|a, \phi)} y \mathbb{P}(Y = y | a, \phi) dy + \frac{1}{s_-(a, \phi)} \int_{\mathbb{F}^{-1}(c_+|a, \phi)}^{+\infty} y \mathbb{P}(Y = y | a, \phi) dy,$$

- The bounds are **valid** wrt. the original CATE and **sharp** wrt. the sensitivity constraint
- The bounds are still conservative, i.e., they do not distinguish instruments from confounders (but to do that we would need the original CATE)
- Yet, other sensitivity models, e.g., outcome sensitivity model, are impractical

<sup>1</sup> Zhiqiang Tan. A distributional approach for causal inference using propensity scores. Journal of the American Statistical Association, 101(476):1619–1637, 2006.

# Partial identification of CATE under the RICB: Neural framework



$\hat{\Gamma}(\phi_i)$  is a maximum over all  $\hat{\Gamma}(\Phi(x_j))$ , where  $\Phi(x_j)$  are the representations of the training sample in  $\delta$ -ball around  $\phi_i$ .  $\delta$  is the only hyper-parameter

# Experiments: Baselines – Evaluation – Datasets

- We evaluate our refutation framework together with SOTA representation-based CATE estimators: TARNet, BNN, CFR, InvTARNet, RCFR, CFR-ISW, BWCFR
- To compare our bounds with the point estimates, we employ an error rate of the policy (ER):
  - a policy based on the point estimate of the CATE applies a treatment whenever the CATE is positive:

## Baselines

$$\hat{\pi}(\phi) = \mathbb{1}\{\widehat{\tau}^{\phi}(\phi) > 0\}$$

## Evaluation

- a policy based on the bounds on the RICB has three decisions:

## Datasets

- (1) to treat  $\widehat{\tau}^{\phi}(\phi) > 0$
- (2) to do nothing  $\widehat{\tau}^{\phi}(\phi) < 0$
- (3) to defer a decision, otherwise
- We used 1 synthetic and 2 semi-synthetic datasets (IHDP100, HC-MNIST)

# Experiments: Results

- Our framework achieves clear improvements in the error rate among all the baselines, without deferring too many patients

## Results

$d_\phi$	ER <sub>out</sub> ( $\Delta$ ER <sub>out</sub> )	
	1	2
TARNet	30.79% (-12.89%)	9.82% (-3.73%)
BNN (MMD; $\alpha = 0.1$ )	34.32% (-15.41%)	16.15% (-4.19%)
CFR (MMD; $\alpha = 0.1$ )	35.01% (-14.27%)	11.92% (-5.54%)
CFR (MMD; $\alpha = 0.5$ )	35.79% (-11.43%)	17.89% (-7.27%)
CFR (WM; $\alpha = 1.0$ )	34.97% (-14.27%)	10.88% (-7.97%)
CFR (WM; $\alpha = 2.0$ )	35.18% (-13.63%)	13.19% (-6.28%)
InvTARNet	29.51% (-0.95%)	5.64% (-0.02%)
RCFR (WM; $\alpha = 1.0$ )	33.02% (-3.58%)	8.00% (-4.27%)
CFR-ISW (WM; $\alpha = 1.0$ )	35.00% (-9.43%)	7.27% (-1.86%)
BWCFR (WM; $\alpha = 1.0$ )	34.97% (-10.02%)	7.44% (-4.57%)

Classical CATE estimators:  $k$ -NN: 8.18%; BART: 17.37%; C-Forest: 16.10%

$d_\phi$	ER <sub>out</sub> ( $\Delta$ ER <sub>out</sub> )		
	7	39	78
TARNet	11.21% (-2.59%)	10.91% (-3.34%)	11.01% (-2.62%)
BNN (MMD; $\alpha = 0.1$ )	12.00% (-4.50%)	11.37% (-5.29%)	20.78% (-2.01%)
CFR (MMD; $\alpha = 0.1$ )	11.40% (-1.89%)	11.05% (-3.13%)	11.73% (-4.67%)
CFR (MMD; $\alpha = 0.5$ )	16.01% (+19.25%)	12.55% (-4.95%)	12.90% (-5.25%)
CFR (WM; $\alpha = 1.0$ )	24.55% (-10.42%)	27.87% (-10.18%)	31.19% (-11.53%)
CFR (WM; $\alpha = 2.0$ )	31.71% (-10.34%)	30.77% (-7.22%)	31.83% (-11.91%)
InvTARNet	12.18% (-1.29%)	11.38% (-3.98%)	11.55% (-4.34%)
RCFR (WM; $\alpha = 1.0$ )	21.51% (-9.17%)	26.97% (-6.17%)	30.14% (-14.26%)
CFR-ISW (WM; $\alpha = 1.0$ )	32.64% (-10.32%)	26.66% (-11.30%)	30.02% (-13.31%)
BWCFR (WM; $\alpha = 1.0$ )	13.62% (-3.96%)	28.18% (+0.24%)	32.54% (-6.75%)

Lower = better. Improvement over the baseline in green, worsening of the baseline in red

Classical CATE estimators:  $k$ -NN: 22.34%; BART: 17.51%; C-Forest: 17.65%

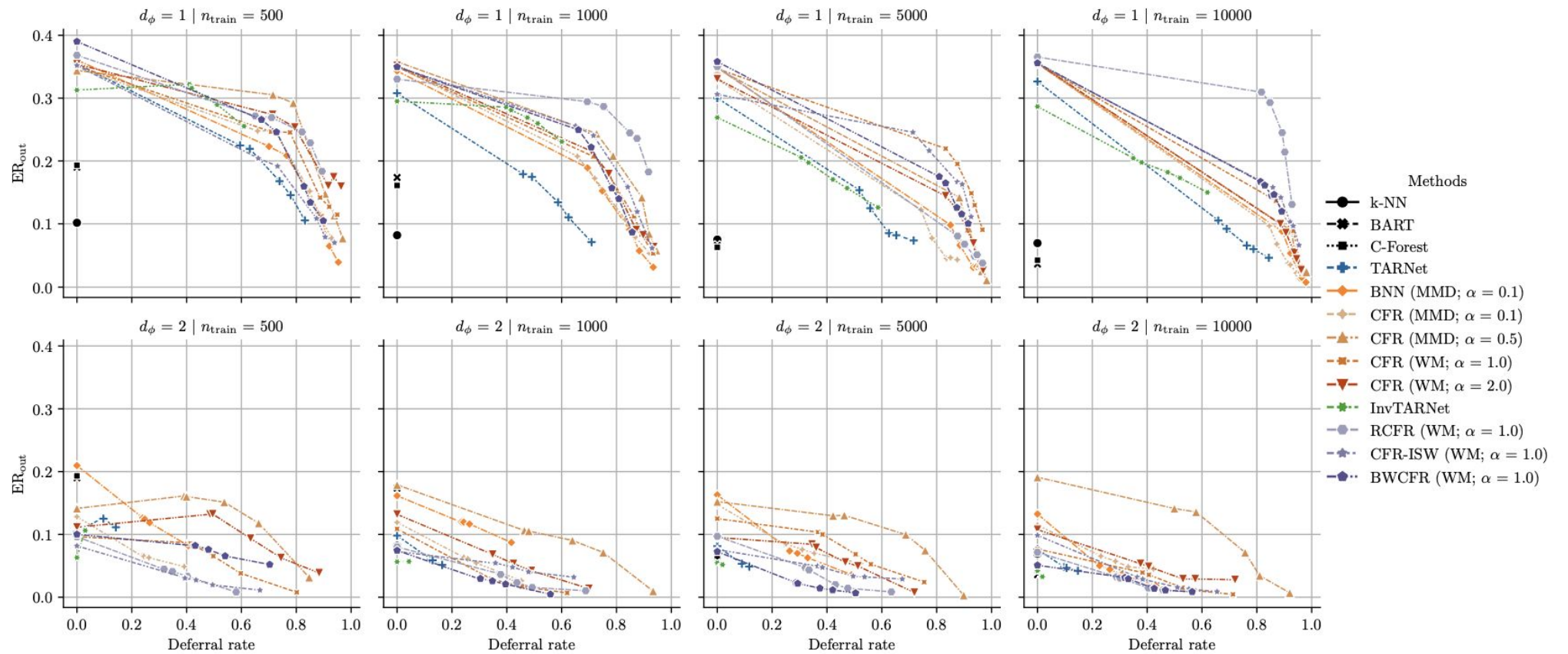
$d_\phi$	ER <sub>out</sub> ( $\Delta$ ER <sub>out</sub> )				
	5	10	15	20	25
TARNet	3.17% (-2.65%)	2.88% (-2.30%)	3.28% (-2.74%)	3.23% (-2.52%)	2.89% (-2.37%)
BNN (MMD; $\alpha = 0.1$ )	2.32% (-1.49%)	2.43% (-1.40%)	2.59% (-2.03%)	2.43% (-1.87%)	2.29% (-1.16%)
CFR (MMD; $\alpha = 0.1$ )	1.77% (-0.89%)	2.09% (-1.03%)	2.23% (-1.63%)	1.88% (-0.48%)	2.04% (-1.46%)
CFR (MMD; $\alpha = 0.5$ )	2.07% (-1.46%)	2.00% (+3.98%)	2.68% (+1.89%)	2.36% (+6.37%)	2.17% (+3.41%)
CFR (WM; $\alpha = 1.0$ )	1.93% (-0.89%)	1.75% (-0.25%)	1.83% (-1.24%)	1.83% (-0.49%)	1.80% (-0.20%)
CFR (WM; $\alpha = 2.0$ )	1.97% (-0.04%)	2.17% (-1.49%)	2.05% (-1.21%)	2.08% (-1.29%)	2.09% (-1.36%)
InvTARNet	2.52% (-1.95%)	3.11% (-2.47%)	2.99% (-2.51%)	2.79% (-2.41%)	2.83% (-2.28%)
RCFR (WM; $\alpha = 1.0$ )	3.36% (-2.84%)	3.45% (-1.52%)	2.67% (-1.57%)	4.69% (-3.83%)	1.95% (+1.06%)
CFR-ISW (WM; $\alpha = 1.0$ )	2.24% (-0.96%)	1.93% (-0.68%)	1.71% (-1.18%)	1.85% (-1.54%)	1.88% (-0.19%)
BWCFR (WM; $\alpha = 1.0$ )	3.57% (-1.49%)	3.52% (-2.16%)	3.88% (-1.10%)	3.80% (-2.38%)	4.07% (-1.18%)

Lower = better. Improvement over the baseline in green, worsening of the baseline in red

Classical CATE estimators:  $k$ -NN: 7.47%; BART: 5.07%; C-Forest: 6.28%

# Experiments: Results

- Ablation study on  $\delta$  shows, that the bounds remain valid under different values



Results



## Conclusion

We studied the validity of representation learning for CATE estimation. The validity may be violated due to low-dimensional representations as these introduce a **representation-induced confounding bias**.

As a remedy, we introduced a novel, **representation-agnostic refutation framework** that estimates bounds on the RICB and thus improves the reliability of their CATEs.



GitHub:  
[github.com/Valentyn1997/  
RICB](https://github.com/Valentyn1997/RICB)



ArXiv Paper:  
[arxiv.org/abs/2311.11321](https://arxiv.org/abs/2311.11321)