

# Robust and causal-oriented prediction models from heterogeneous data

Xinwei Shen<sup>†</sup>, Peter Bühlmann<sup>†</sup>, and Armeen Taeb<sup>‡</sup>

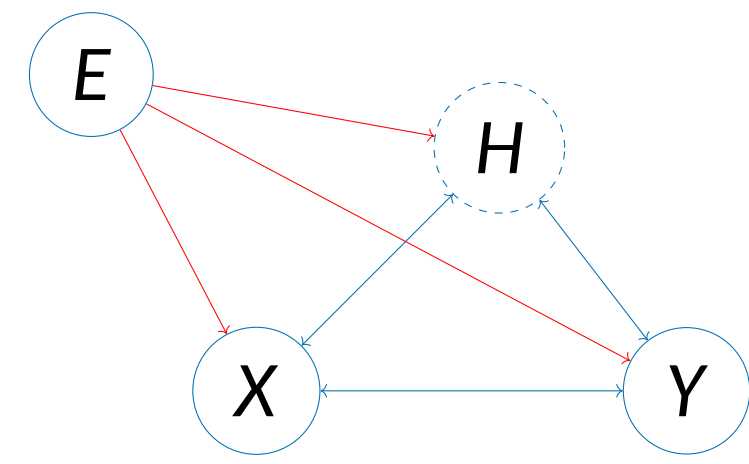
<sup>†</sup> Seminar for Statistics, ETH Zürich, <sup>‡</sup>Department of Statistics, University of Washington

## Introduction

- Heterogeneous data provides opportunities for causal inference and for learning prediction models that generalize to unseen environments.
- Perturbations may affect both means and variances of the variables, while previous methods only exploit shifts in the means.
- We propose Distributionally Robust predictions via Invariant Gradients (DRIG), a method that leverages perturbations in the form of both mean and variance shifts for robust predictions.
- Viewing causality as an extreme case of distributional robustness, we investigate the causal identifiability of DRIG under various scenarios of interventions and causal structures.

## Linear structural causal model

Covariates  $X \in \mathbb{R}^p$  and response variable  $Y \in \mathbb{R}$  with latent variables  $H$ .



- Training data from multiple environments  $e \in \mathcal{E}$ :

$$\begin{pmatrix} X^e \\ Y^e \end{pmatrix} = B^* \begin{pmatrix} X^e \\ Y^e \end{pmatrix} + \varepsilon + \delta^e, \quad (1)$$

where  $B^* := B_{p+1,p}^*$  denotes the causal effects and  $\varepsilon \perp\!\!\!\perp \delta^e$ .

- *Reference environment*:  $0 \in \mathcal{E}$  such that  $\sum_{e \in \mathcal{E}} \omega^e \mathbb{E}[\delta^e \delta^{eT}] \succeq \mathbb{E}[\delta^0 \delta^{0T}]$ . e.g., an observational environment with no intervention, i.e.,  $\delta^0 \equiv 0$ .
- Test distribution under new interventions:

$$\begin{pmatrix} X^v \\ Y^v \end{pmatrix} = B^* \begin{pmatrix} X^v \\ Y^v \end{pmatrix} + \varepsilon + v, \quad (2)$$

## Our method DRIG

Given a scalar  $\gamma \geq 0$ , the population DRIG:  $b_\gamma^{\text{opt}} = \operatorname{argmin}_b \mathcal{L}_\gamma(b)$  where

$$\mathcal{L}_\gamma(b) := \mathbb{E}[\ell(X^0, Y^0; b)] + \gamma \sum_{e \in \mathcal{E}} \omega^e (\mathbb{E}[\ell(X^e, Y^e; b)] - \mathbb{E}[\ell(X^0, Y^0; b)]), \quad (3)$$

where  $\ell(x, y; b) := (y - b^T x)^2$  and  $\gamma$  is a hyperparameter.

Special cases: interpolation between OLS and the causal parameter:

- $\gamma = 0$ : observational OLS
- $\gamma = 1$ : pooled OLS
- $\gamma \rightarrow \infty$ : causal parameter (when identifiable)
- $\delta^e$ 's are deterministic: anchor regression with categorical anchors

In the limit of  $\gamma \rightarrow \infty$ , the DRIG solution  $b_\gamma^{\text{opt}}$  satisfies gradient invariance:

## Definition (Gradient invariance)

A regression parameter  $b$  is said to satisfy the gradient invariance (GI) condition if  $\sum_{e \in \mathcal{E}} \omega^e \nabla_b \mathbb{E}[\ell(X^e, Y^e; b)] = \nabla_b \mathbb{E}[\ell(X^0, Y^0; b)]$ .

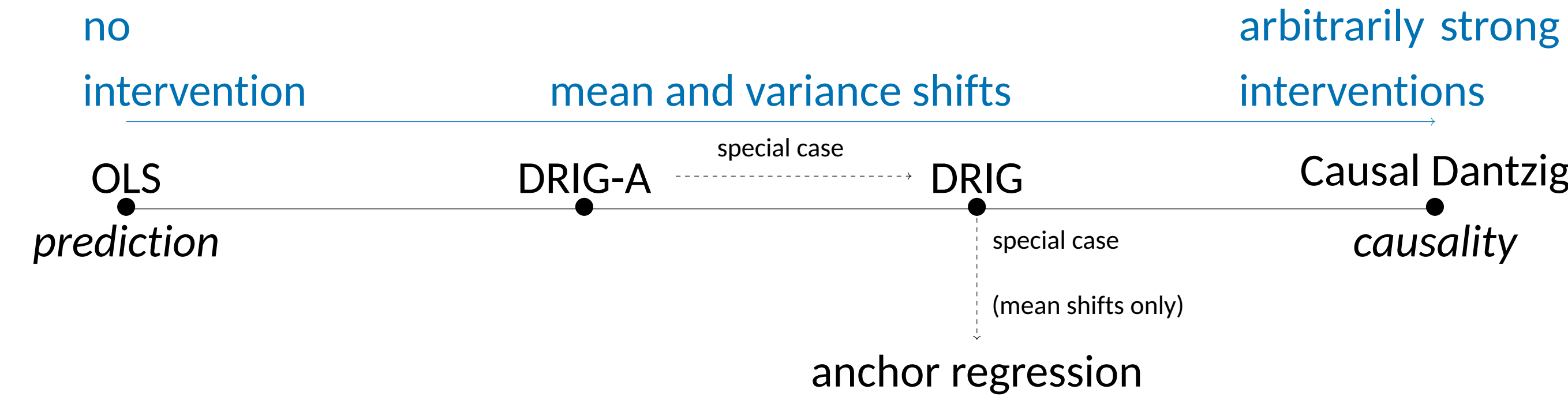


Figure: Trade-off between prediction and causality

## Distributional robustness

Minimizing the worst-case risk over test perturbations  $v \in \mathcal{C} \subseteq \mathbb{R}^{p+1}$ , i.e.

$$\operatorname{argmin}_{b \in \mathbb{R}^p} \sup_{v \in \mathcal{C}} \mathbb{E}_v[\ell(X^v, Y^v; b)]. \quad (4)$$

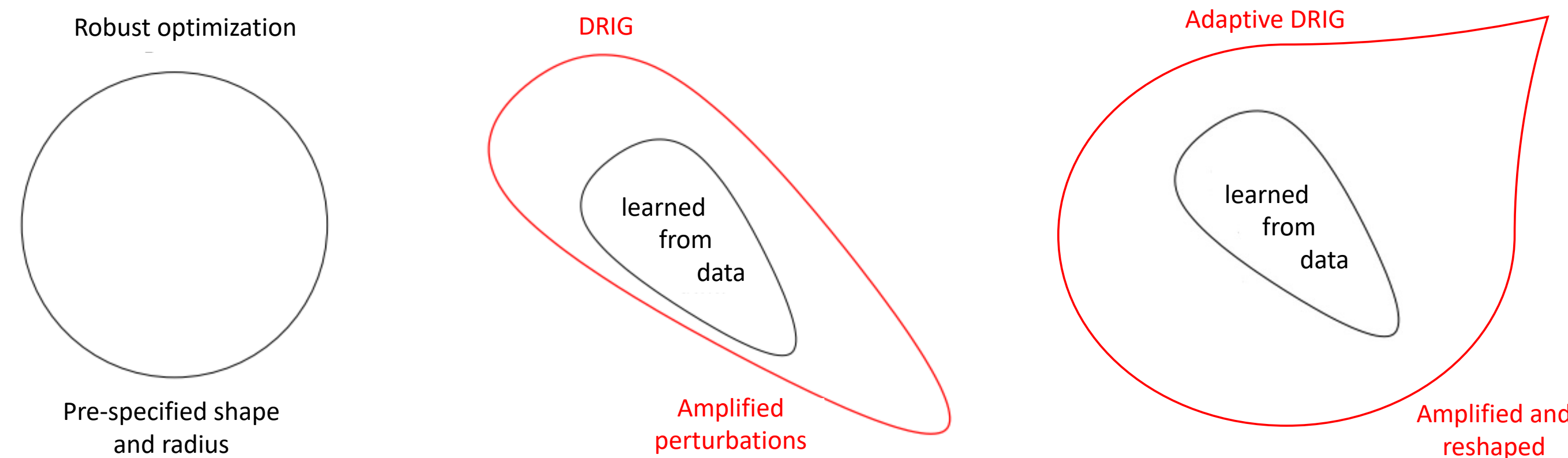
## Theorem (Robustness guarantee)

The prediction model by DRIG is the solution to the worst-case risk (4) with

$$\mathcal{C}_{\text{DRIG}}^\gamma := \left\{ v \in \mathbb{R}^{p+1} : \mathbb{E}[v v^T] \preceq \gamma \sum_{e \in \mathcal{E}} \omega^e (S^e - S^0 + \mu^e \mu^{eT}) \right\},$$

where  $S^e := \operatorname{cov}(\delta^e)$  and  $\mu^e := \mathbb{E}[\delta^e]$

- Strength of perturbations: controlled by  $\gamma$
  - Directions of perturbations: row and column spaces of  $\sum_{e \in \mathcal{E}} \omega^e (S^e - S^0 + \mu^e \mu^{eT})$
- Comparing to traditional distributionally robust optimization (DRO):



Comparing to other methods: for  $\gamma \geq 1$ , we have  $\mathcal{C}_{\text{OLS}} \subseteq \mathcal{C}_{\text{pOLS}} \subseteq \mathcal{C}_{\text{anchor}}^\gamma \subseteq \mathcal{C}_{\text{DRIG}}^\gamma$ .

- observational OLS:  $\mathcal{C}_{\text{OLS}} = \{v \in \mathbb{R}^{p+1} : \mathbb{E}[v v^T] \preceq \sum_{e \in \mathcal{E}} \omega^e (S^e - S^0 + \mu^e \mu^{eT})\}$
- pooled OLS:  $\mathcal{C}_{\text{pOLS}} = \{v \in \mathbb{R}^{p+1} : \mathbb{E}[v v^T] \preceq \sum_{e \in \mathcal{E}} \omega^e (S^e - S^0 + \mu^e \mu^{eT})\}$
- anchor regression:  $\mathcal{C}_{\text{anchor}}^\gamma = \{v \in \mathbb{R}^{p+1} : \mathbb{E}[v v^T] \preceq \sum_{e \in \mathcal{E}} \omega^e (S^e - S^0 + \gamma \mu^e \mu^{eT})\}$
- causal parameter:  $\mathcal{C}_{\text{causal}} = \{v \in \mathbb{R}^{p+1} : v_{p+1} \equiv 0\}$

## Causal identification

DRIG solution with  $\gamma \rightarrow \infty$ :

$$b_\infty^{\text{opt}} := \lim_{\gamma \rightarrow \infty} b_\gamma^{\text{opt}} = \operatorname{argmin}_{b \text{ satisfies GI}} \mathbb{E}[(Y^0 - b^T X^0)^2].$$

Identifiable cases:  $b_\infty^{\text{opt}} = b^*$

- Sufficient interventions on  $X$  & no interventions on  $Y$  or  $H$
  - Sufficient interventions on  $X$  & independent interventions on  $Y$  &  $Y$  is childless.
- Unidentifiable cases: approximate identifiability  $\|b_\infty^{\text{opt}} - b^*\| \leq c$
- Interventions on the latent variables with dense latent effects
  - Insufficient interventions on  $X$

## DRIG-A: adaptive DRIG in semi-supervised domain adaptation settings

- DRIG with matrix  $\Gamma$  for more flexible robustness:  $b_\Gamma^{\text{opt}}$  is minimizing 
$$\mathcal{L}_\Gamma(b) := \mathbb{E}[(Y^0 - b^T X^0)^2] + \sum_{e \in \mathcal{E}} \omega^e (\mathbb{E}[\gamma_y Y^e - b^T \Gamma_x X^e]^2 - \mathbb{E}[\gamma_y Y^0 - b^T \Gamma_x X^0]^2).$$

with a closed-form solution

$$b_\Gamma^{\text{opt}} = [\mathbb{E}X^0 X^{0T} + \Gamma_x \Delta_x \Gamma_x]^{-1} [\mathbb{E}X^0 Y^0 + \gamma_y \Gamma_x \Delta_{xy}].$$

where  $\Delta_x = \sum_{e \in \mathcal{E}} \omega^e [\mathbb{E}X^e X^{eT} - \mathbb{E}X^0 X^{0T}]$  and  $\Delta_{xy} = \sum_{e \in \mathcal{E}} \omega^e [\mathbb{E}X^e Y^e - \mathbb{E}X^0 Y^0]$ .

- Test distribution  $P_{\text{test}}$  according to SCM (2)

- a small labeled sample  $\{(X_i^Y, Y_i^Y) \sim P_{\text{test}}, i = 1, \dots, n_l\}$
- a large unlabeled test samples  $\{X_i^X \sim P_{\text{test}}, i = 1, \dots, n_u\}$ .

- Test OLS  $\hat{b}_{\text{OLS}} := (\frac{1}{n_u} \sum_{i=1}^{n_u} X_i^X X_i^{XT})^{-1} (\frac{1}{n_l} \sum_{i=1}^{n_l} X_i^X Y_i^Y)$

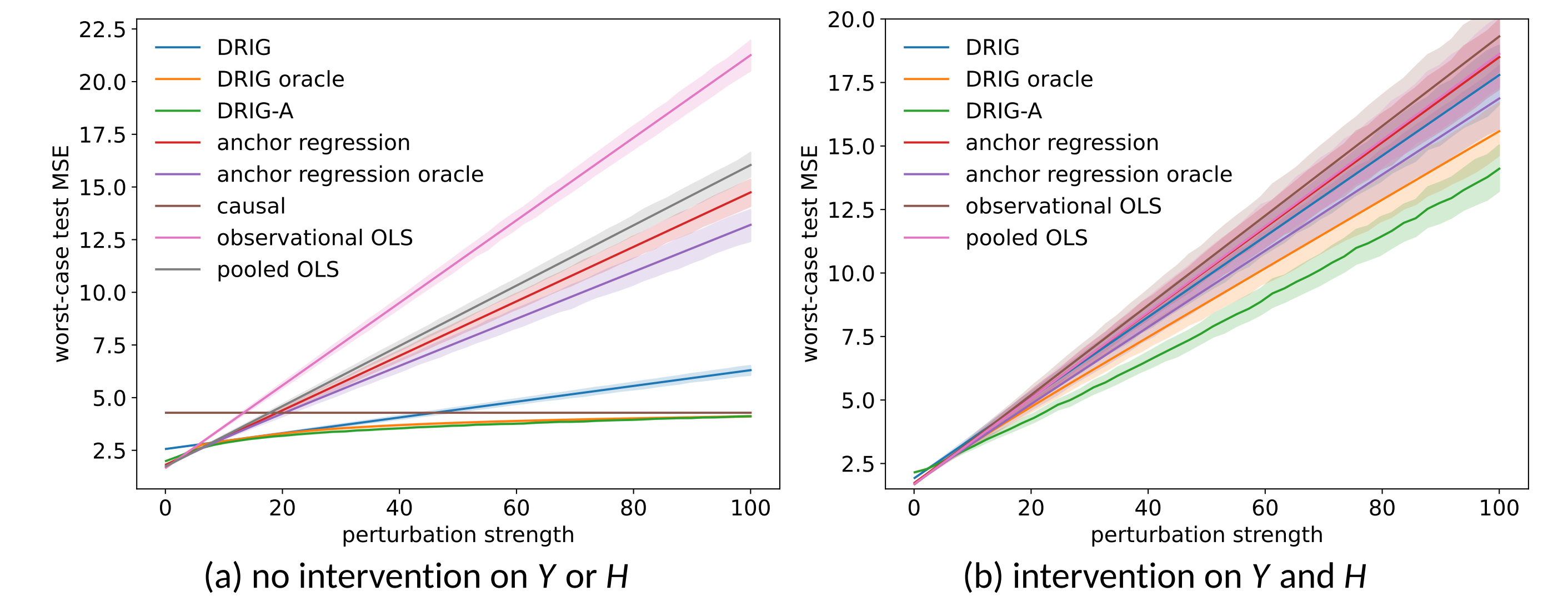
- Choosing  $\Gamma$  based on the semi-supervised test sample: in population

$$\begin{aligned} \min_{\Gamma_x, \gamma_y} \quad & \mathbb{E}[(Y^v - b_\Gamma^{\text{opt}T} X^v)^2] \\ \text{s.t.} \quad & \mathbb{E}X^0 X^{0T} + \Gamma_x \Delta_x \Gamma_x = \mathbb{E}X^v X^{vT} \end{aligned}$$

**Theorem.** Assume  $p > 1$ ;  $\operatorname{Var}(X^u Y^u) \succ (\mathbb{E}[X^u Y^u] - \mathbb{E}[X^0 Y^0])^{\otimes 2}$ . Then  $\exists N_u, N_l > 0$  such that when  $n_u \geq N_u$  and  $n_l \leq N_l$ , we have  $\mathbb{E}[\mathcal{L}_{\text{test}}(b_\Gamma^{\text{opt}})] < \mathbb{E}[\mathcal{L}_{\text{test}}(\hat{b}_{\text{OLS}})]$ , where the expectation is taken over all test samples.

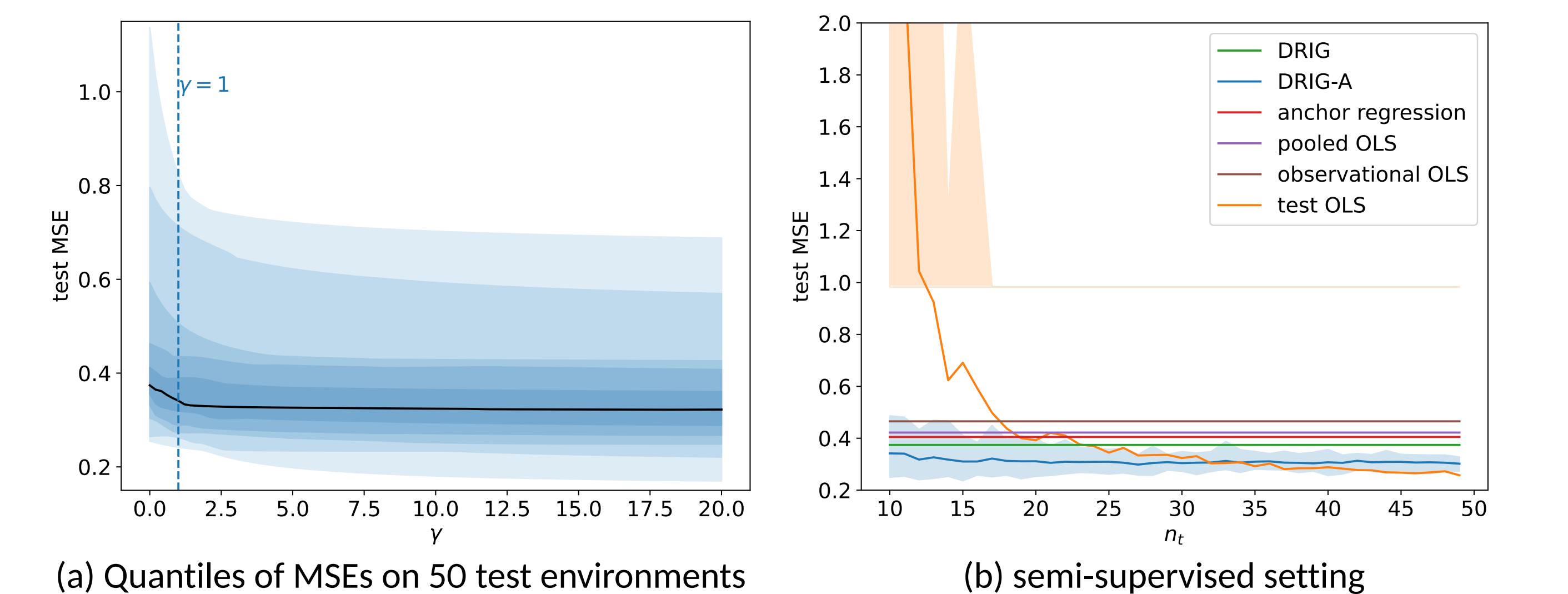
## Simulations

Worst-case MSEs over 20 randomly simulated test environments:



## Single-cell data

- 10 genes (1 response), 11,485 observational data, 10 interventional environments.
- Hundreds of test environments; on each of them, one hidden gene is intervened.



(a) Quantiles of MSEs on 50 test environments

(b) semi-supervised setting