

Achievable distributional robustness when the robust risk is only partially identified

Julia Kostin, Nicola Gnecco, Fanny Yang

ETH Zurich, University College London

ETH zürich



Team



Julia Kostin

ETH zürich



Nicola Gnecco



Fanny Yang

ETH zürich

Out-of-domain generalization



Training distribution $\mathbb{P}_{\text{train}}$

Test distribution \mathbb{P}_{test}

Out-of-domain generalization



$$\beta_{\text{train}} = \arg \min_{\beta} \mathcal{R}(\beta; \mathbb{P}_{\text{train}})$$

Low training risk

Training distribution $\mathbb{P}_{\text{train}}$

$$\mathcal{R}(\beta_{\text{train}}; \mathbb{P}_{\text{test}}) \gg \min_{\beta} \mathcal{R}(\beta; \mathbb{P}_{\text{test}})$$

High test risk

Test distribution \mathbb{P}_{test}



Distributional robustness

Goal: given training data, generalize to a set of feasible test distributions, called **robustness set**, by computing a minimiser of the **robust risk**

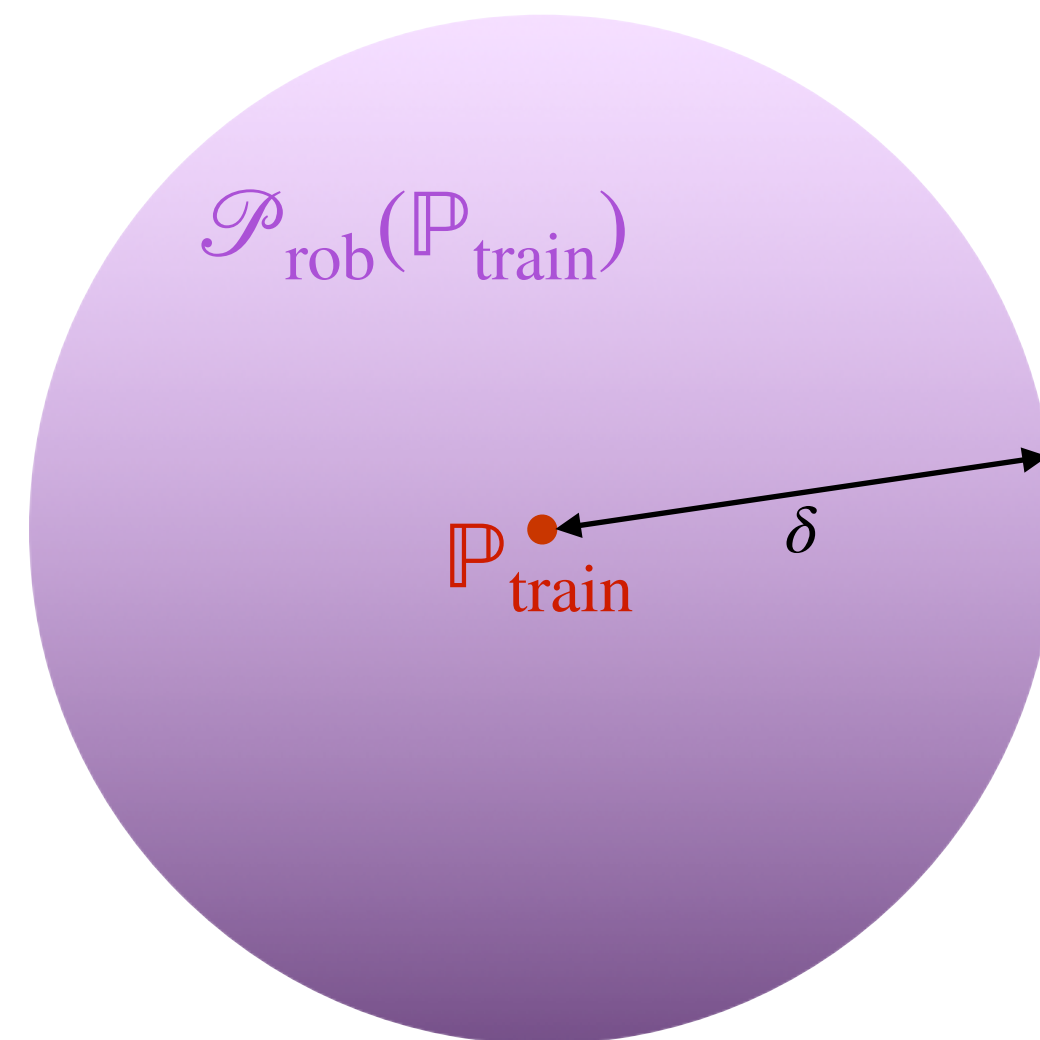
$$\beta_{\text{rob}} = \arg \min_{\beta} \left[\mathcal{R}_{\text{rob}}(\beta; \mathcal{P}_{\text{rob}}(\theta_{\star})) := \sup_{\mathbb{P} \in \mathcal{P}_{\text{rob}}(\theta_{\star})} \mathcal{R}(\beta; \mathbb{P}) \right]$$

In previously considered robustness scenarios, the parameters θ_{\star} and/or the robustness set $\mathcal{P}_{\text{rob}}(\theta_{\star})$ are considered to be **known**:

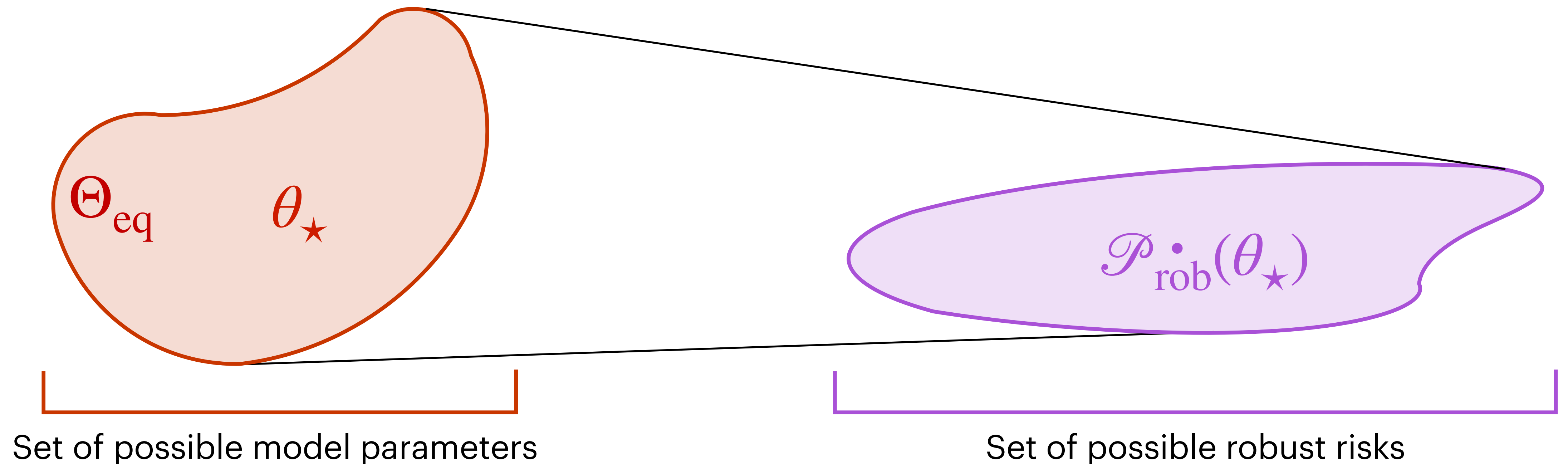
Distributionally robust optimization

$$\theta_{\star} = \mathbb{P}_{\text{train}}$$

$$\mathcal{P}_{\text{rob}}(\theta_{\star}) = \{\mathbb{P} : D(\mathbb{P}, \mathbb{P}_{\text{train}}) \leq \delta\}$$



Often, θ_{\star} and/or $\mathcal{P}_{\text{rob}}(\theta_{\star})$ are neither known nor computable from training data



Instead, they can be merely set identified.

**We propose to minimise a new objective called the
identifiable robust risk:**

$$\mathcal{R}_{\text{rob,ID}}(\beta; \Theta_{\text{eq}}) := \sup_{\theta \in \Theta_{\text{eq}}} \sup_{\mathbb{P} \in \mathcal{P}_{\text{rob}}(\theta)} \mathcal{R}(\beta, \mathbb{P})$$

Best achievable distributional robustness:

$$\mathfrak{M}(\Theta_{\text{eq}}) = \inf_{\beta \in \mathbb{R}^d} \mathcal{R}_{\text{rob,ID}}(\beta; \Theta_{\text{eq}})$$

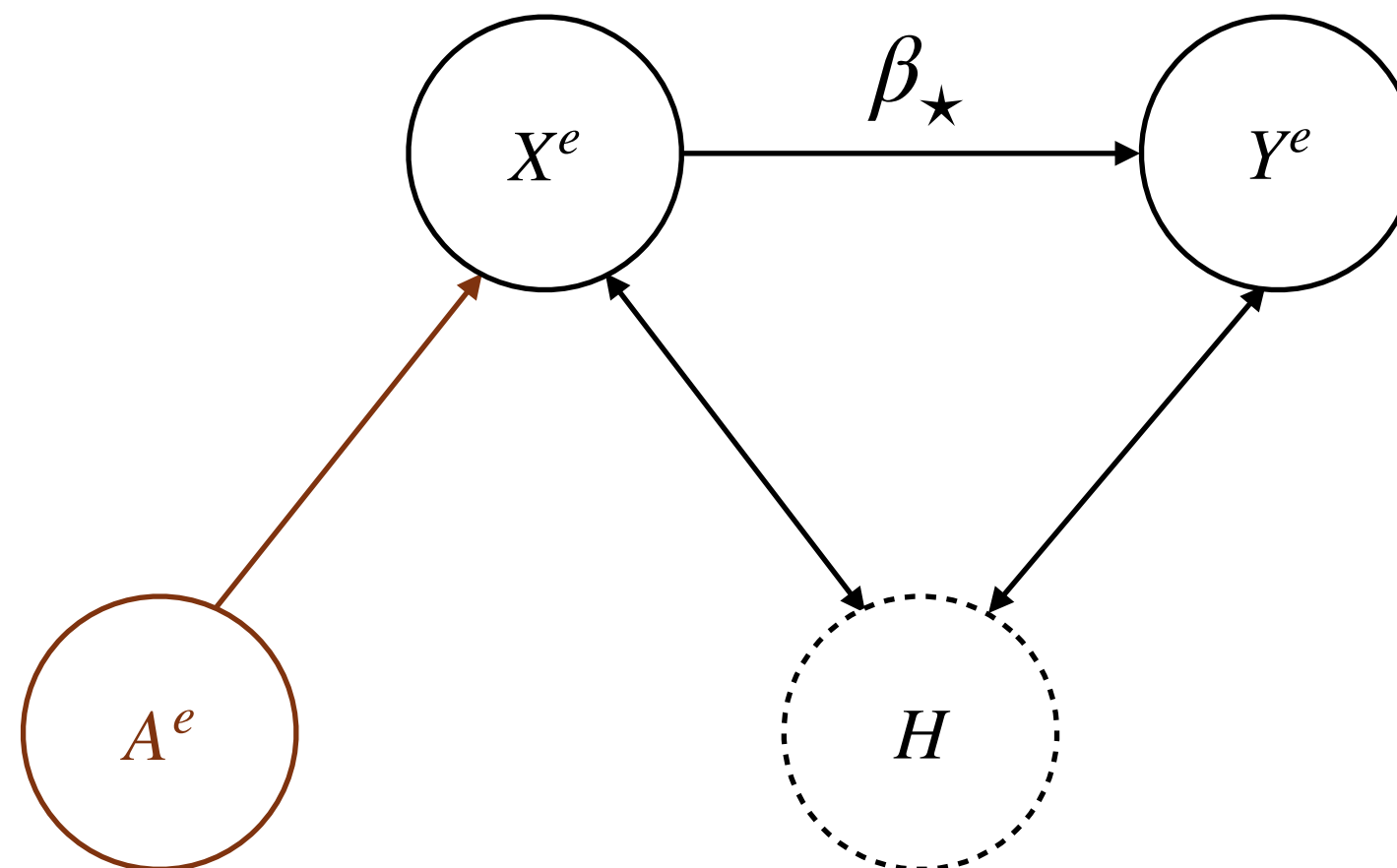
Setting of structural causal models

Data model: linear structural causal model (SCM) with unobserved confounding, environments differ via **additive shifts** A^e :

$$X^e = A^e + \eta;$$

$$Y^e = \beta_{\star}^{\top} X^e + \xi,$$

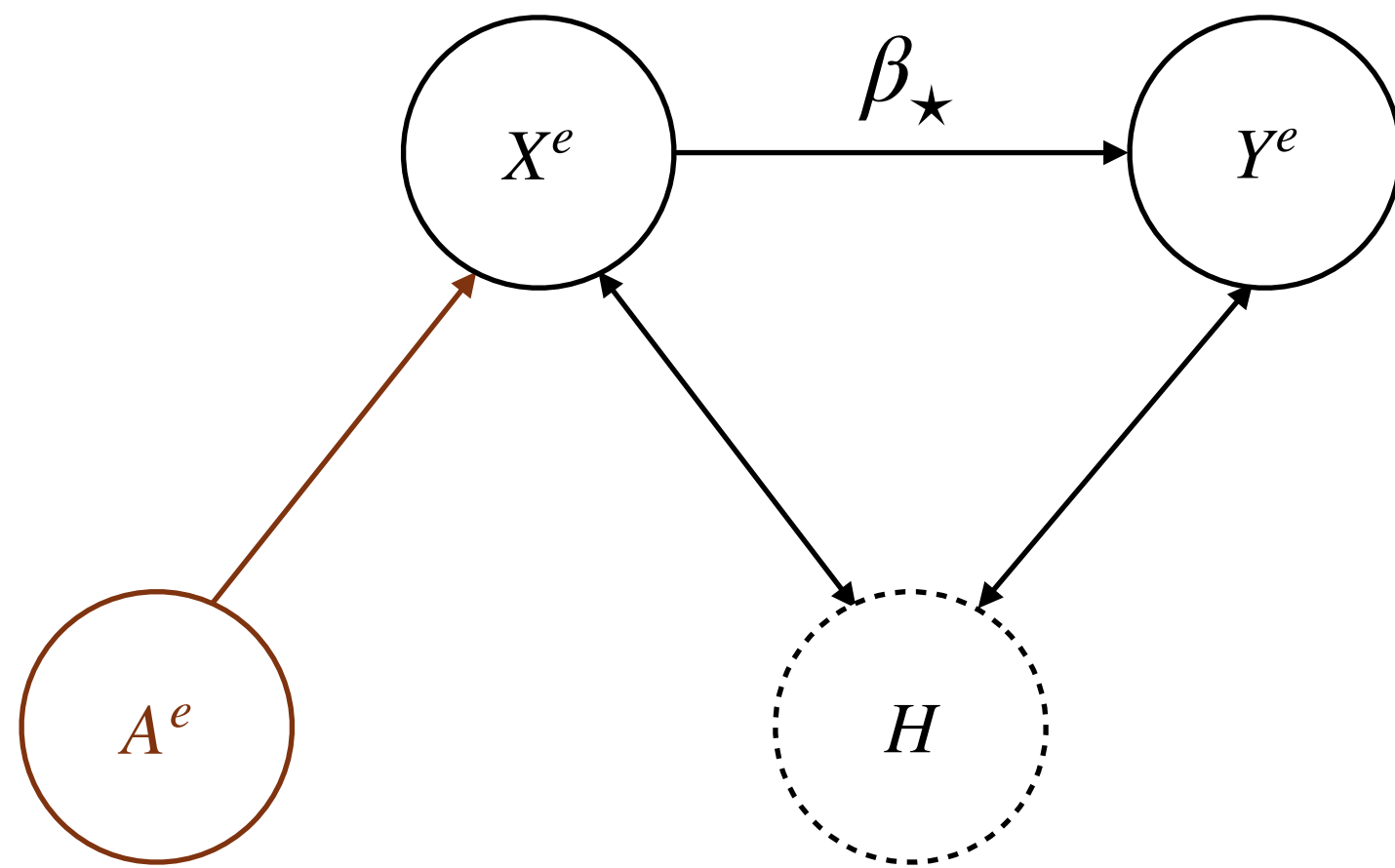
where $(\eta, \xi) \sim \mathcal{N}(0, \Sigma_{\star})$ and $\theta_{\star} = (\Sigma_{\star}, \beta_{\star})$ are the model parameters.



Setting of structural causal models

Some structural knowledge about the strength and direction of the test shift:

$$\mathbb{E}[A^{\text{test}} A^{\text{test}\top}] \preceq M_{\text{test}} = \gamma \Pi_{\mathcal{M}}.$$



- Infinite **robustness** to arbitrary shifts **only possible if β_\star known** (requires $\mathcal{O}(d)$ env's)
- However, β_\star only identified on

$$\mathcal{S} = \text{range} \left(\sum_{e \in \mathcal{E}_{\text{train}}} \mathbb{E}[A^e A^{e\top}] \right)$$

Identifiable robustness for the SCM setting

We **compute** the identifiable robust risk explicitly:

$$\mathcal{R}_{rob,ID}(\beta; \Theta_{eq}, \gamma \Pi_{\mathcal{M}}) = \underbrace{\mathcal{R}(\beta; \theta_{\star})}_{\text{Reference risk}} + \underbrace{\gamma \|S^{\top}(\beta^{\mathcal{S}} - \beta)\|_2^2}_{\text{Invariance term}} + \underbrace{\gamma (C_{ker} + \|R^{\top} \beta\|_2)^2}_{\text{Non-identifiability term}},$$

where:

- S : test shift directions along which the causal model can be identified
- R : test shift directions along which the model is non-identifiable
- C_{ker} : max. norm of the model along non-identified directions

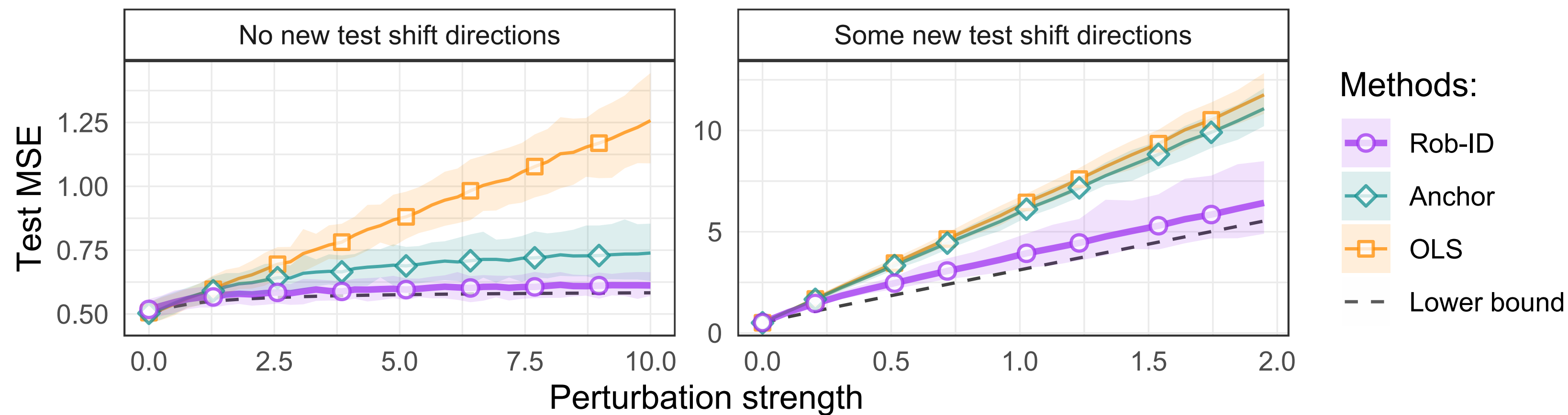
Identifiable robustness for the SCM setting

We **compute** the identifiable robust risk explicitly:

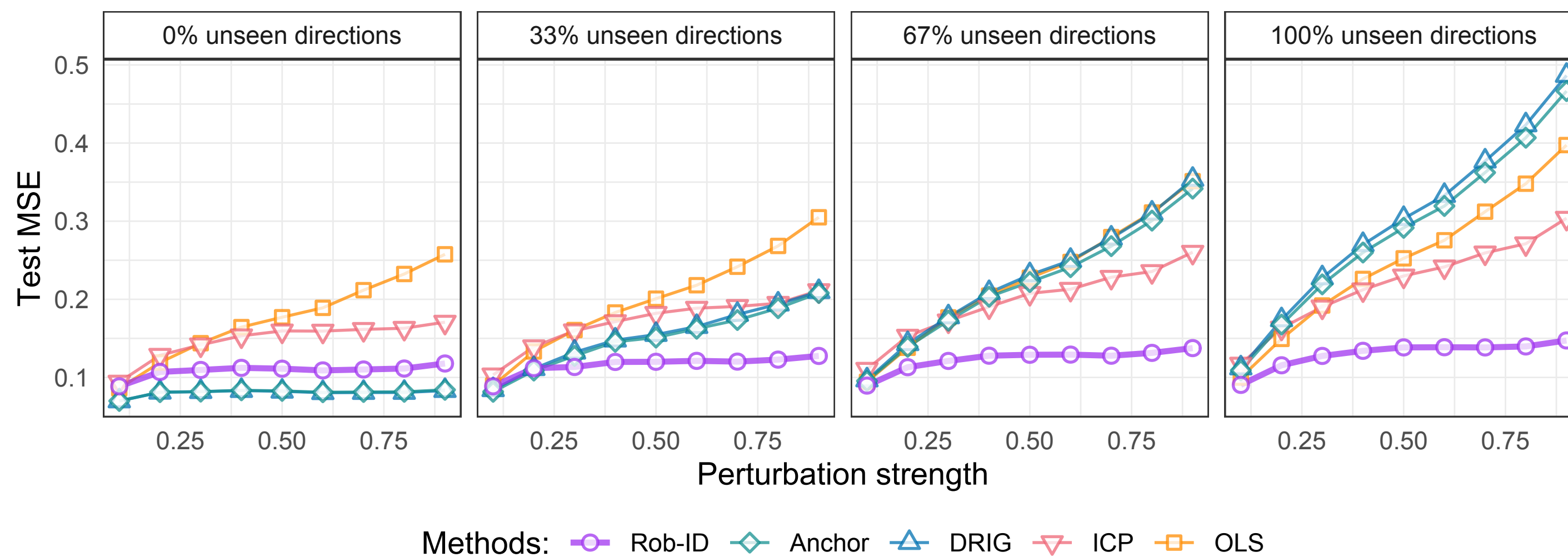
$$\mathcal{R}_{rob,ID}(\beta; \Theta_{eq}, \gamma \Pi_{\mathcal{M}}) = \underbrace{\mathcal{R}(\beta; \theta_{\star})}_{\text{Reference risk}} + \underbrace{\gamma \|S^{\top}(\beta^{\mathcal{S}} - \beta)\|_2^2}_{\text{Invariance term}} + \underbrace{\gamma(C_{ker} + \|R^{\top}\beta\|_2)^2}_{\text{Non-identifiability term}},$$

- We prove a **lower bound** for the id. robust risk which is tight for large γ ;
- For large γ , we prove **suboptimality of existing robustness methods** such as **anchor regression** [Rothenhäusler et al. 2021] and **DRIG** [Shen et al. 2023].

Simulations on **Gaussian SCM data**:



Experiments on **real-world gene expression dataset** [Replogle et al. 2022]:



Outlook

- Extension to classification
- Nonlinear models
- Use for active intervention selection
- Partially identifiable framework beyond causality



Thank you!

