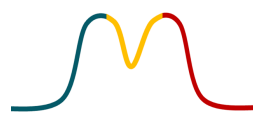# Robust prediction beyond the identifiable case

Oberwolfach Workshop January 2025
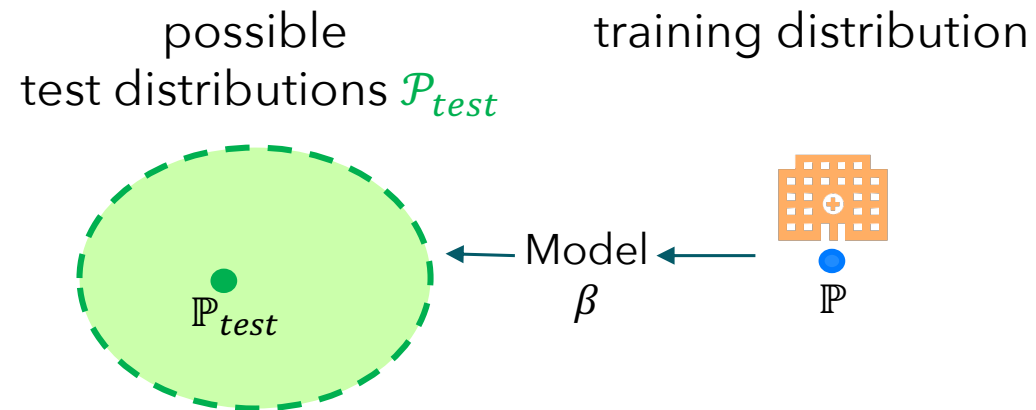on Overparameterization, **Regularization, Identifiability** and Uncertainty

Fanny Yang, joint with Julia Kostin, Nicola Gnecco
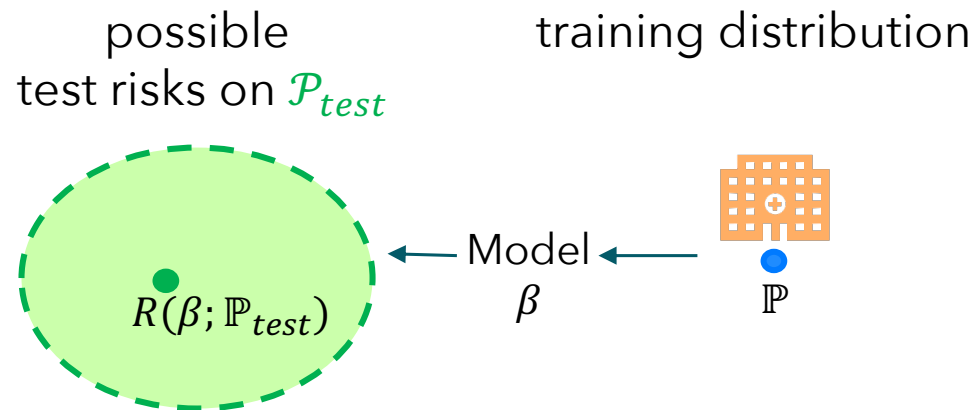
Statistical Machine Learning group,
Department of Computer Science, ETH Zurich

# Robust prediction for safety purposes

possible
test distributions $\mathcal{P}_{test}$

training distribution



$\mathbb{P}_{test}$

Model
$\beta$

$\mathbb{P}$

# Robust prediction for safety purposes

possible
test risks on $\mathcal{P}_{test}$

training distribution

$R(\beta; \mathbb{P}_{test})$

Model
$\beta$

$\mathbb{P}$

# Robust prediction for safety purposes

worst-case
risk in $\mathcal{P}_{test}$ for $\beta$

possible
test risks on $\mathcal{P}_{test}$

training distribution



$$R_{rob}(\beta) = \sup_{\mathbb{P} \in \mathcal{P}_{test}} R(\beta; \mathbb{P})$$
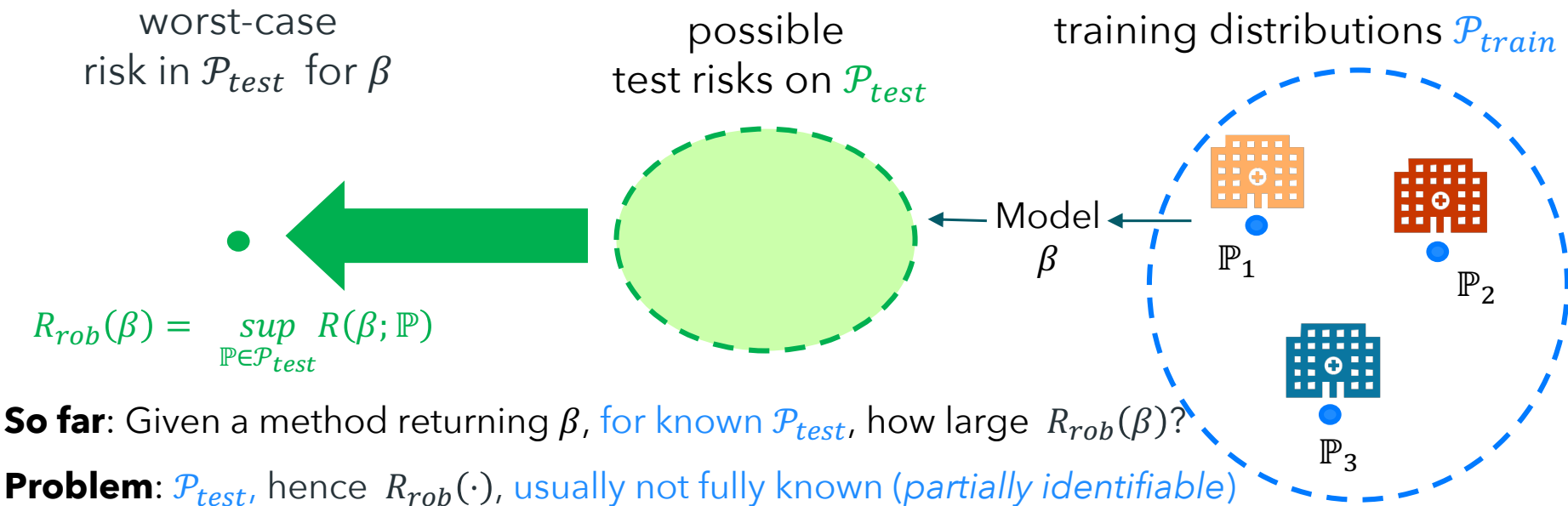
- A model $\beta$ is *more robust* if it has smaller $R_{rob}(\beta)$

# Robust prediction for safety purposes

worst-case
risk in $\mathcal{P}_{test}$ for $\beta$

possible
test risks on $\mathcal{P}_{test}$

training distributions $\mathcal{P}_{train}$



$$R_{rob}(\beta) = \sup_{\mathbb{P}\in\mathcal{P}_{test}} R(\beta; \mathbb{P})$$

- A model $\beta$ is *more robust* if it has smaller $R_{rob}(\beta)$

- Any robustness gains from observing multiple heterogeneous training distributions?

# Robustness analysis of methods – what's missing?



worst-case risk in $\mathcal{P}_{test}$ for $\beta$

possible test risks on $\mathcal{P}_{test}$

training distributions $\mathcal{P}_{train}$

Model $\beta$

$\mathbb{P}_1$

$\mathbb{P}_2$

$\mathbb{P}_3$

$$R_{rob}(\beta) = \sup_{\mathbb{P} \in \mathcal{P}_{test}} R(\beta; \mathbb{P})$$

**So far**: Given a method returning $\beta$, for known $\mathcal{P}_{test}$, how large $R_{rob}(\beta)$?

**Problem**: $\mathcal{P}_{test}$, hence $R_{rob}(\cdot)$, usually not fully known (*partially identifiable*)

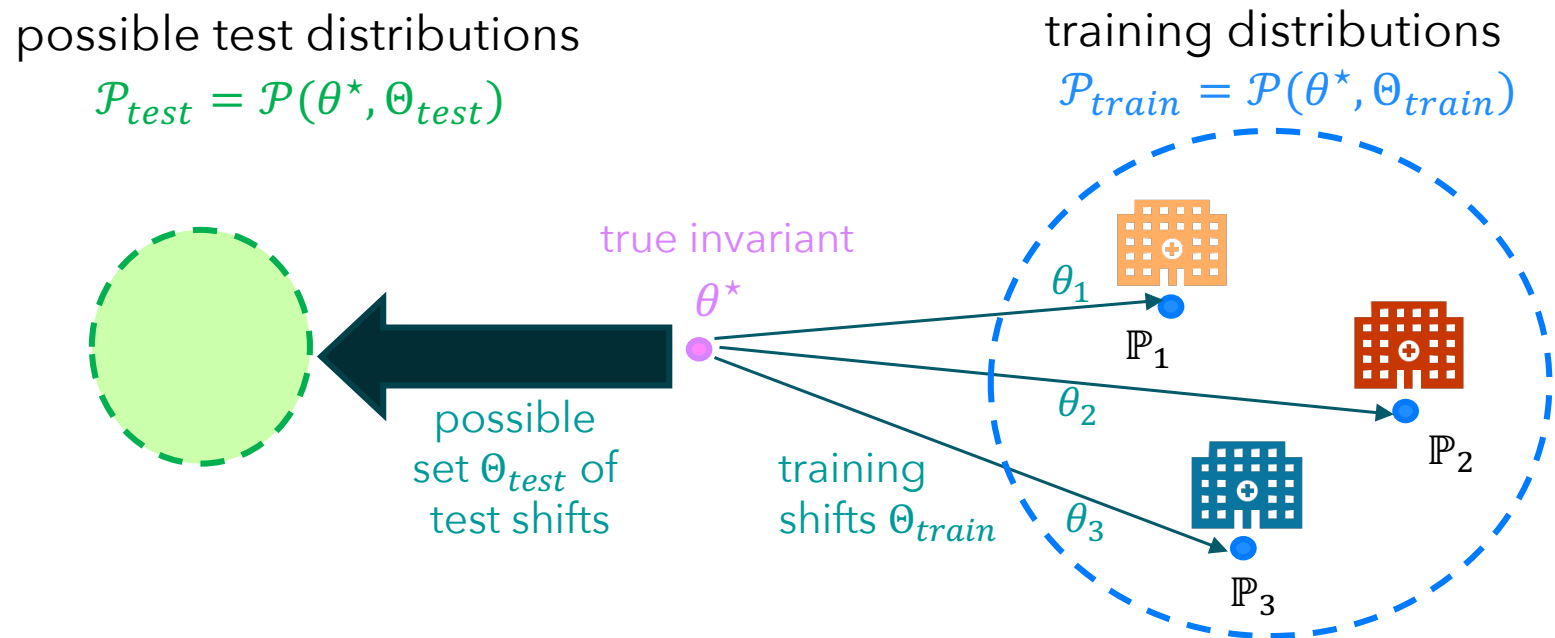**Neglected question**: Given partial knowledge about $\mathcal{P}_{test}$,

- how robust can *any algorithm* be, i.e. what is the "information-theoretic" (population) limit?

- how well do *existing algorithms perform*, and how close to optimal/adaptive are they?

our work

# How can we model partial knowledge of $\mathcal{P}_{test}$ / $R_{rob}$ via its relationship to $\mathcal{P}_{train}$?

- Setting up unified shift robustness view via invariance (+ one example)

- From fully identifiable (prior work) to partially identifiable (our work) $R_{rob}$

- Measure of robustness and hardness in partially identifiable case

# Unified view of shift robustness using invariance

possible test distributions
$$\mathcal{P}_{test} = \mathcal{P}(\theta^\star, \Theta_{test})$$

training distributions
$$\mathcal{P}_{train} = \mathcal{P}(\theta^\star, \Theta_{train})$$

true invariant
$\theta^\star$

possible
set $\Theta_{test}$ of
test shifts

training
shifts $\Theta_{train}$

$\theta_1$

$\mathbb{P}_1$

$\theta_2$

$\mathbb{P}_2$

$\theta_3$

$\mathbb{P}_3$

- Assume that $(\theta^\star, \theta_e)$ parameterize distributions $\mathbb{P}_e$ with $\theta^\star$ invariant and $\theta_e$ varying with $e$

- Viewpoint includes traditional shift concepts (covariate shift, spurious correlations, domain mixtures, neighborhood) & causality-based ones (IRM-related or next slide)

# Imagine simple linear example for concreteness...

Assume that joint distributions in each "environment" $e$ in train and test environments are defined by
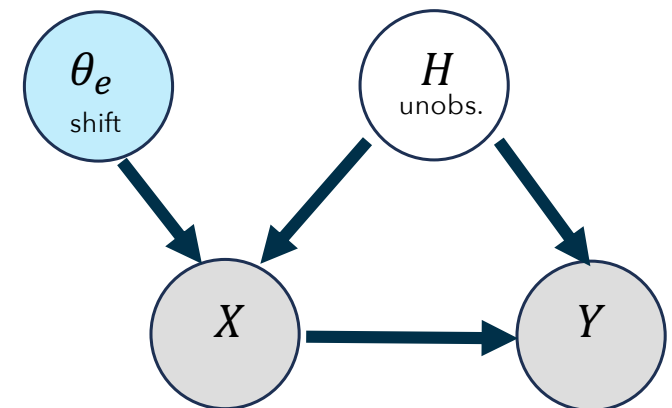
mean shifts varying with $e$
(assume ref. env has $\theta_e = 0$)

$$X^e = \theta_e + \eta$$
$$Y^e = \beta_\star^\top X^e + \xi$$

exogeneous noise
$(\eta, \xi) \sim N(0, \Sigma_\star)$
invariant covariance

with invariant $\theta_\star = (\beta_\star, \Sigma_\star)$ same across environments

Possible underlying causal model (most simplified version)



We allow cross-covariance $\Sigma_{\star,\eta\xi} \neq 0$ corresponding to confounding

$\Longrightarrow$ allows not only covariate shift, but also shift in $\mathbb{E}[Y|X]$!
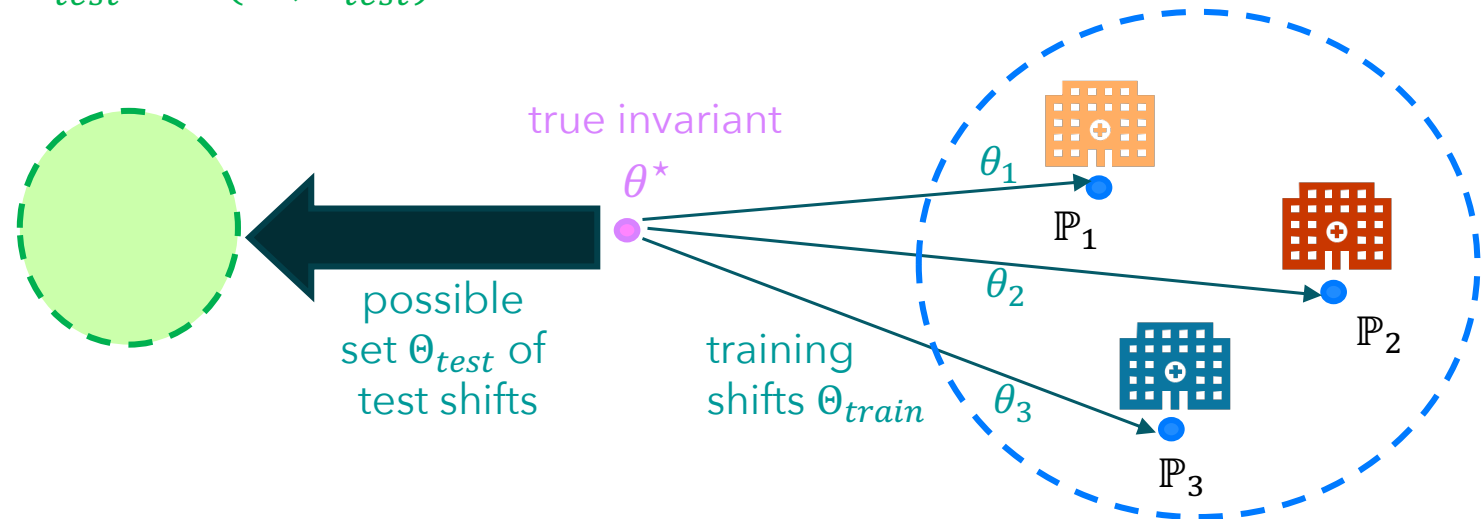
# Unified view of shift robustness using invariance

possible test distributions
$$\mathcal{P}_{test} = \mathcal{P}(\theta^{\star}, \Theta_{test})$$

training distributions
$$\mathcal{P}_{train} = \mathcal{P}(\theta^{\star}, \Theta_{train})$$



true invariant
$\theta^{\star}$

possible set $\Theta_{test}$ of test shifts

training shifts $\Theta_{train}$

$\theta_1$  $\mathbb{P}_1$

$\theta_2$  $\mathbb{P}_2$

$\theta_3$  $\mathbb{P}_3$

**Remember: we're interested in answering,** given some invariance assumption & any $\Theta_{\text{test}}, \Theta_{train}$
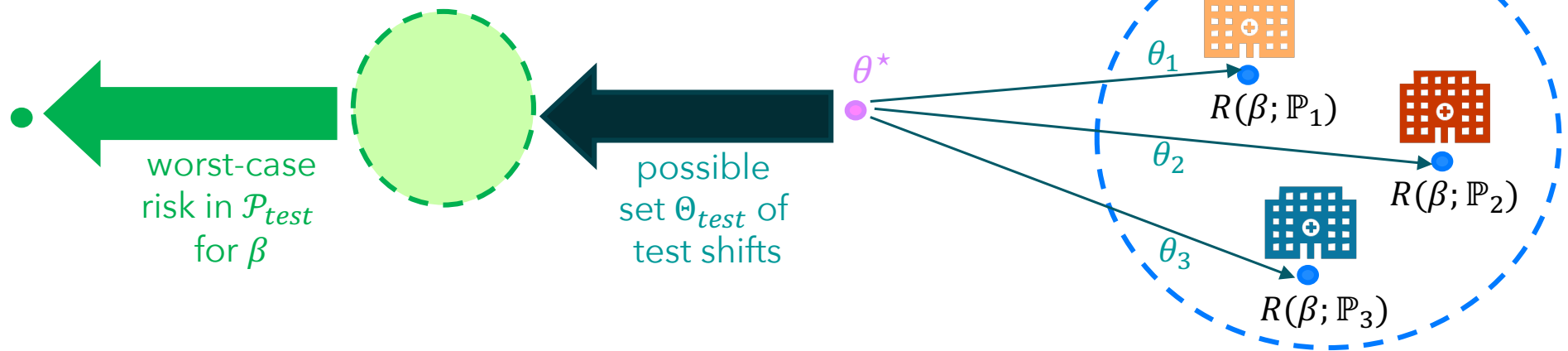
- how robust can *any algorithm* be, i.e. what is the "information-theoretic" (population) limit?

- how do *existing algorithms perform*, and how close to optimal/adaptive are they?

# Measuring robustness via robust risk



robust risk
$R_{rob}(\beta; \theta^\star, \Theta_{test})$

possible test risks $R(\beta; \mathbb{P})$
for $\mathbb{P} \in \mathcal{P}(\theta^\star, \Theta_{test})$

training risks
$\{R(\beta, \mathbb{P}): \mathbb{P} \in \mathcal{P}(\theta^\star, \Theta_{train})\}$

worst-case
risk in $\mathcal{P}_{test}$
for $\beta$

possible
set $\Theta_{test}$ of
test shifts

$\theta^\star$

$\theta_1$

$\theta_2$

$\theta_3$

$R(\beta; \mathbb{P}_1)$

$R(\beta; \mathbb{P}_2)$
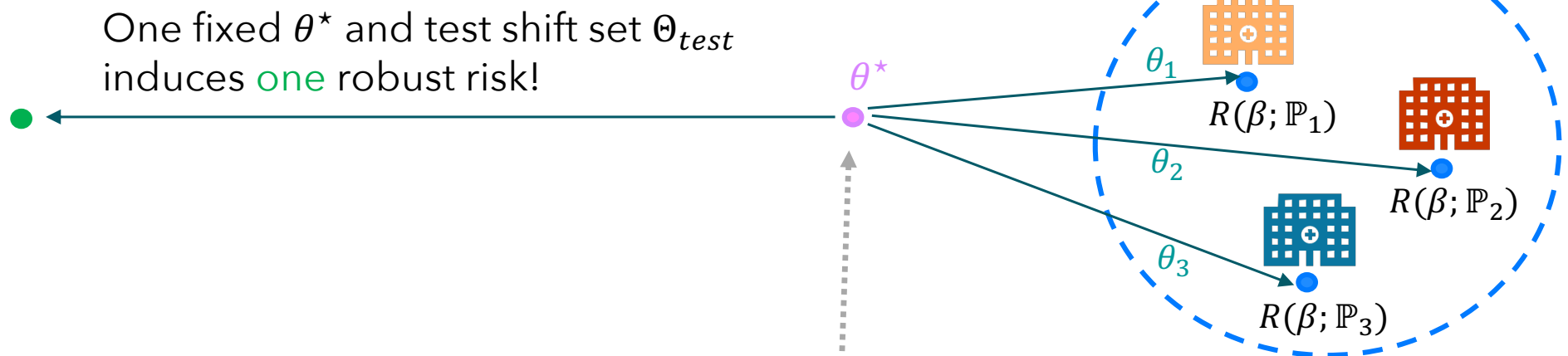
$R(\beta; \mathbb{P}_3)$

# Measuring robustness via robust risk

robust risk
$R_{rob}(\beta; \theta^\star, \Theta_{test})$

invariant
parameter

training risks
$\{R(\beta, \mathbb{P}): \mathbb{P} \in \mathcal{P}(\theta^\star, \Theta_{train})\}$



One fixed $\theta^\star$ and test shift set $\Theta_{test}$
induces one robust risk!

$\theta^\star$

$\theta_1$

$R(\beta; \mathbb{P}_1)$

$R(\beta; \mathbb{P}_2)$

$\theta_2$

$\theta_3$

$R(\beta; \mathbb{P}_3)$

*Invariant parameter is unknown/unobserved!*

# Prior work: Assuming identifiable robust risk

robust risk
$R_{rob}(\beta; \theta^\star, \Theta_{test})$

training distributions
$\mathcal{P}_{train} = \mathcal{P}(\theta^\star, \Theta_{train})$

Robust risk identifiable, i.e. computable

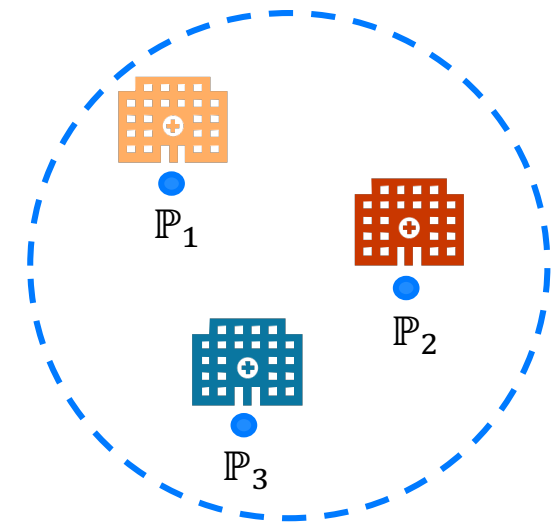using observed $\mathcal{P}_{train}, \Theta_{train}$ and $\Theta_{test}$?

Previous work on robustness only considers identifiable case

- for invariance-based shift models this only holds for

  specific combinations of $\Theta_{test}, \Theta_{train}$

- any other combination naturally corresponds to some kind

  of partial knowledge of $\mathcal{P}_{test}$

$\mathbb{P}_1$

$\mathbb{P}_2$

$\mathbb{P}_3$

(always true for DRO/mixture of domains with no $\theta^\star$ such as in e.g. Mansour et al. '08, Sagawa et al. '19)　　13
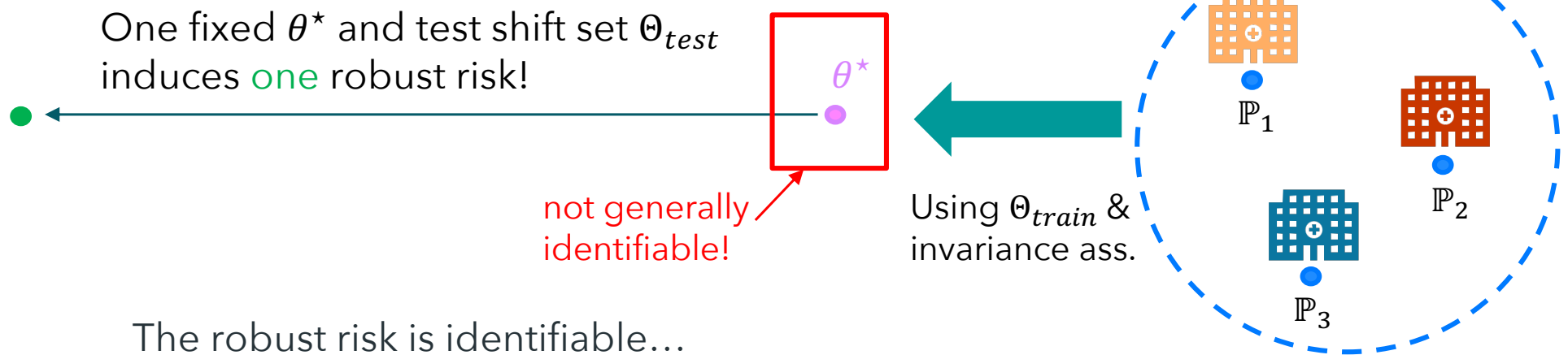
# Prior work I: Identifiable invariant mechanism $\theta^\star$

robust risk
$R_{rob}(\beta; \theta^\star, \Theta_{test})$
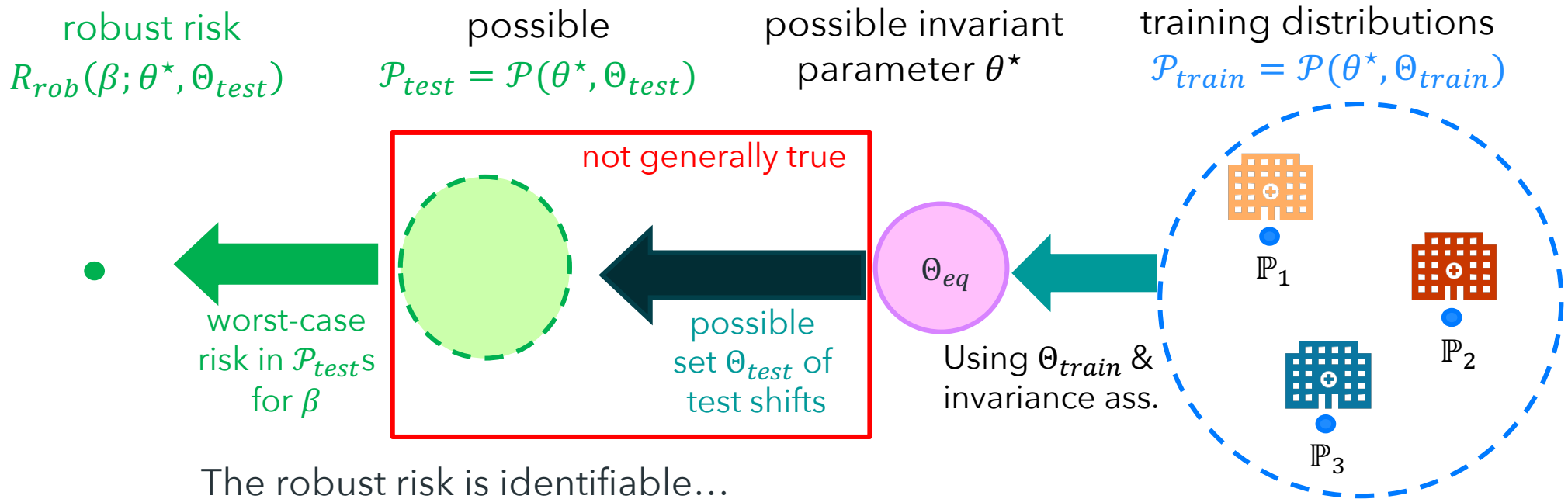
possible invariant parameter $\theta^\star$

training distributions
$\mathcal{P}_{train} = \mathcal{P}(\theta^\star, \Theta_{train})$

One fixed $\theta^\star$ and test shift set $\Theta_{test}$ induces one robust risk!

$\theta^\star$

not generally identifiable!

Using $\Theta_{train}$ & invariance ass.

$\mathbb{P}_1$

$\mathbb{P}_2$

$\mathbb{P}_3$

The robust risk is identifiable…

- when $\Theta_{train}$ is heterogeneous enough to identify $\theta^\star$

e.g. Peters et al. '16, Rojas-Carulla et al. '18, Arjovsky et al. '19, Krueger '20

# Prior work II: Only identifiable robust risk

robust risk
$R_{rob}(\beta; \theta^\star, \Theta_{test})$

possible
$\mathcal{P}_{test} = \mathcal{P}(\theta^\star, \Theta_{test})$

possible invariant
parameter $\theta^\star$

training distributions
$\mathcal{P}_{train} = \mathcal{P}(\theta^\star, \Theta_{train})$

not generally true

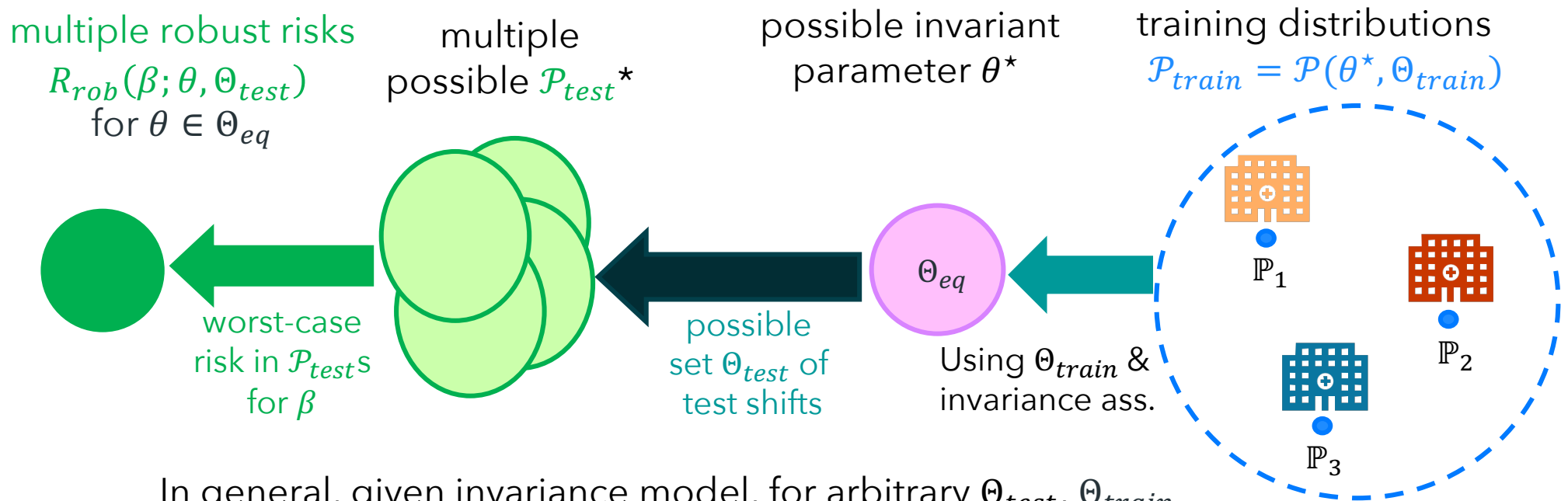$\Theta_{eq}$

$\mathbb{P}_1$

$\mathbb{P}_2$

$\mathbb{P}_3$

worst-case
risk in $\mathcal{P}_{test}$s
for $\beta$

possible
set $\Theta_{test}$ of
test shifts

Using $\Theta_{train}$ &
invariance ass.

The robust risk is identifiable...

- when $\Theta_{train}$ is heterogeneous enough to identify $\theta^\star$

- when $\Theta_{test}$ similar to $\Theta_{train}$ if one can't identify $\theta^\star$

e.g. Rothenhaeusler et al. '21, Shen et al. '23

# Our work: general partially identifiable robust risk

multiple robust risks
$R_{rob}(\beta; \theta, \Theta_{test})$
for $\theta \in \Theta_{eq}$

multiple
possible $\mathcal{P}_{test}*$

possible invariant
parameter $\theta^{\star}$

training distributions
$\mathcal{P}_{train} = \mathcal{P}(\theta^{\star}, \Theta_{train})$



$\Theta_{eq}$

$\mathbb{P}_1$

$\mathbb{P}_2$

$\mathbb{P}_3$

worst-case
risk in $\mathcal{P}_{test}$s
for $\beta$

possible
set $\Theta_{test}$ of
test shifts

Using $\Theta_{train}$ &
invariance ass.

In general, given invariance model, for arbitrary $\Theta_{test}$, $\Theta_{train}$

we end up only with partially/set-identifying the robust risk!

# Our work: general partially identifiable robust risk
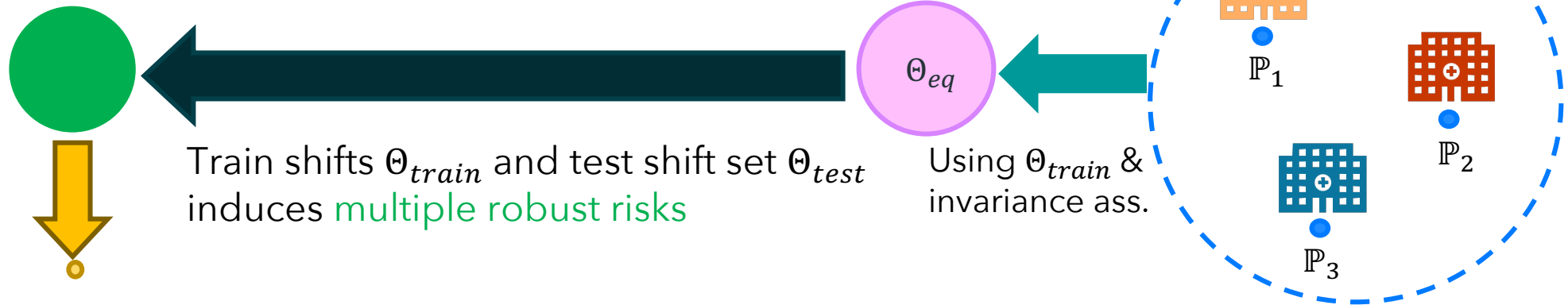
multiple robust risks
$R_{rob}(\beta; \theta, \Theta_{test})$
for $\theta \in \Theta_{eq}$

possible invariant
parameter $\theta^\star$

training distributions
$\mathcal{P}_{train} = \mathcal{P}(\theta^\star, \Theta_{train})$



Train shifts $\Theta_{train}$ and test shift set $\Theta_{test}$
induces multiple robust risks

Using $\Theta_{train}$ &
invariance ass.

In general, given invariance model, for arbitrary $\Theta_{test}$, $\Theta_{train}$

we end up only with partially/set-identifying the robust risk!

How do we even measure robustness in this case?

# Quantifying robustness in partial identifiable setting

multiple robust risks
$R_{rob}(\beta; \theta, \Theta_{test})$
for $\theta \in \Theta_{eq}$

possible invariant
parameter $\theta^\star$

training distributions
$\mathcal{P}_{train} = \mathcal{P}(\theta^\star, \Theta_{train})$



$\Theta_{eq}$

$\mathbb{P}_1$

$\mathbb{P}_2$

$\mathbb{P}_3$

Train shifts $\Theta_{train}$ and test shift set $\Theta_{test}$
induces multiple robust risks

Using $\Theta_{train}$ &
invariance ass.

worst-case robust risk
$\mathfrak{R}_{rob}(\beta; \Theta_{train}, \Theta_{test}) =$
$\max_{\theta \in \Theta_{eq}} R_{rob}(\beta; \theta, \Theta_{test})$

…want small robust risk even
for the hardest true $\theta \in \Theta_{eq}$ that
could have induced $\mathcal{P}_{train}$

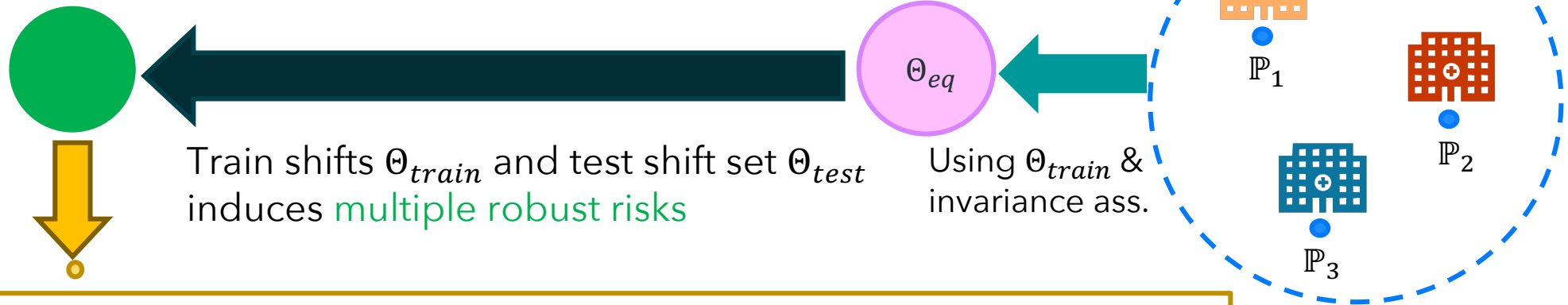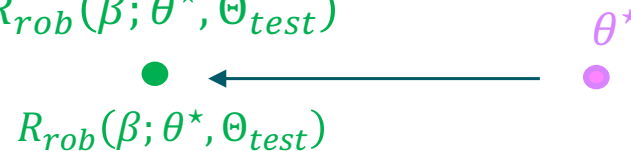# Achievable robustness in partial identifiable setting

multiple robust risks
$R_{rob}(\beta; \theta, \Theta_{test})$
for $\theta \in \Theta_{eq}$

possible invariant
parameter $\theta^{\star}$

training distributions
$\mathcal{P}_{train} = \mathcal{P}(\theta^{\star}, \Theta_{train})$

$\Theta_{eq}$

$\mathbb{P}_1$

$\mathbb{P}_2$

$\mathbb{P}_3$

Train shifts $\Theta_{train}$ and test shift set $\Theta_{test}$
induces multiple robust risks

Using $\Theta_{train}$ &
invariance ass.

worst-case robust risk
$\mathfrak{R}_{rob}(\beta; \Theta_{train}, \Theta_{test}) =$
$\max\limits_{\theta \in \Theta_{eq}} R_{rob}(\beta; \theta, \Theta_{test})$

and achievable worst-case robust risk
$\mathfrak{M}(\Theta_{train}, \Theta_{test}) = \min\limits_{\beta} \mathfrak{R}_{rob}(\beta; \Theta_{train}, \Theta_{test})$
$= \min\limits_{\beta} \max\limits_{\theta \in \Theta_{eq}} R_{rob}(\beta; \theta, \Theta_{test})$

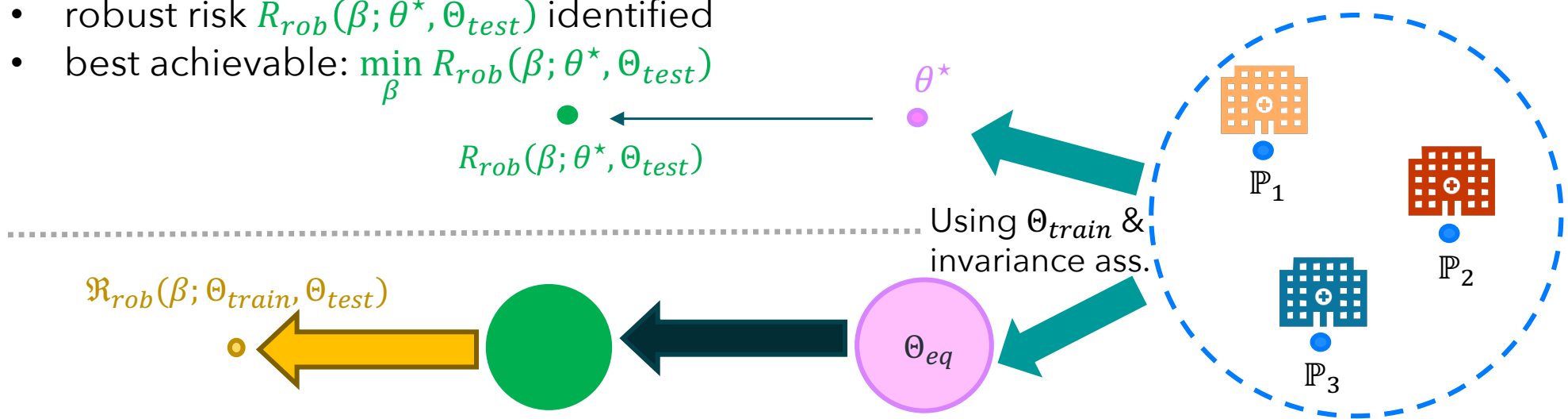allow us to quantify
robustness in the
non-identifiable case

# Summary of differences in identifiability

**Identifiable case** (prior work):

- robust risk $R_{rob}(\beta; \theta^\star, \Theta_{test})$ identified
- best achievable: $\min_\beta R_{rob}(\beta; \theta^\star, \Theta_{test})$

training distributions
$\mathcal{P}_{train} = \mathcal{P}(\theta^\star, \Theta_{train})$

$\theta^\star$

$R_{rob}(\beta; \theta^\star, \Theta_{test})$

$\mathbb{P}_1$

$\mathbb{P}_2$

$\mathbb{P}_3$

Using $\Theta_{train}$ & invariance ass.

$\mathfrak{R}_{rob}(\beta; \Theta_{train}, \Theta_{test})$

$\Theta_{eq}$

**Partially identifiable case** (ours):

- only worst-case robust risk $\mathfrak{R}_{rob}(\beta; \Theta_{train}, \Theta_{test}) = \max_{\theta \in \Theta_{eq}} R_{rob}(\beta; \theta, \Theta_{test})$ identified

- best-achievable: $\mathfrak{M}(\Theta_{train}, \Theta_{test}) = \min_\beta \mathfrak{R}_{rob}(\beta; \Theta_{train}, \Theta_{test})$

**Remember: we're interested in answering,** given some invariance assumption & any $\Theta_{\text{test}}, \Theta_{train}$

- how robust can *any algorithm* be, i.e. what is the "information-theoretic" (population) limit?

- how do *existing algorithms perform*, and how close to optimal/adaptive are they?

We quantify for invariance-based methods in this general setting

- the best achievable robustness $\mathfrak{M}(\Theta_{train}, \Theta_{test})$

- how ranking of different methods wrt $\mathfrak{R}_{rob}(\beta; \Theta_{train}, \Theta_{test})$ changes drastically with varying $\Theta_{test}, \Theta_{train}$

theoretically for linear model
empirically for real data

# Simple linear example for concreteness

Assume that joint distributions in each "environment" $e$

in train and test environments are defined by

mean shifts varying with $e$
(assume ref. env has $\theta_e = 0$)

exogeneous noise
$(\eta, \xi) \sim N(0, \Sigma_\star)$
invariant covariance

$$X^e = \theta_e + \eta$$
$$Y^e = \beta_\star^\top X^e + \xi$$

with invariant $\theta_\star = (\beta_\star, \Sigma_\star)$ same across environments

**Test time shifts assumptions**

$M_{seen}$: covariance with range

in span of seen shift directions

$range(M_{seen}) \subset span \{\theta_e\}_{e \in [k]}$

$M_{unseen}$: projection matrix onto

unseen directions with

$range(M_{seen}) \perp span \{\theta_e\}_{e \in [k]}$

Mean shifts during test time assumed to lie in $\Theta_{test} = \{\theta_{test} : \theta_{test}\theta_{test}^\top \preccurlyeq \gamma M_{seen} + \gamma' M_{unseen}\}$

shift strengths

# Achievable and achieved robustness

$$X^e = A^e + \eta$$
$$Y^e = \beta_\star^\mathsf{T} X^e + \xi$$
$$\Theta_{test} = \{\theta_{test}: \theta_{test}\theta_{test}^\mathsf{T} \preccurlyeq \gamma M_{seen} + \gamma' M_{unseen}\}$$

**Trends concluded from formal statement**

In partially identifiable case & new test shift directions $\gamma'$ large, anchor regression and OLS

- are far from achievable robustness $\mathfrak{M}(\Theta_{train}, \Theta_{test}) = \min_\beta \mathfrak{R}_{rob}(\beta; \Theta_{eq}, \Theta_{test})$

- have similar linear robustness when term with unseen directions $\gamma'$ dominates

**Corollary [KGY' 24] (informal) – Performance comparison in the partially identifiable setting**

For large $\gamma'$ fixed $\gamma$, $\mathfrak{M}(\Theta_{train}, \Theta_{test}) = \min_\beta \mathfrak{R}_{rob}(\beta; \Theta_{eq}, \Theta_{test}) = C^2 \gamma' + c_1$

vs. Anchor regression*: $\mathfrak{R}_{rob}(\beta_{anchor}; \Theta_{train}, \Theta_{test}) = (C + h(\gamma))^2 \gamma' + c_2$   ($h$ is decreasing

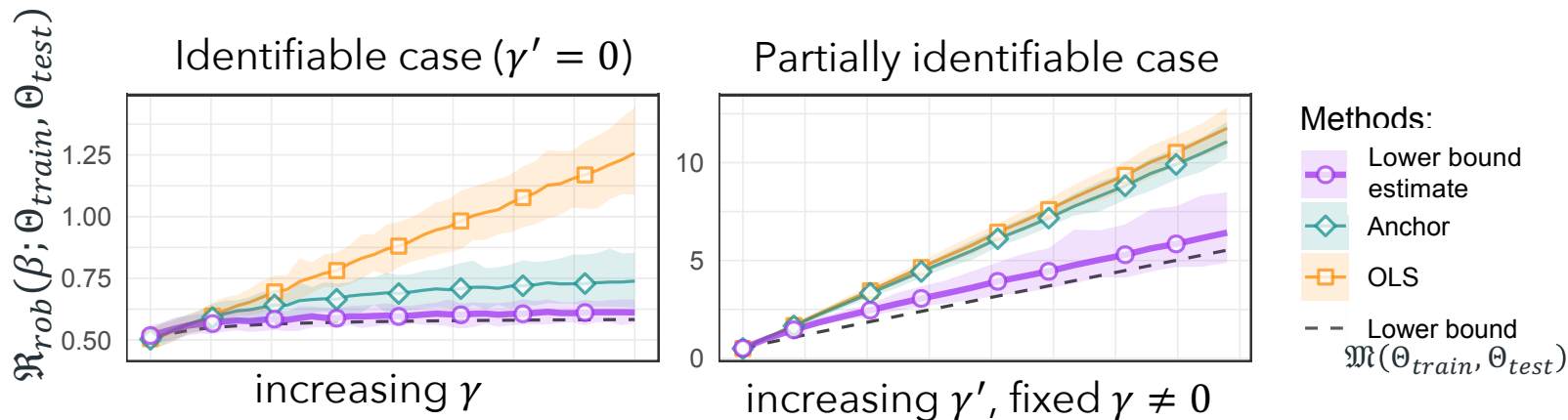vs. Ordinary least squares: $\mathfrak{R}_{rob}(\beta_{OLS}; \Theta_{train}, \Theta_{test}) = (C + h(1))^2 \gamma' + c_3$     function, $c$ are constant in $\gamma'$)

* Rothenhaeusler et al. '21

$X^e = A^e + \eta$
$Y^e = \beta_\star^\mathsf{T} X^e + \xi$
$\Theta_{test} = \{\theta_{test}: \theta_{test}\theta_{test}^\mathsf{T}$
$\preccurlyeq \gamma M_{seen} + \gamma' M_{unseen}\}$

# Experimental comparison in linear setting

**Trends concluded from formal statement**

In partially identifiable case & new test shift directions $\gamma'$ large, anchor regression and OLS

- are far from achievable robustness $\mathfrak{M}(\Theta_{train}, \Theta_{test}) = \min_\beta \mathfrak{R}_{rob}(\beta; \Theta_{eq}, \Theta_{test})$

- have similar linear robustness when term with unseen directions $\gamma'$ dominates



Identifiable case ($\gamma' = 0$)

Partially identifiable case

increasing $\gamma$

increasing $\gamma'$, fixed $\gamma \neq 0$

Methods:
- Lower bound estimate
- Anchor
- OLS
- Lower bound $\mathfrak{M}(\Theta_{train}, \Theta_{test})$

- Using correct $\Theta_{test}$

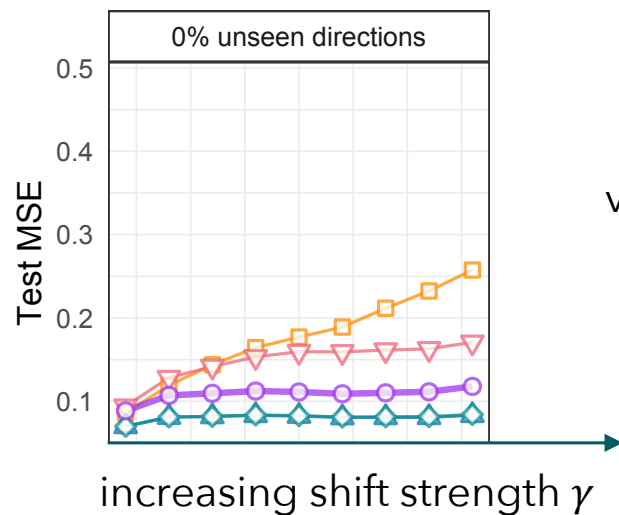- Average over random draws of $\theta^\star$

# Real-world data setting

Single–cell gene expression dataset with

- prediction task: expression of one gene (target) as a function of expressions of 3 others

- per target, 3 different environments ($\triangleq$ individual gene knocked out) + observational
  - training environments: observational + 1 knocked-out environment
  - shift strengths $\gamma, \gamma' \triangleq$ distance of covariates to mean in observational environment
  - partially identified setting: test data also includes some percentage of samples from knock-out environments not seen during training

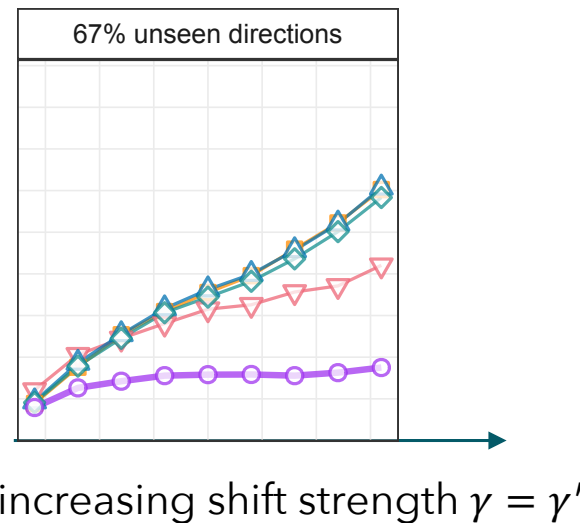- Results are averaged across these scenarios

# Experimental comparison in real-world setting

identifiable case ($\gamma' = 0$)   partially identifiable case ($\gamma' \neq 0$)



vs.

Test set here consists of:

- 67% from knock-outs unseen during training

- 33% held-out data from seen knock-out env

Different invariance-based Methods: — Rob-ID — Anchor — DRIG — ICP — OLS

$\underset{\beta}{\mathrm{argmin}}\ \widehat{\Re}_{rob}(\beta; \Theta_{train}, \Theta_{test})$ assuming previous synthetic setting where $\mathrm{M}_{unseen}$

is most conservatively chosen to be entire $\left(span\ \{\theta_e\}_{e \in [k]}\right)^{\perp}$ (no reason to do well!)

# Summary

How analyze the more general partially identifiable

setting (vs. focusing on identifiable vs. non-identifiable)

- introduced measures of (achievable) robustness

- computed them for a linear example and

  compared achievable robustness with prior methods

J. Kostin, N. Gnecco, F. Yang "*Achievable distributional robustness when the robust
risk is only partially identified*", NeurIPS 2024

# Future work

- Apply on other types of invariant mechanisms

  (see e.g. Francesco's, Arthur's talk, and

  beyond causality)

- Use achievable robustness for

  active selection of training distributions