



Berkeley
UNIVERSITY OF CALIFORNIA
Electrical Engineering
and Computer Sciences

Computational
PRECISION HEALTH

The Data Addition Dilemma

Irene Y. Chen

Assistant Professor, Computational Precision Health
Electrical Engineering and Computer Science
Berkeley AI Research
UC Berkeley and UCSF



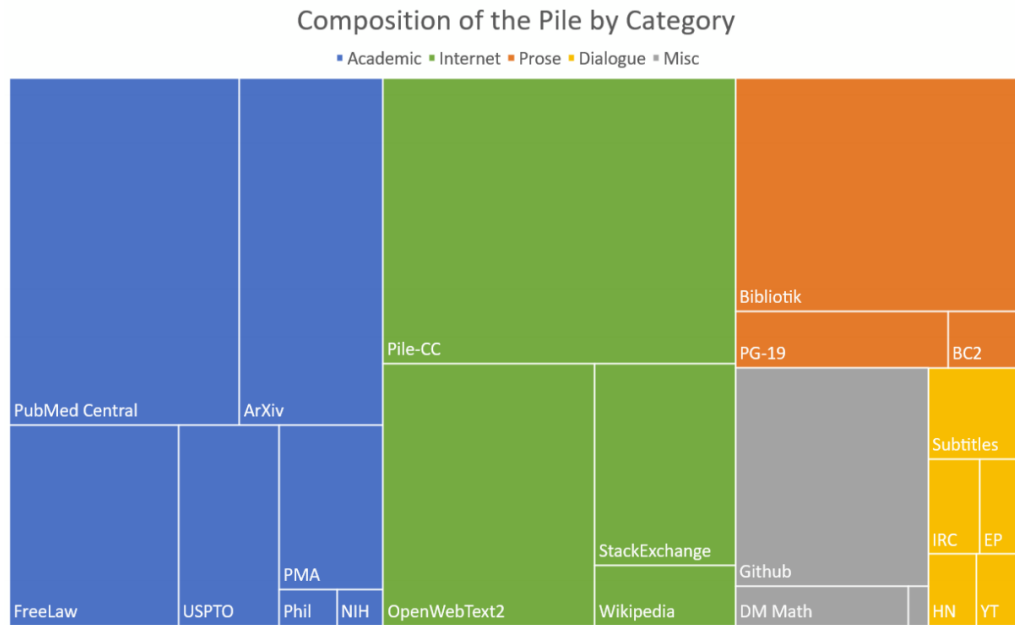
irenechen.net



@irenetrampoline



Age of bigger and bigger datasets



Pile dataset
(825 GB, 22 sources)

text string	timestamp string	url string
Beginners BBQ Class Taking Place in Missoula! Do you want to get better at making delicious BBQ? Yo...	2019-04-25T12:57:54Z	https://klyq.com/beginners-bbq-class-taking-p-in-missoula/
Discussion in 'Mac OS X Lion (10.7)' started by axboi87, Jan 20, 2012. I've got a 500gb internal...	2019-04-21T10:07:13Z	https://forums.macrumors.com/threads/restore-larger-disk-to-smaller-disk.1311329/
Foil plaid lycra and spandex shortall with metallic slinky insets. Attached metallic elastic belt with...	2019-04-25T10:40:23Z	https://awishcometrue.com/Catalogs/Clearance/V1960-Find-A-Way
How many backlinks per day for new site? Discussion in 'Black Hat SEO' started by Omoplata, Dec 3,...	2019-04-21T12:46:19Z	https://www.blackhatworld.com/seo/how-many-backlinks-per-day-for-new-site.258615/
The Denver Board of Education opened the 2017-18 school year with an update on projects that includ...	2019-04-20T14:33:21Z	http://bond.dpsk12.org/category/news/

C4 (750 GB)

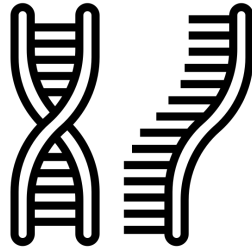


StarCoder (783 GB,
86 programming languages) ²

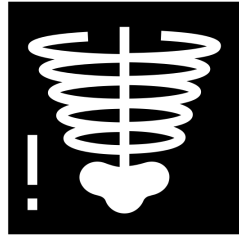
The pursuit of more health data



Electronic Medical
Records



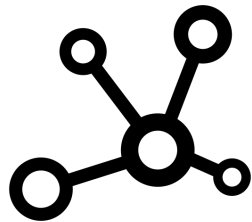
Genomics



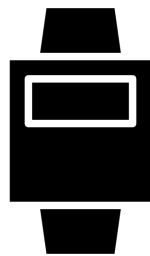
Medical Imaging



Signals



Molecular Data



Wearable Data

The pursuit of more health data

All of Us
RESEARCH PROGRAM

800k patients
50 clinical sites

^{uk} **biobank**

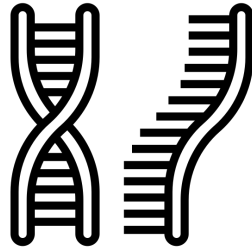
500k individuals
22 recruitment centers

**eICU Collaborative
Research Database**

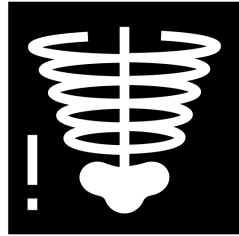
200k patient stays
208 hospitals



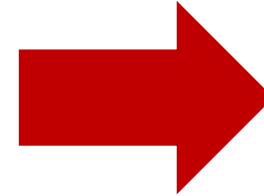
Electronic Medical
Records



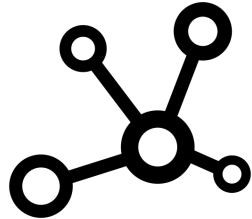
Genomics



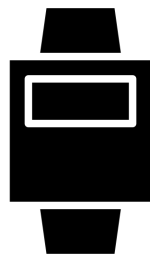
Medical Imaging



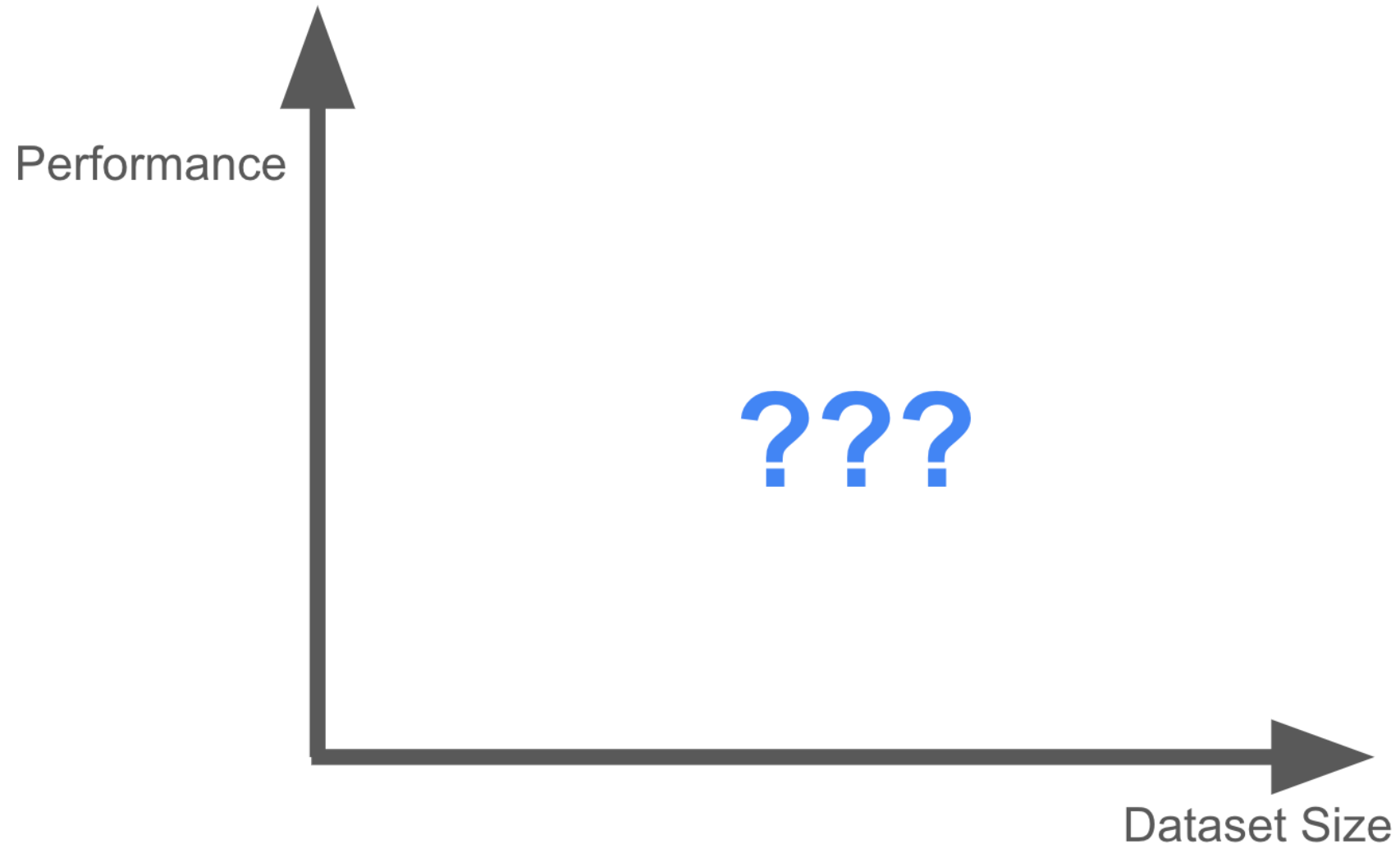
Signals

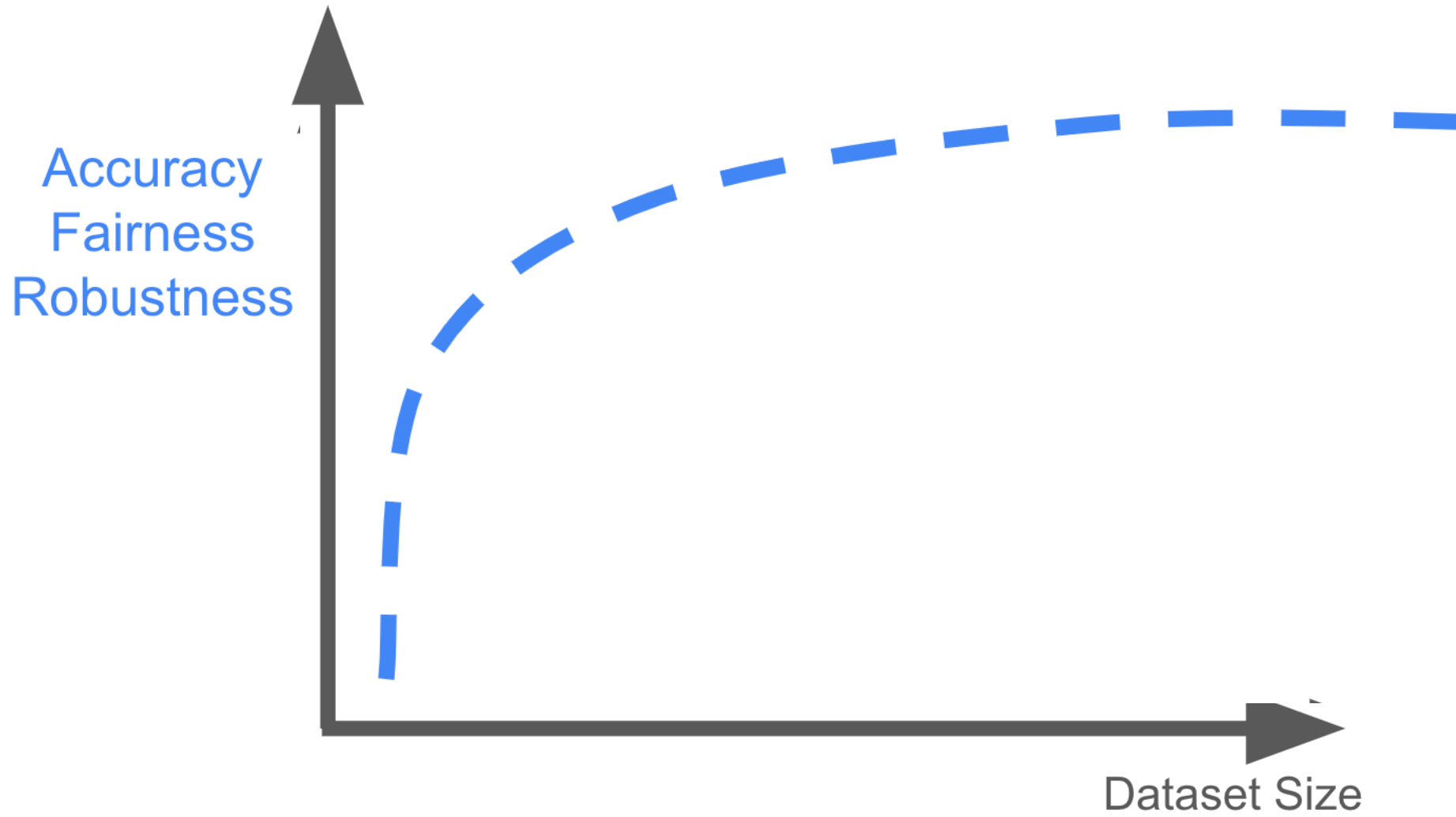


Molecular Data



Wearable Data







Judy Hanwen Shen
(Stanford)

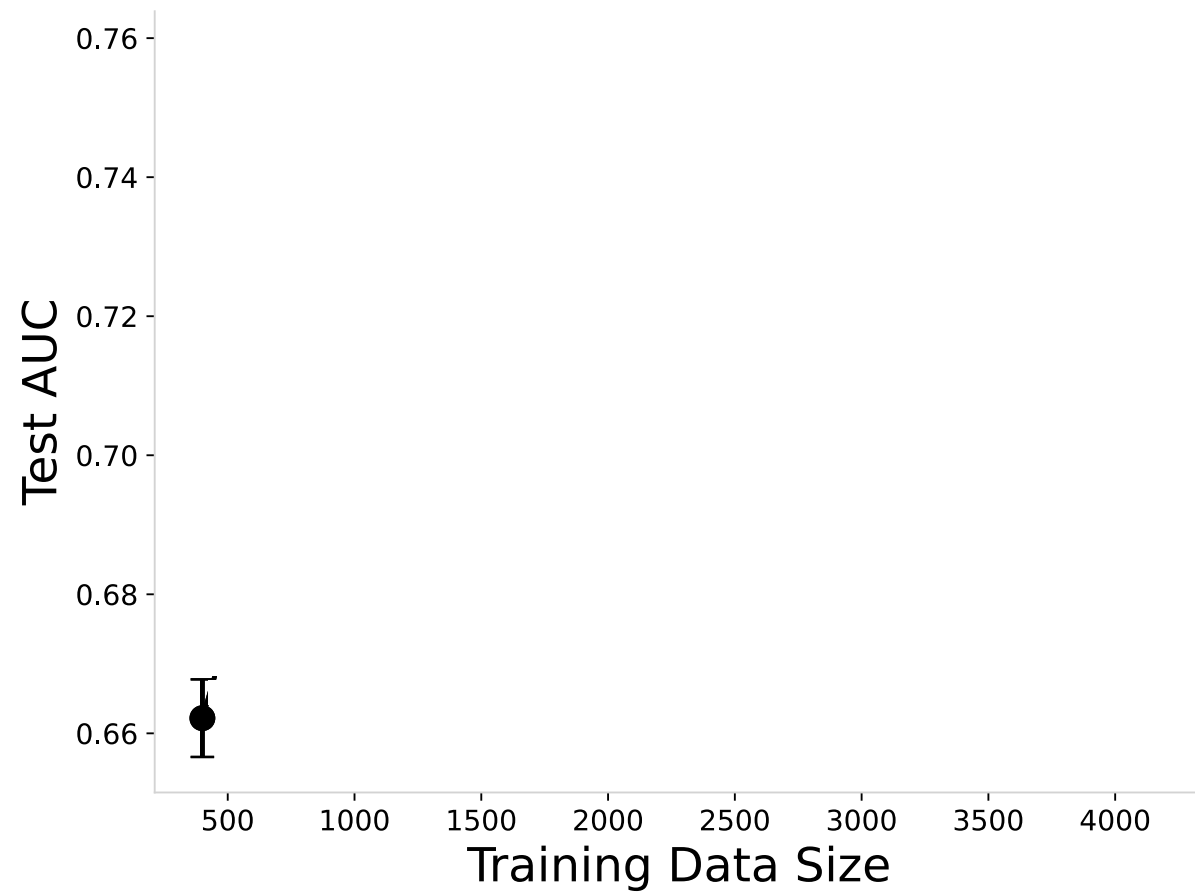


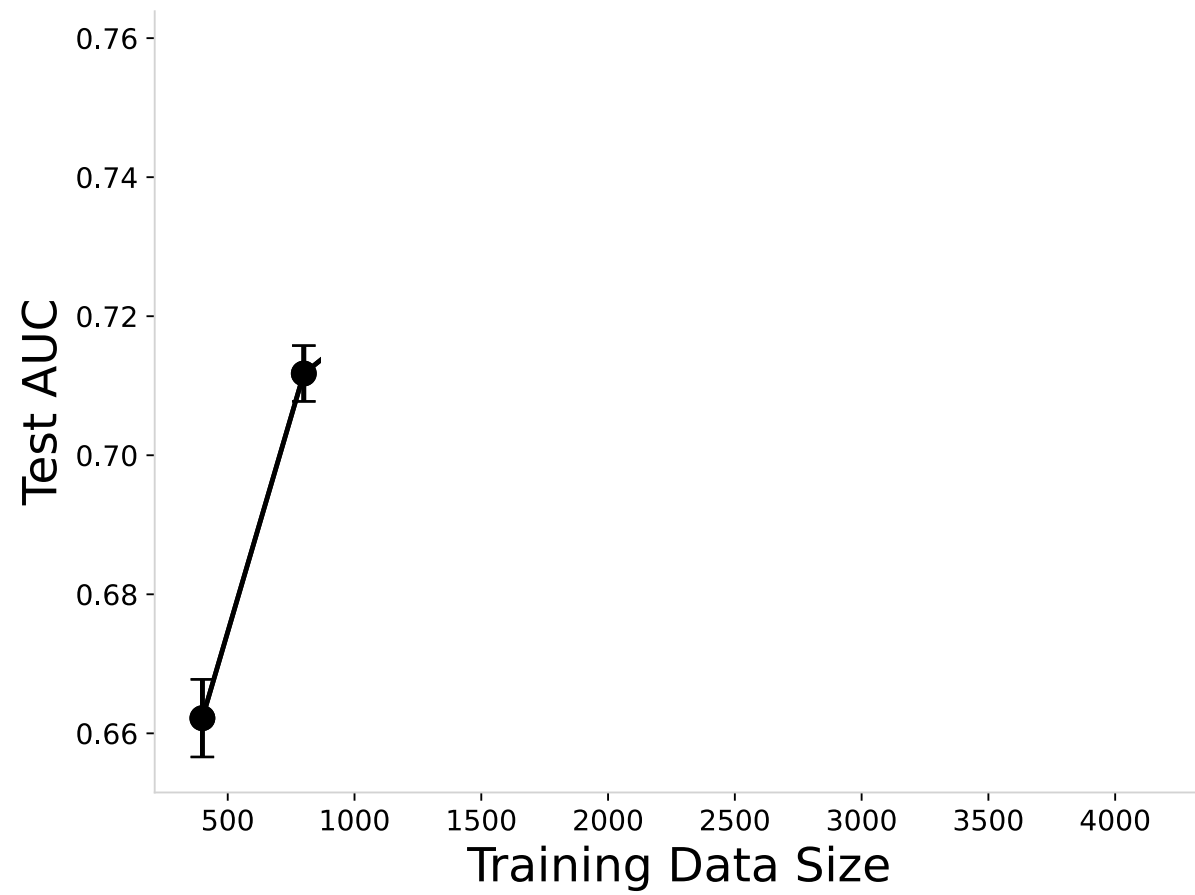
Inioluwa Deborah Raji
(UC Berkeley)

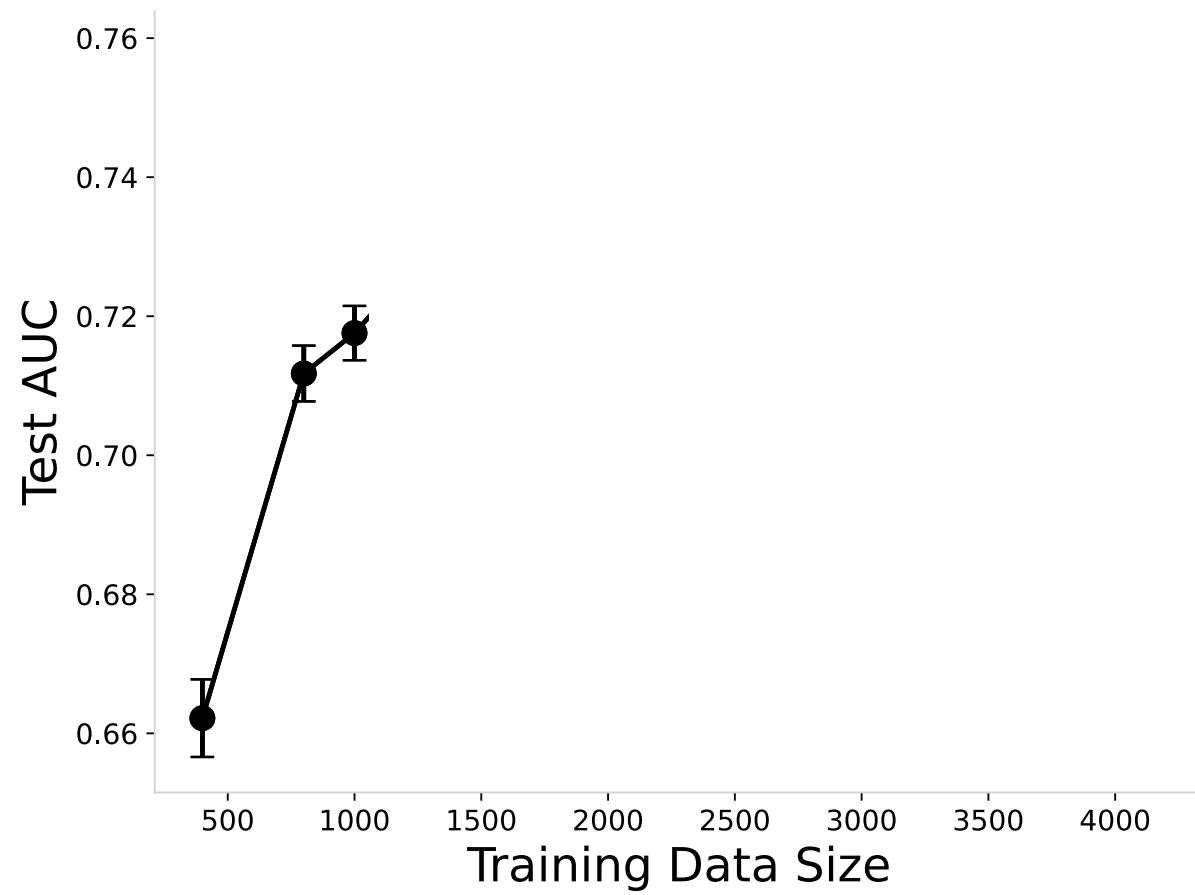
How do we balance data accumulation and distribution shift?

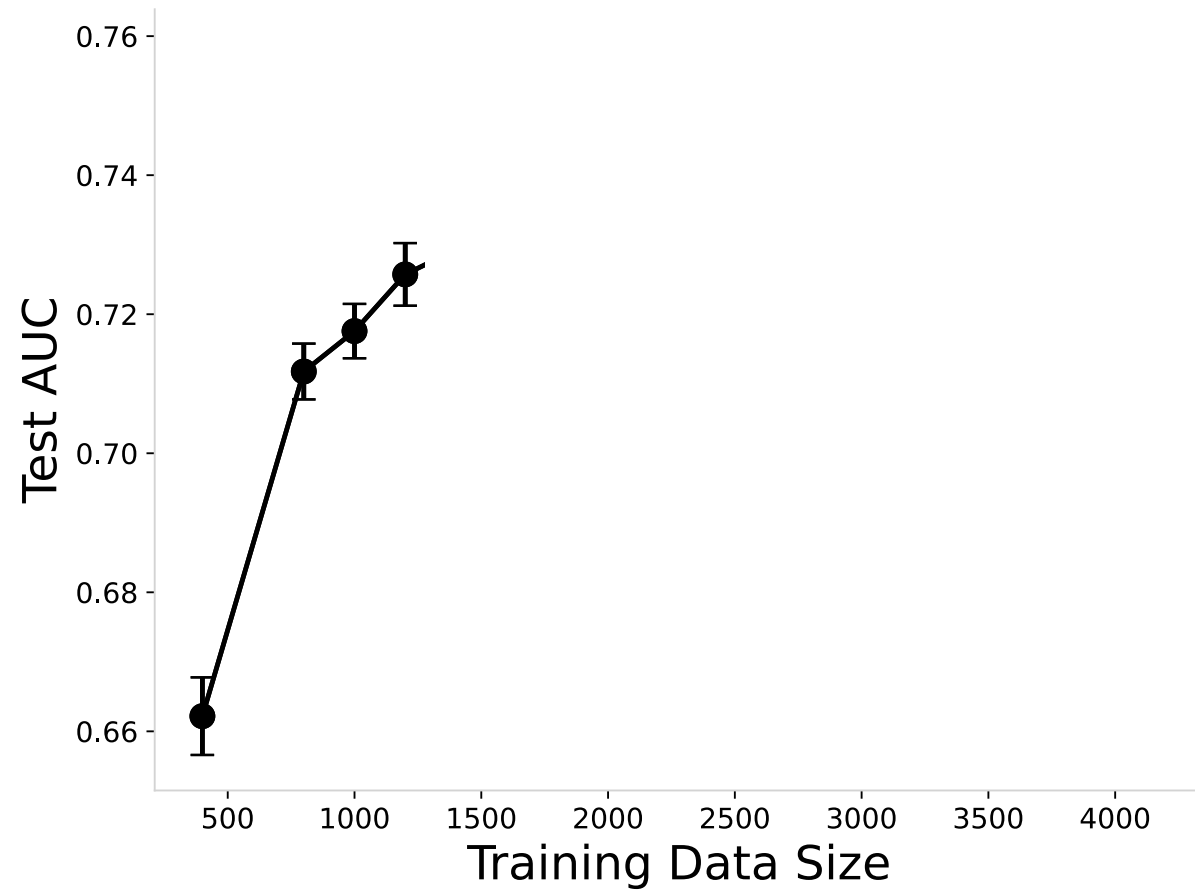
Shen et al, "The Data Addition Dilemma." MLHC 2024

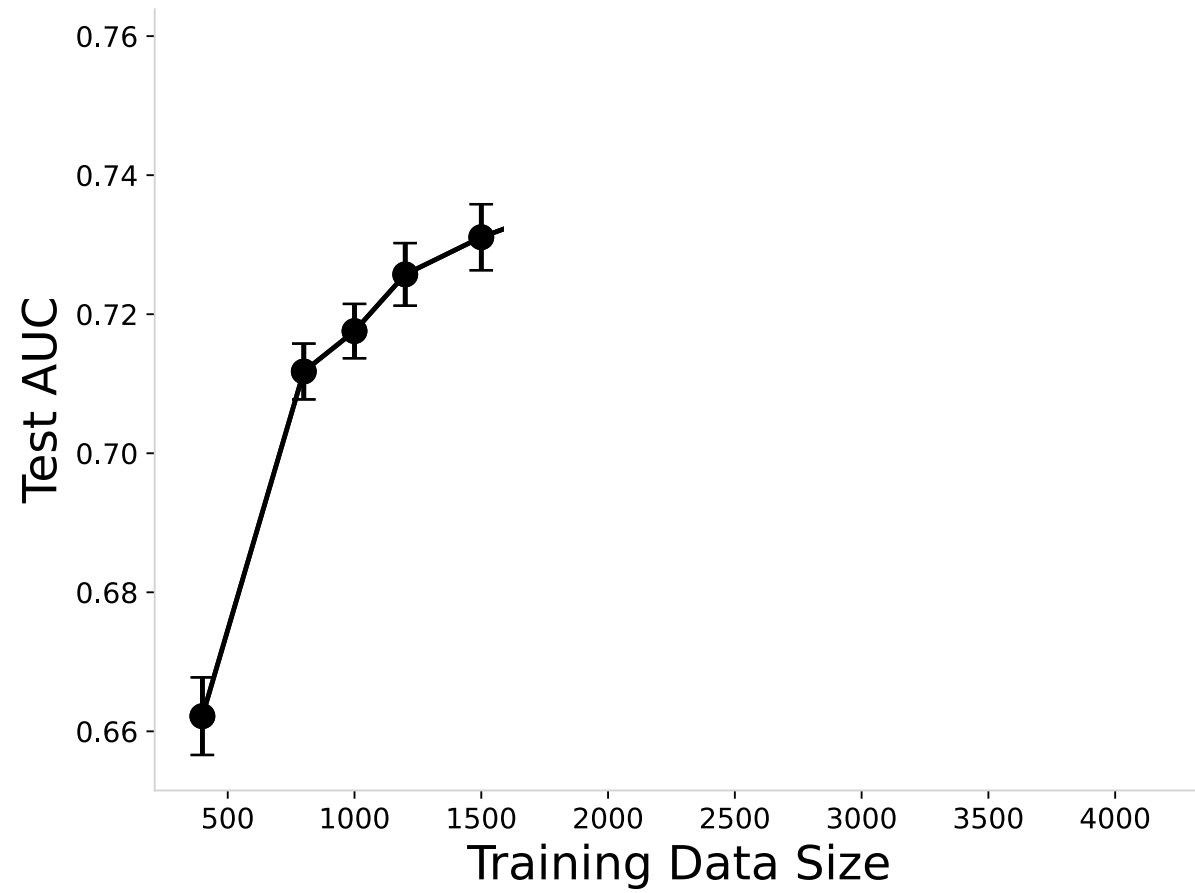


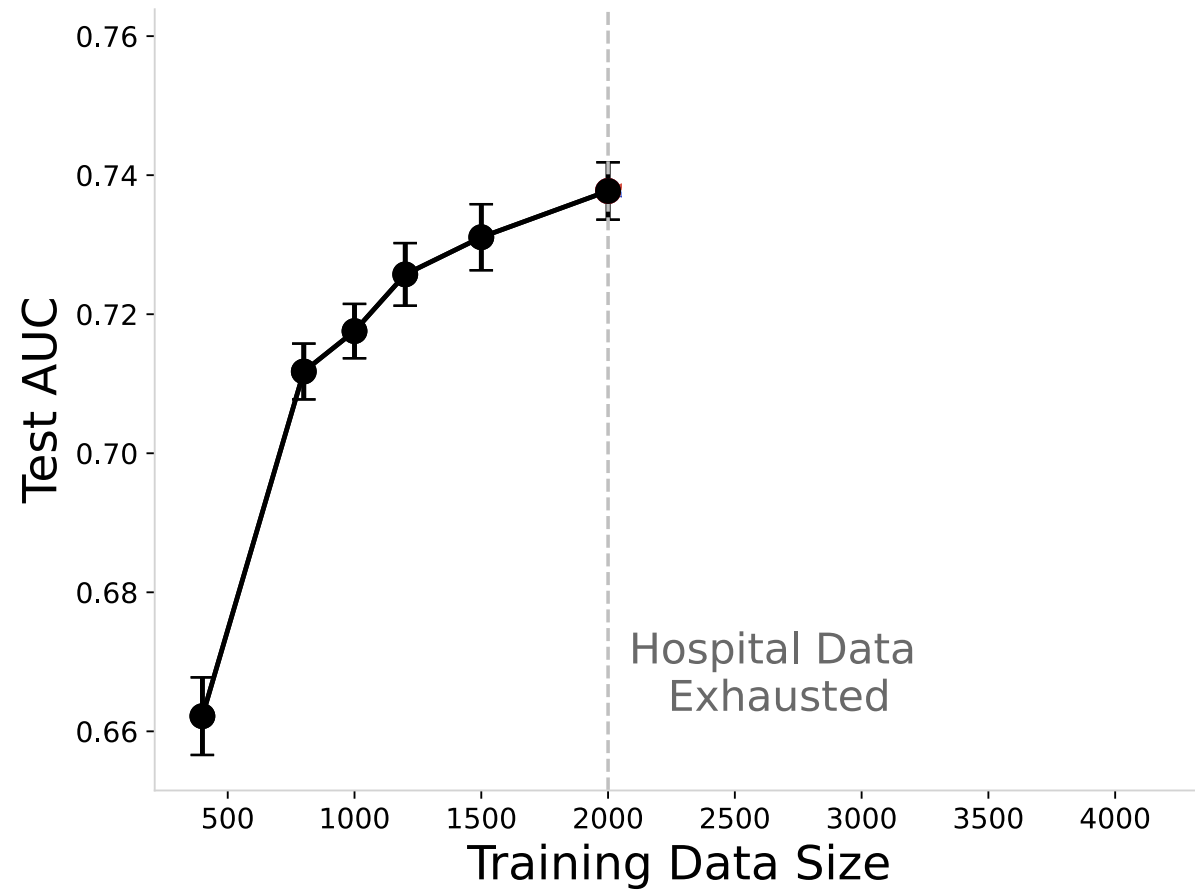


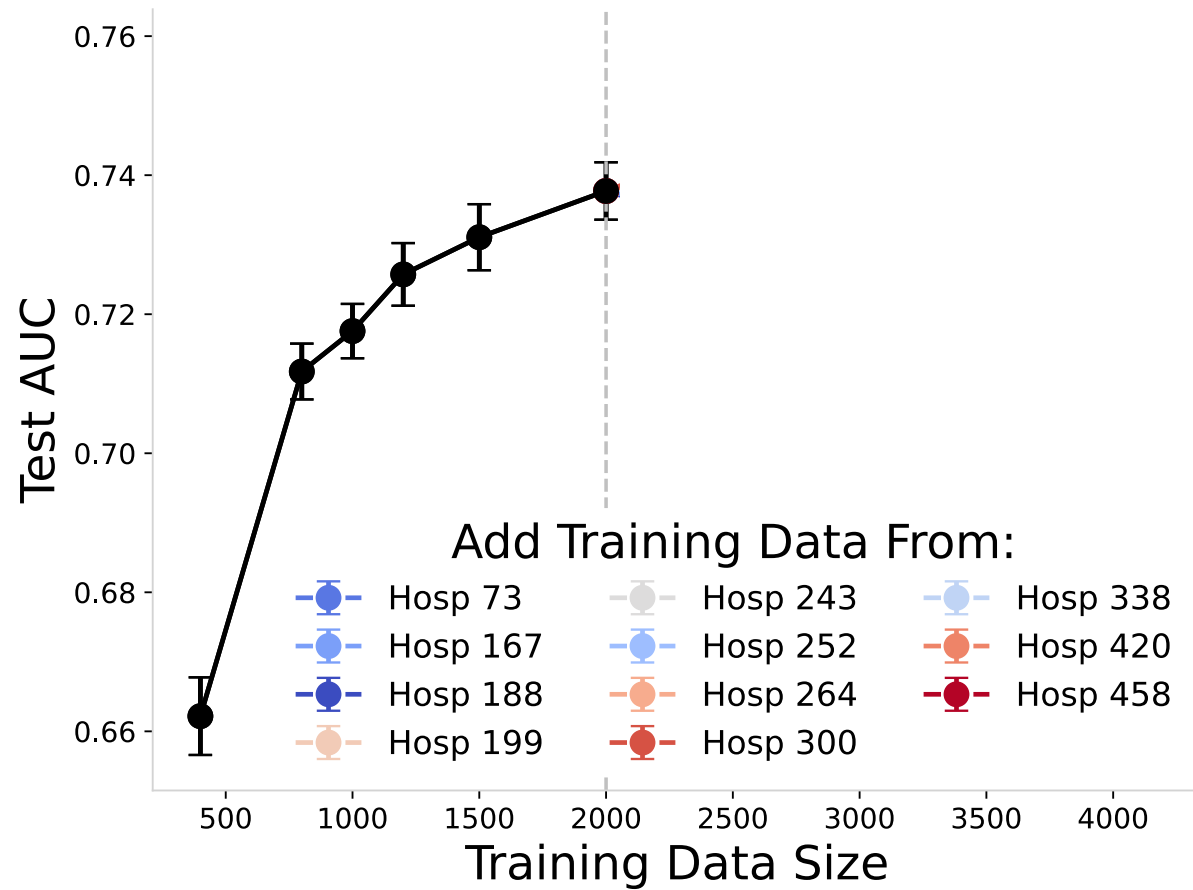


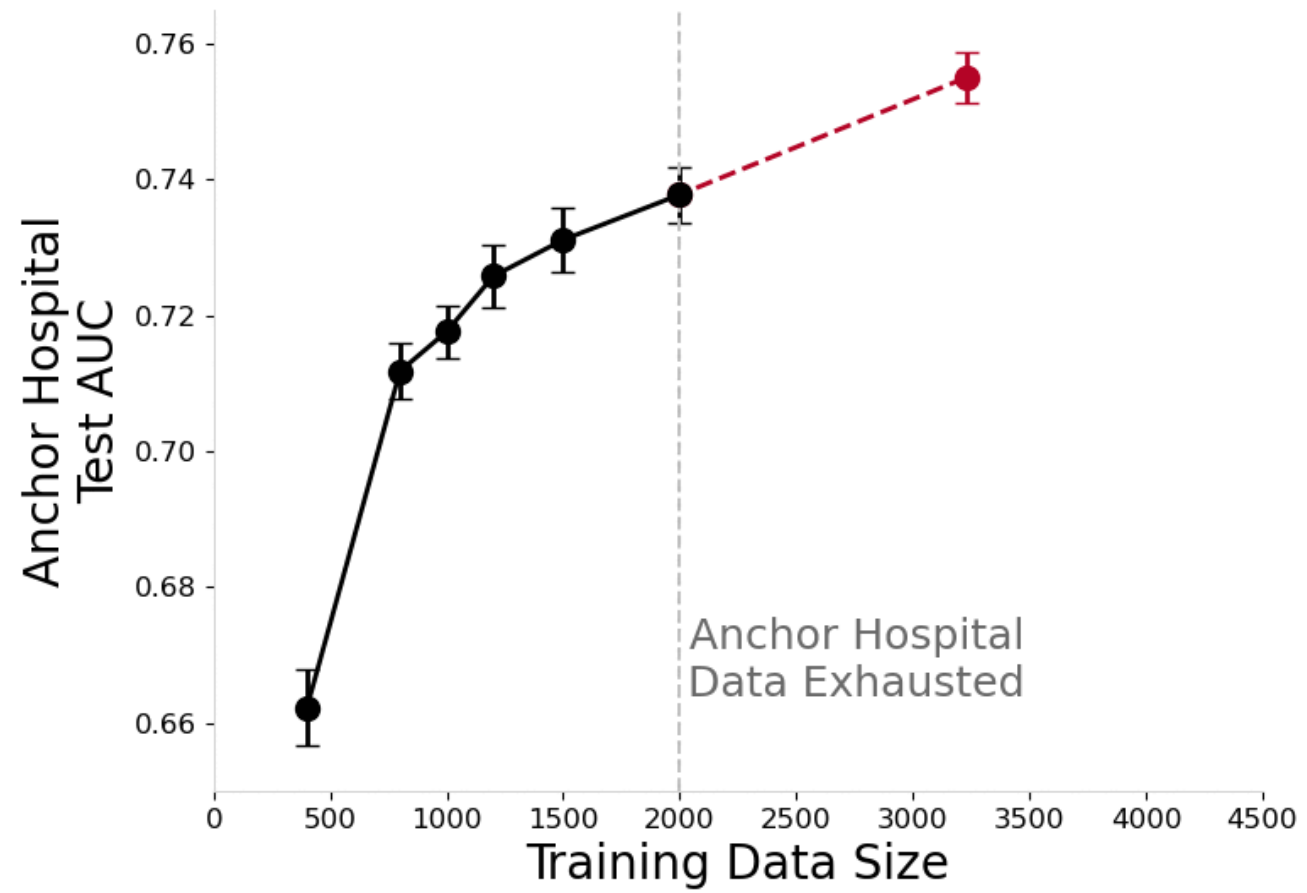


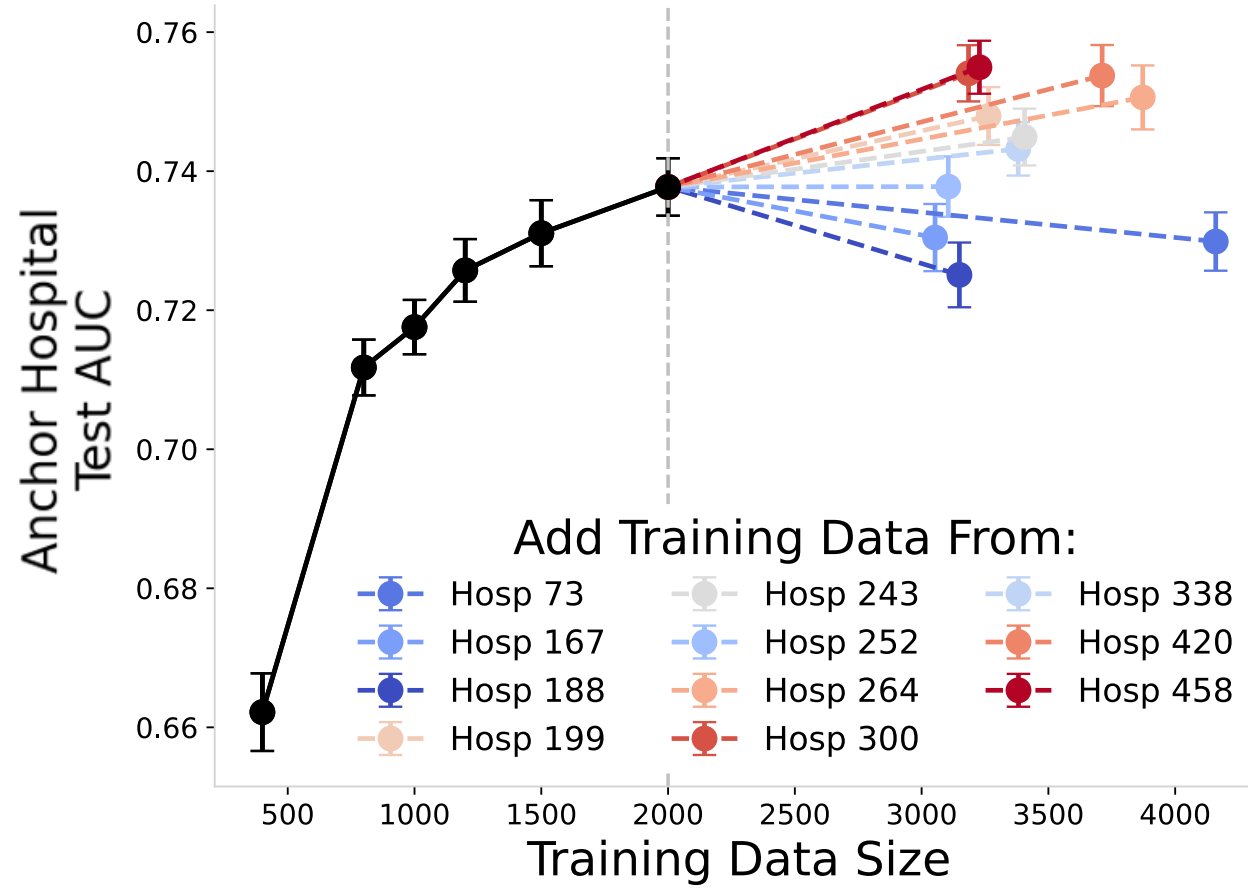


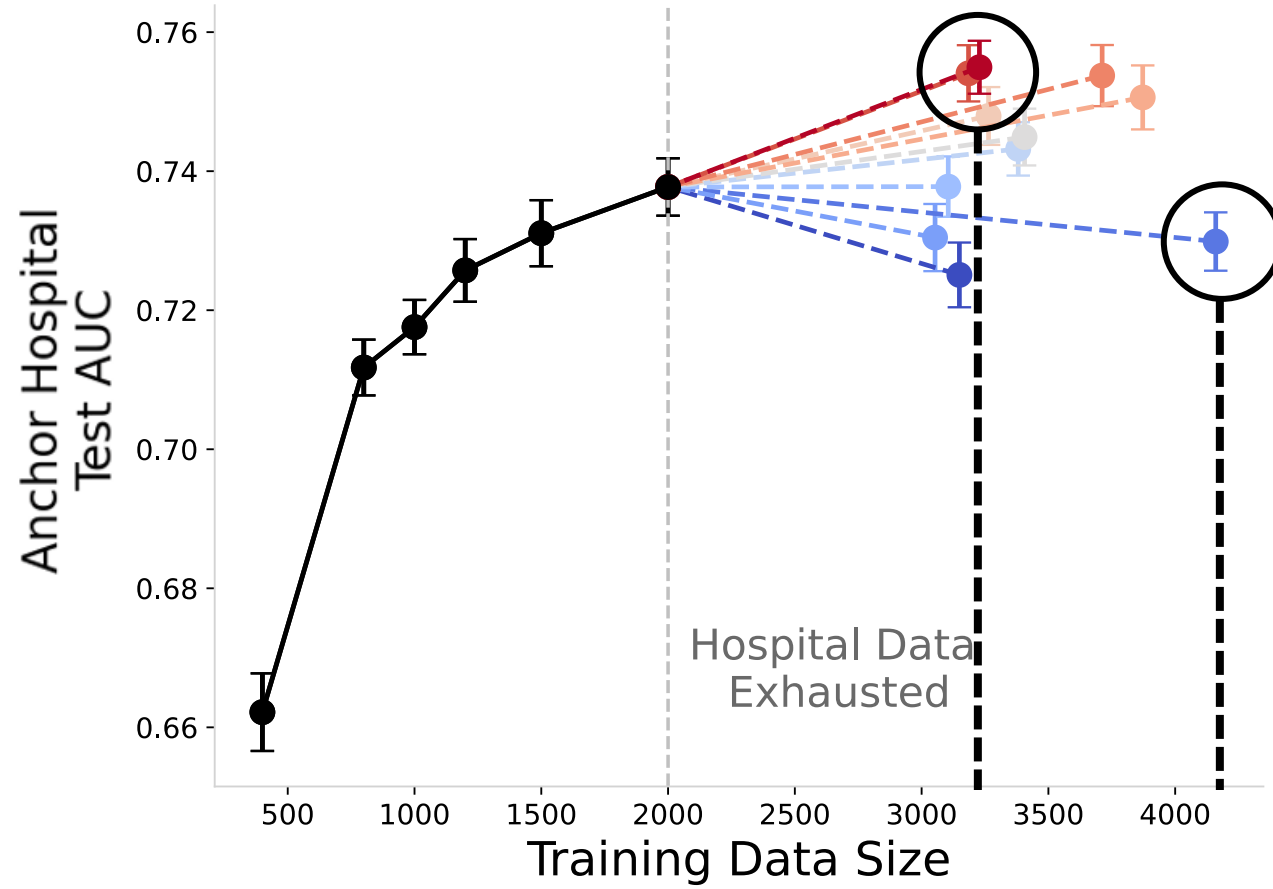




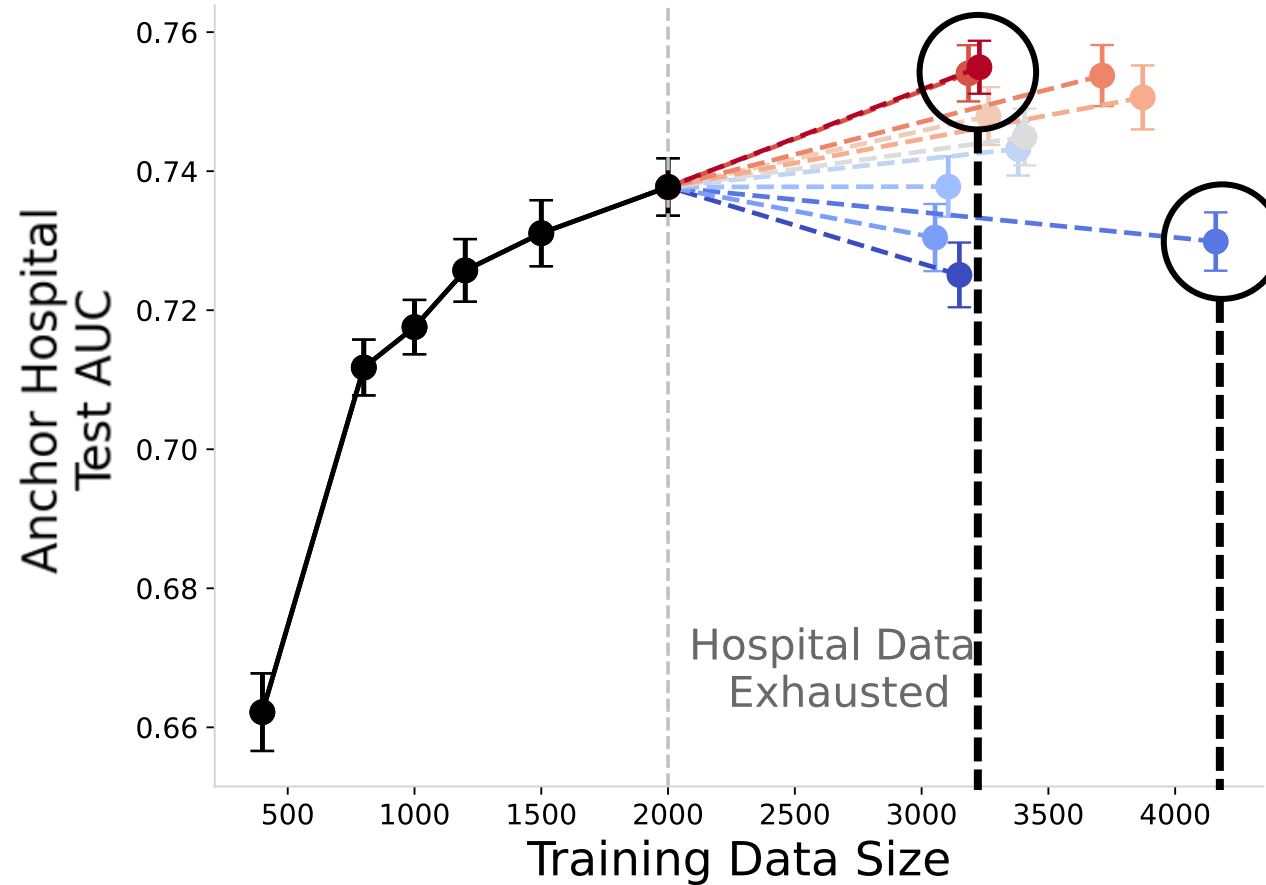








We call this the **Data Addition Dilemma**

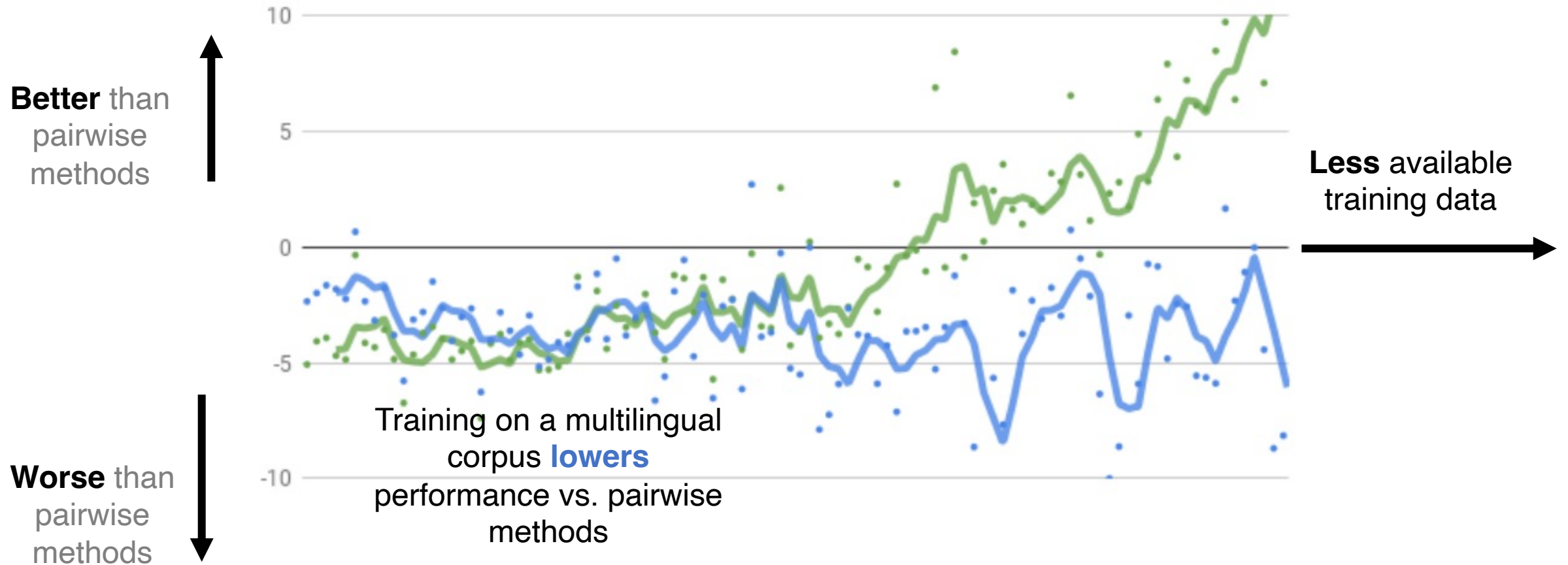


Also seen in: Multilingual Translation



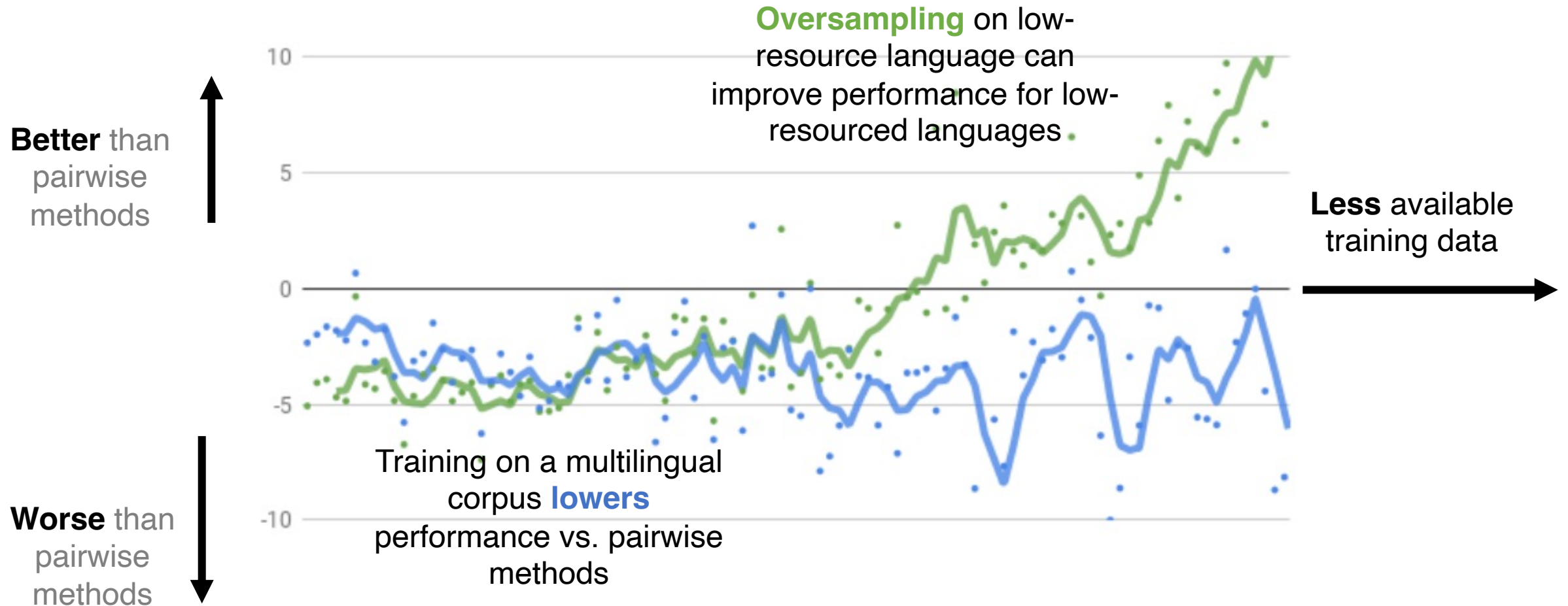
Arivazhagan et al, "Massively Multilingual Neural Machine Translation in the Wild: Findings and Challenges" NAACL 2019

Also seen in: Multilingual Translation



Arivazhagan et al, "Massively Multilingual Neural Machine Translation in the Wild: Findings and Challenges" NAACL 2019

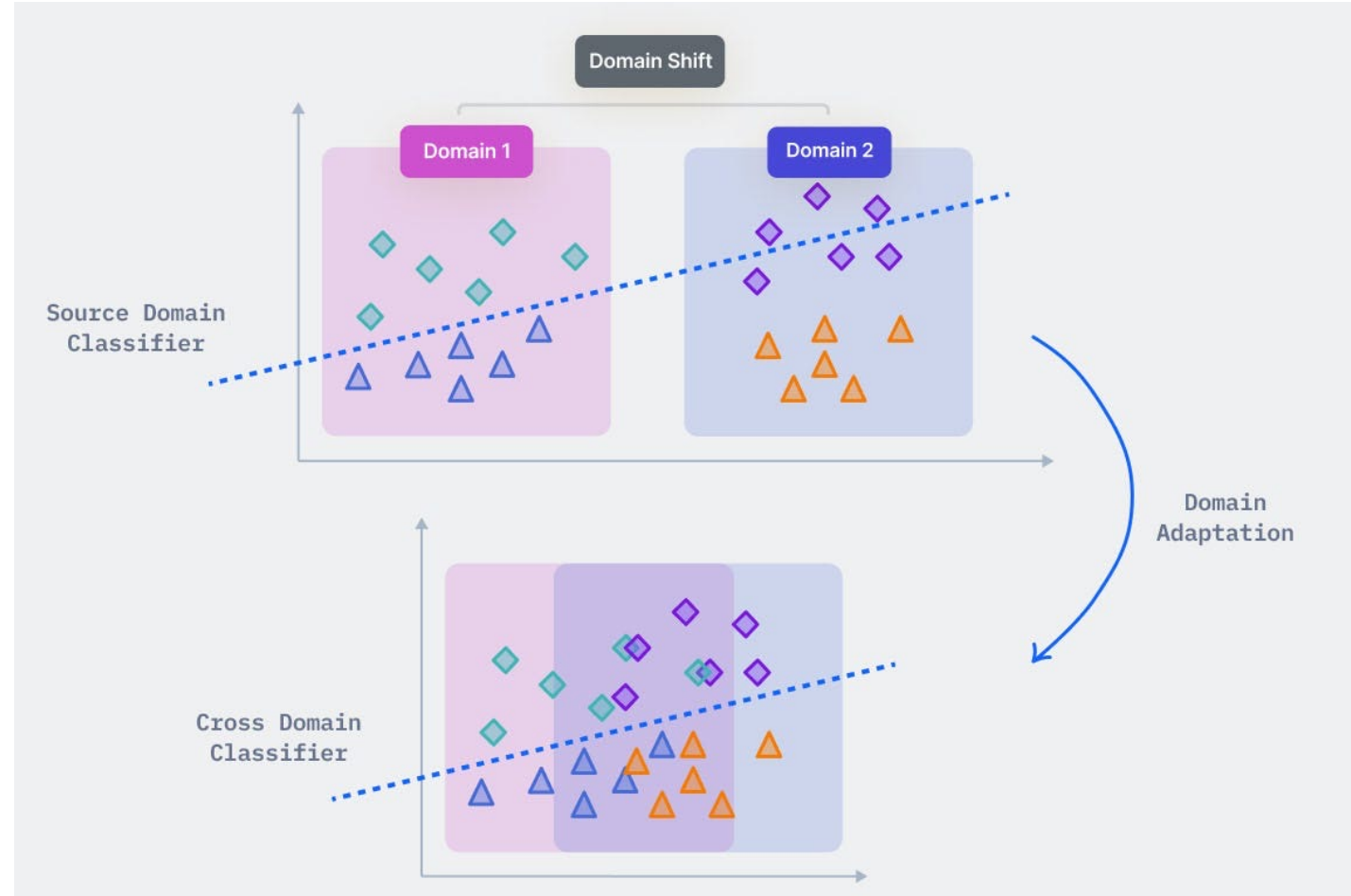
Also seen in: Multilingual Translation



Arivazhagan et al, "Massively Multilingual Neural Machine Translation in the Wild: Findings and Challenges" NAACL 2019

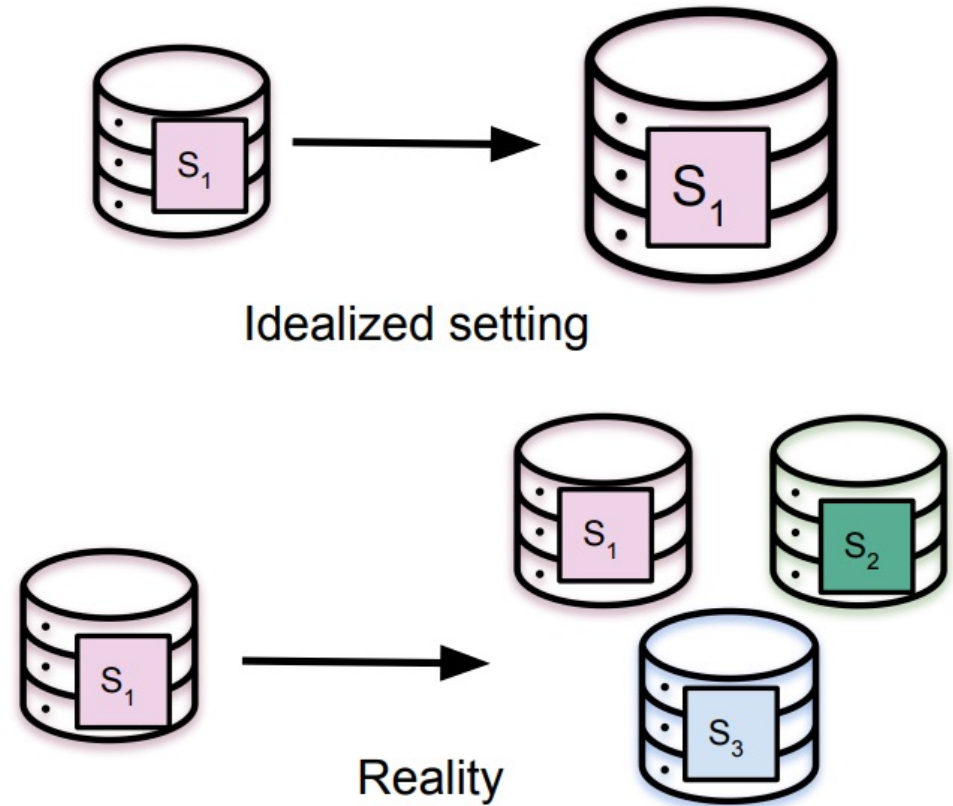
Related: Domain Adaptation

- DA goal is to make a model that has **high performance across multiple domains**
- Our problem formulation only is interested in performance in **one domain**



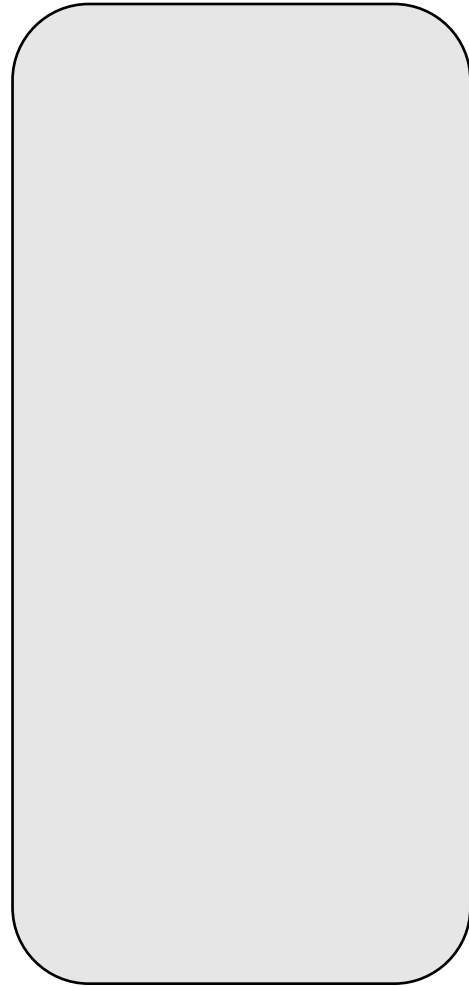
The Data Addition Dilemma

Framework to weigh **cost**
and benefit of adding
more data

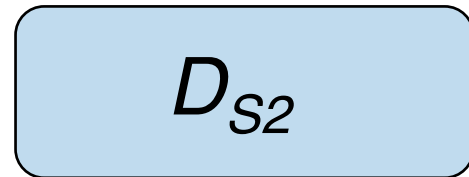
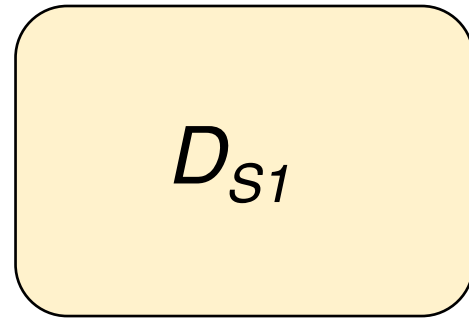


Modeling Data Accumulation

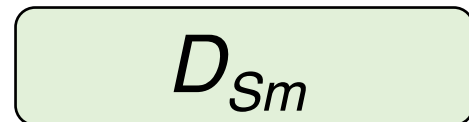
$$\{x, y\}^n \sim D$$



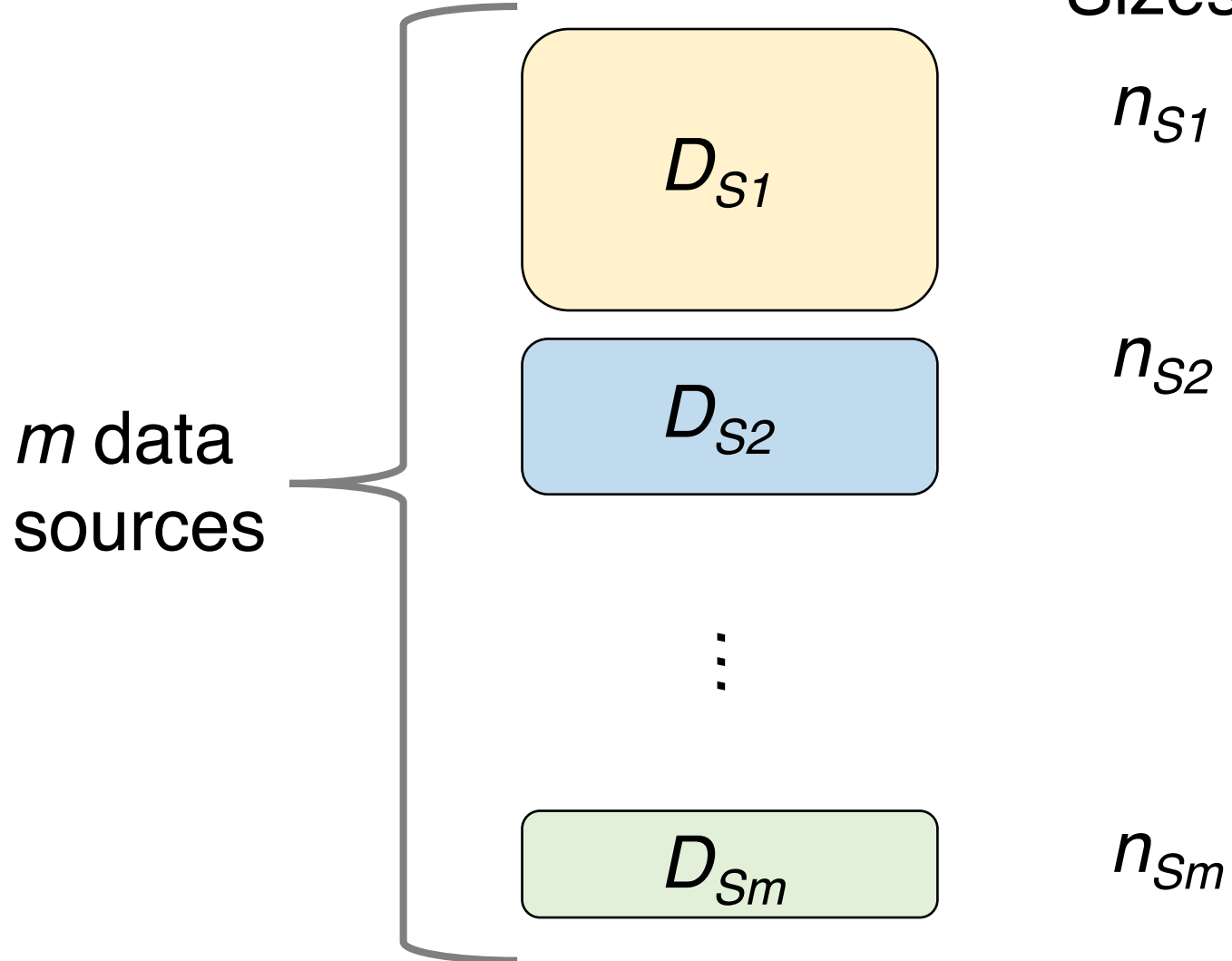
Modeling Data Accumulation



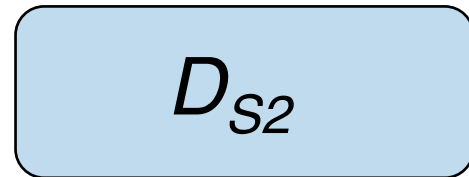
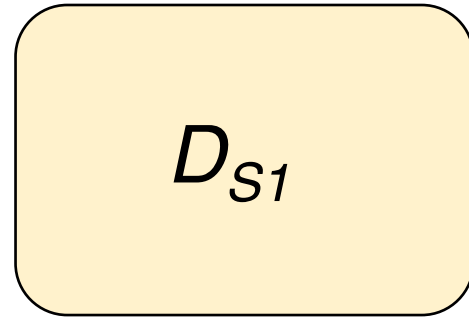
⋮



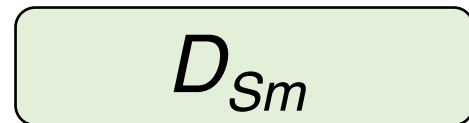
Modeling Data Accumulation Sizes



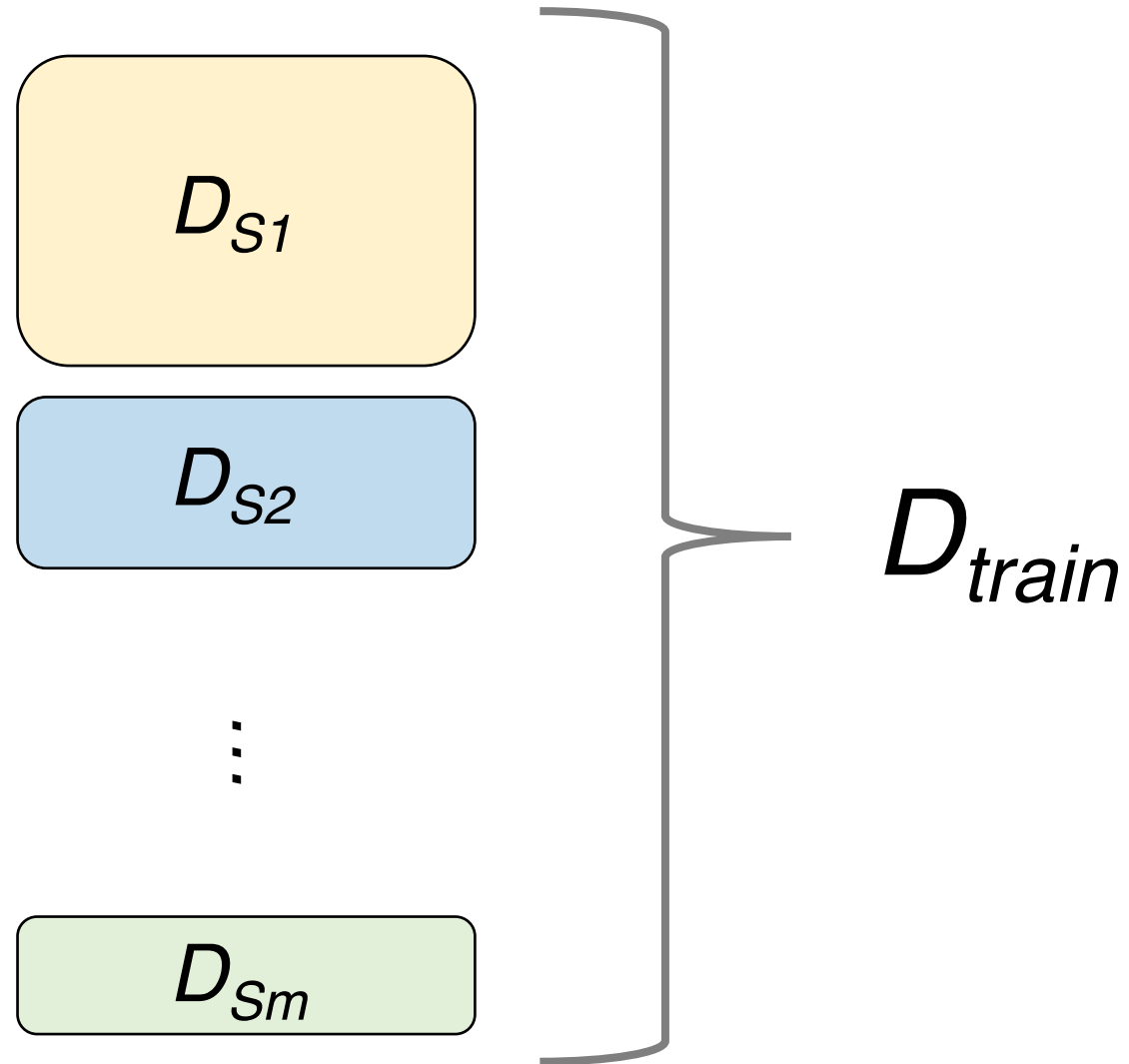
Modeling Data Accumulation



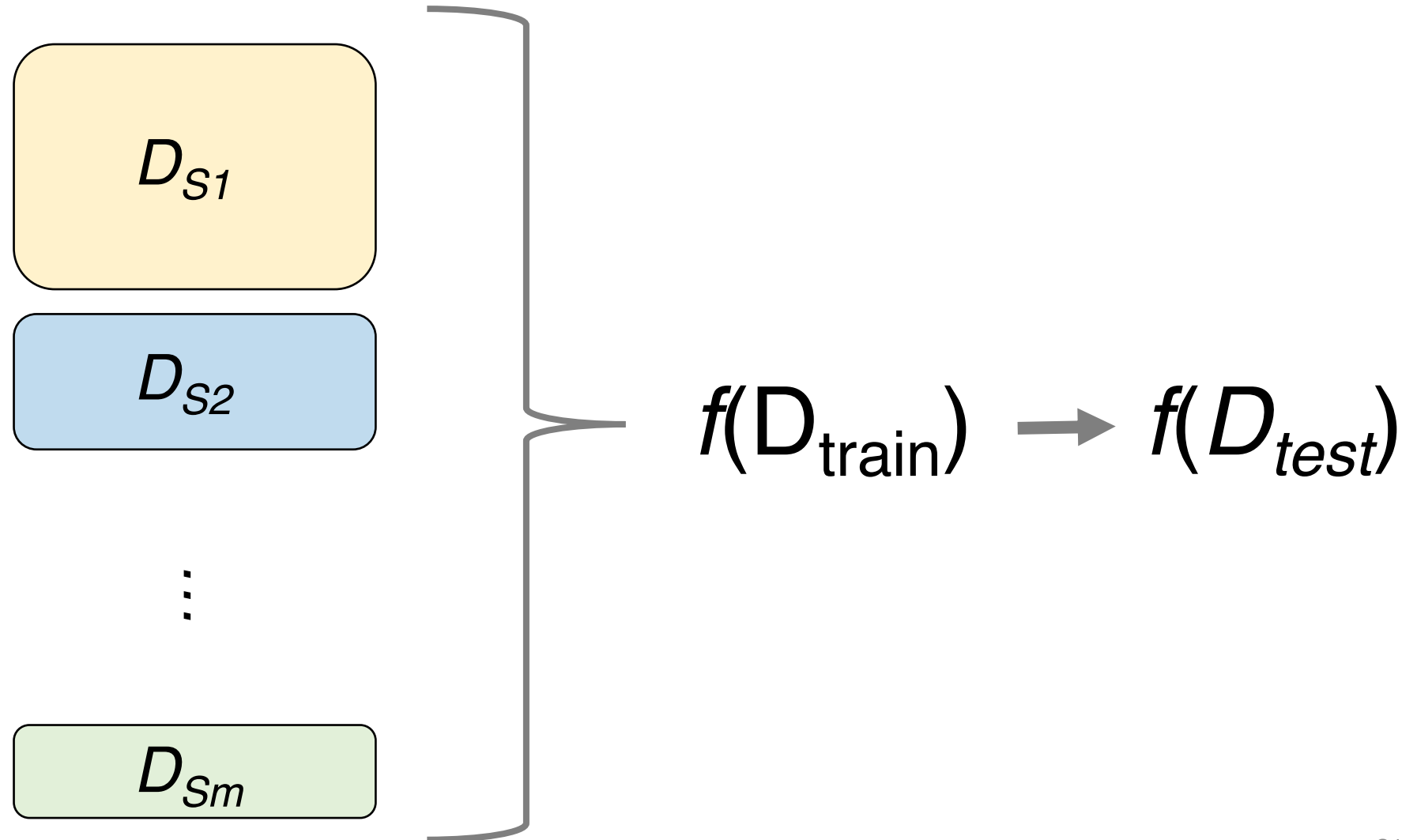
⋮



Modeling Data Accumulation



Modeling Data Accumulation



$$\text{AUC}(D_{H_0} \cup \dots \cup D_{H_M}, D_{H_{\text{test}}}) - \text{AUC}(D_{H_0} \cup \dots \cup D_{H_{M-1}}, D_{H_{\text{test}}})$$

AUC with adding last
data source

–

AUC without adding last
data source

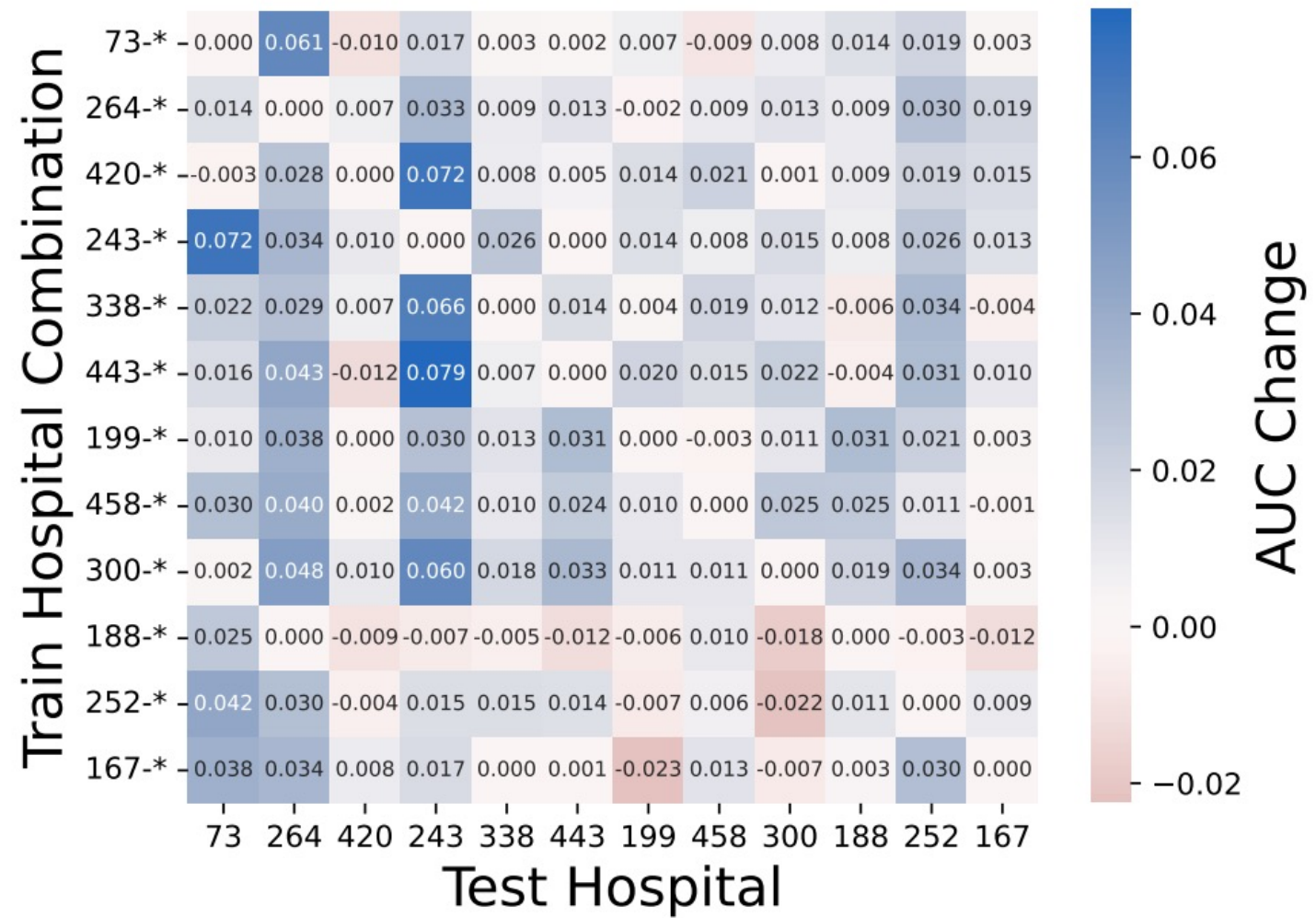
$$\text{AUC}(D_{H_0} \cup \dots \cup D_{H_M}, D_{H_{\text{test}}}) - \text{AUC}(D_{H_0} \cup \dots \cup D_{H_{M-1}}, D_{H_{\text{test}}})$$

\approx

$$\text{AUC}(D_{H_0} \cup H_1, D_{H_0}) - \text{AUC}(D_{H_0}, D_{H_0})$$

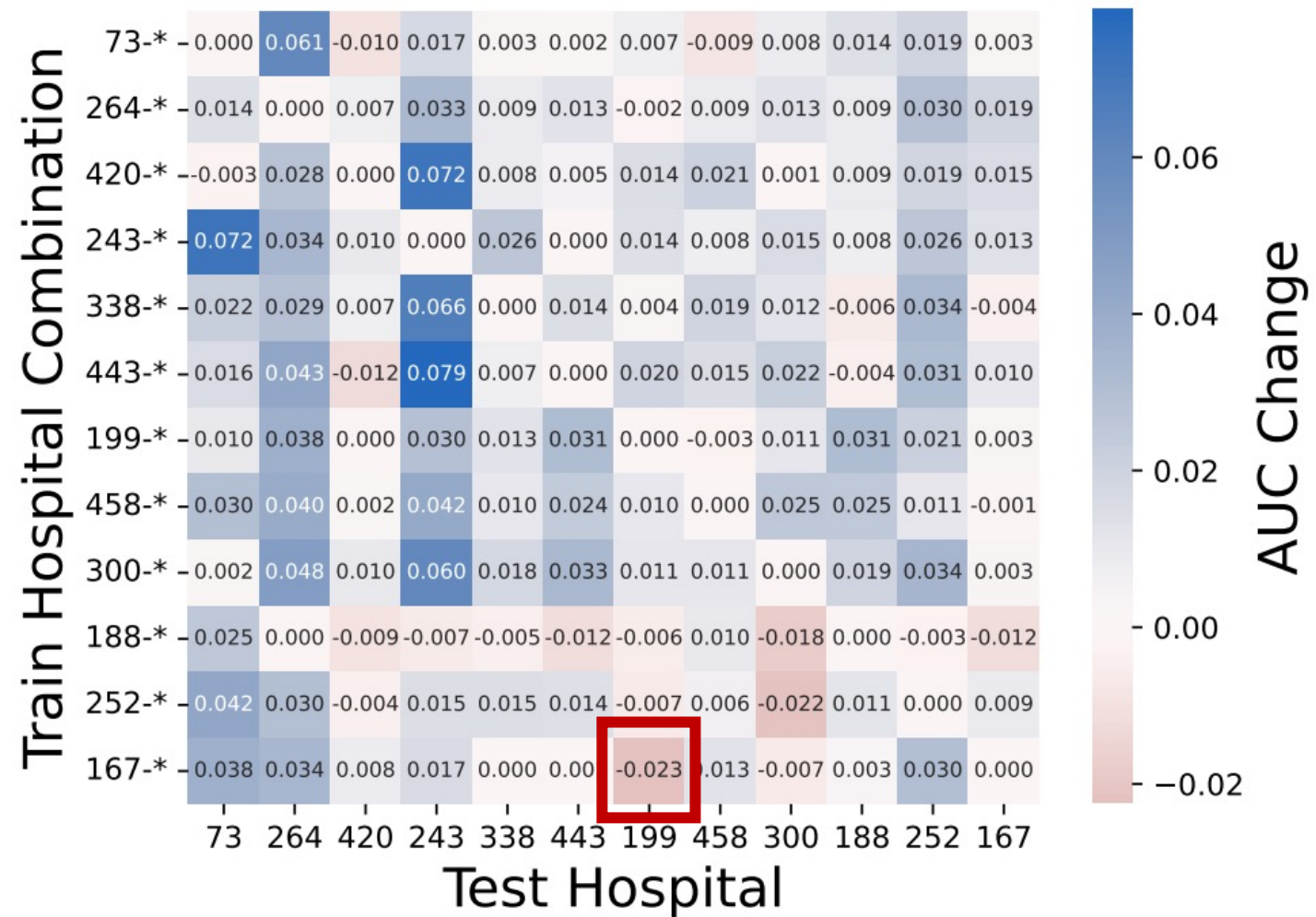
AUC with adding a
data source

AUC without adding a
data source



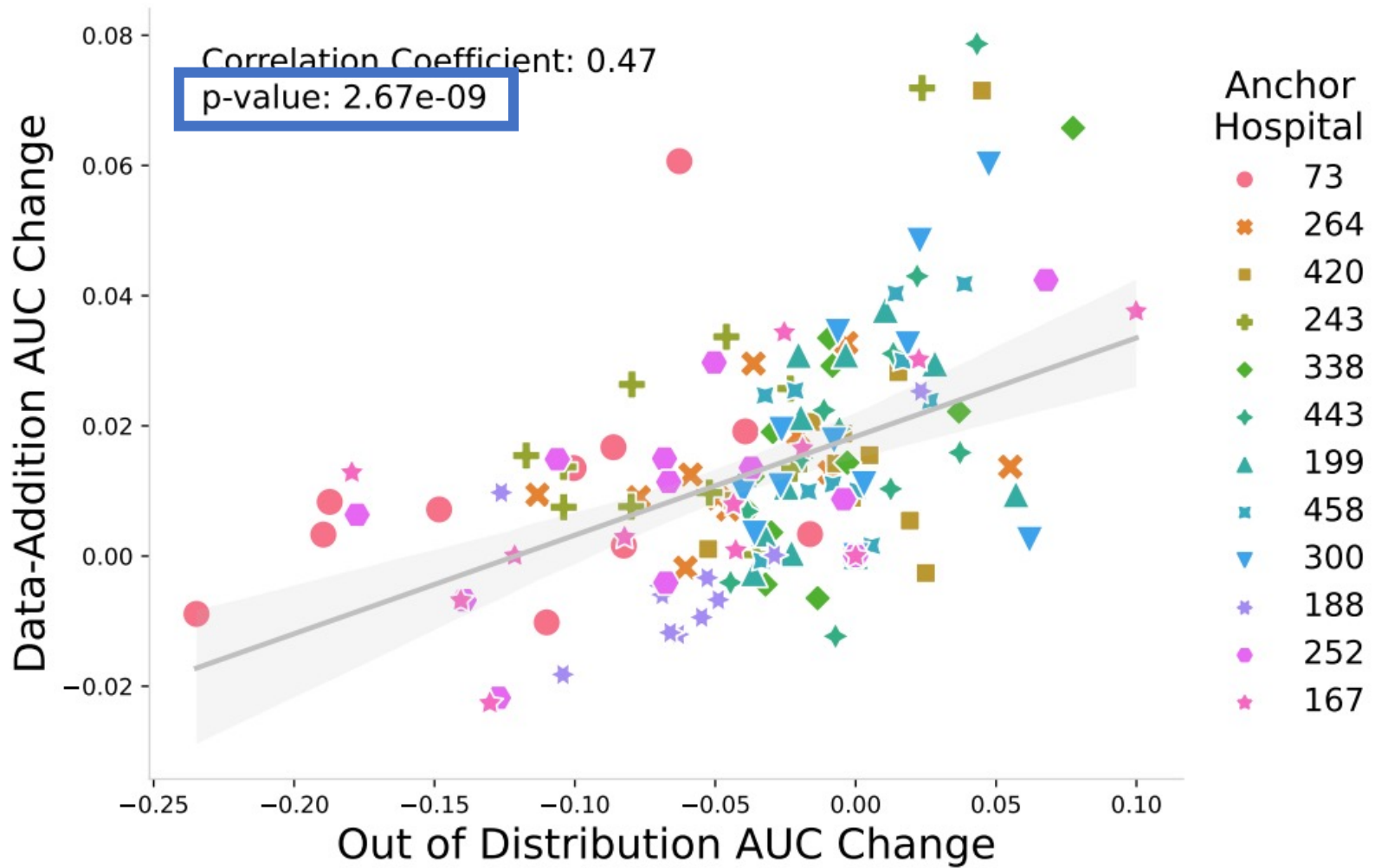
24-hour mortality prediction task on eICU dataset

Adding
hospital 167
 makes
 performance
worse for
hospital 199



24-hour mortality prediction task on eICU dataset

AUC change from adding hospitals is **related** to AUC change from out of distribution



$$\text{AUC}(D_{H_0} \cup \dots \cup D_{H_M}, D_{H_{\text{test}}}) - \text{AUC}(D_{H_0} \cup \dots \cup D_{H_{M-1}}, D_{H_{\text{test}}})$$

$$\approx$$

$$\text{AUC}(D_{H_0} \cup D_{H_1}, D_{H_0}) - \text{AUC}(D_{H_0}, D_{H_0})$$

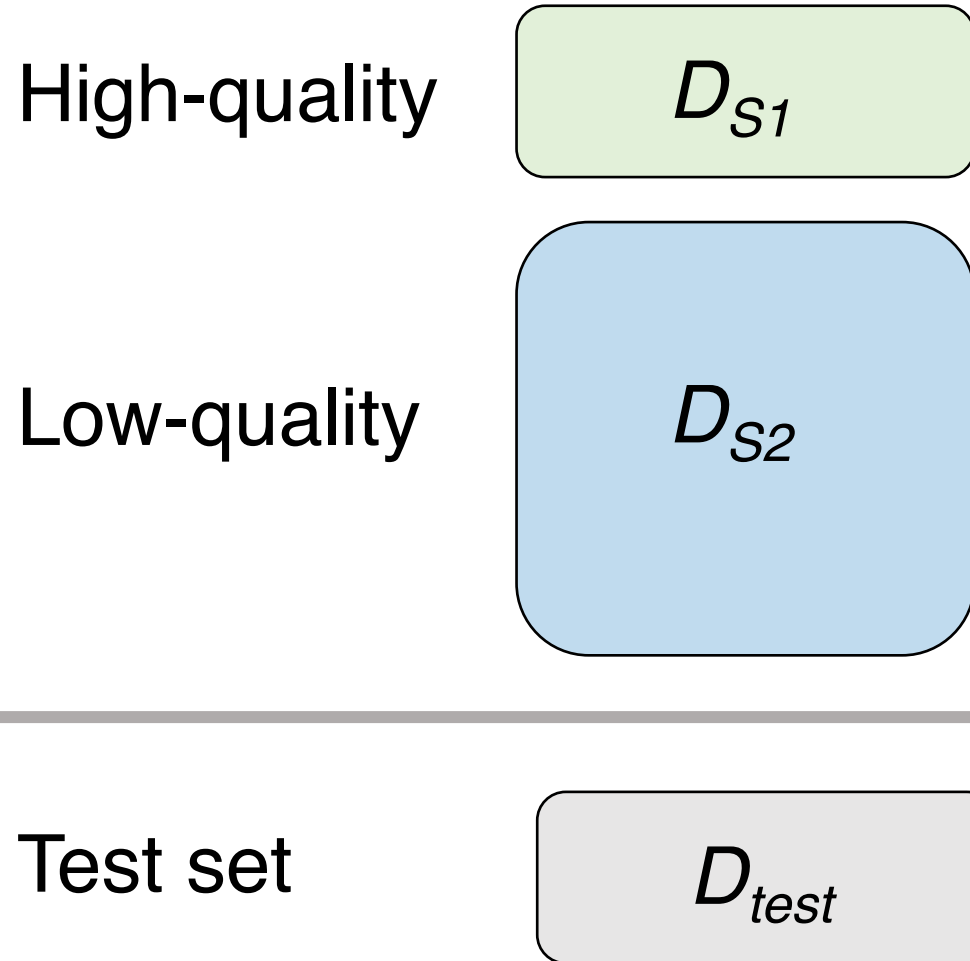
$$\approx$$

$$\delta(D_{H_0}, D_{H_1})$$

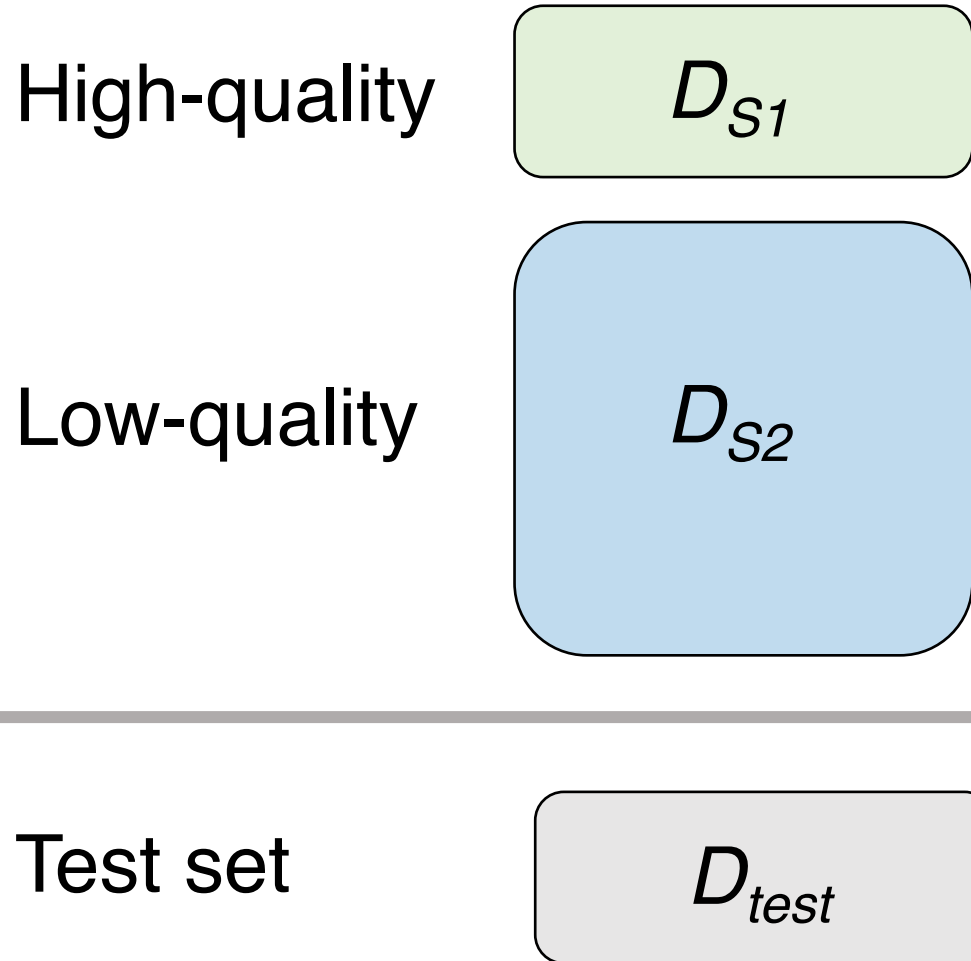
Divergence between
datasets

Toy Example

What is $\delta(D_{train}, D_{test})$?



Toy Example

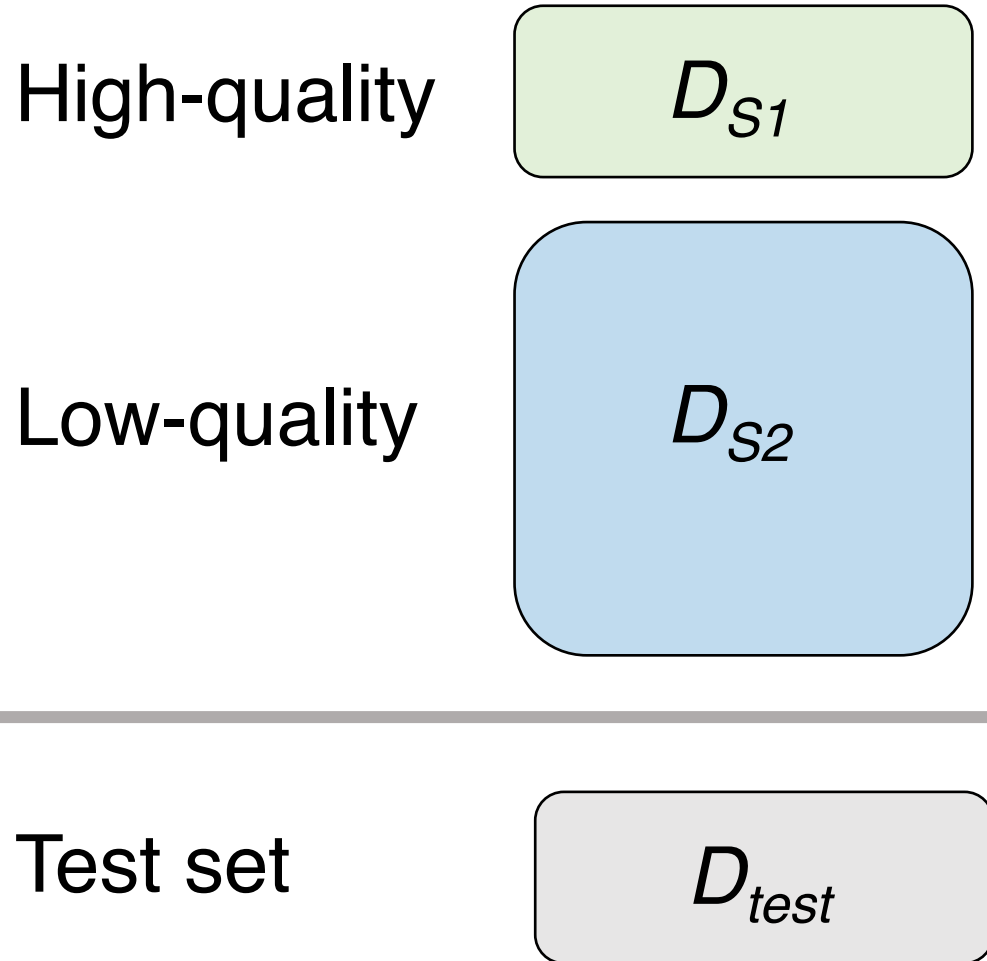


What is $\delta(D_{train}, D_{test})$?

If training data size $\leq n_{S1}$:

$$\delta(D_{S1}, D_{test})$$

Toy Example



What is $\delta(D_{train}, D_{test})$?

If training data size $\leq n_{S1}$:

$$\delta(\mathbf{D}_{S1}, D_{test})$$

If training data size $> n_{S1}$:

$$\frac{n_{S1}}{n} \delta(\mathbf{D}_{S1}, D_{test}) + \left(1 - \frac{n_{S1}}{n}\right) \delta(\mathbf{D}_{S2}, D_{test})$$

*If divergence δ is composed linearly

Contribution: When is adding more data hurtful?

Lemma: If $\delta(D_{Sk}, D_{test}) - \frac{cn}{n_{Sk}} \geq \delta(D_{train}, D_{test})$:

$$\delta(D_{train,n}, D_{test}) \geq \delta(D_{train,n-n_{Sk}}, D_{test})$$

D_{Sk} = data from k -th source

δ = divergence in family of f -divergences
(not assumed linear)

c = divergence-dependent constant

Related: Distances, divergences, and discrepancies

- Relationship between increased divergence to empirical risk for f -divergences by giving a generalization bound
- Different discrepancy measures including L1 distance
- $H\Delta H$ divergence
- Margin disparity discrepancy

Practical Strategies for Data Accumulation

Practical Strategies for Data Accumulation

Idea 1: Try all combinations

Practical Strategies for Data Accumulation

Idea 1: Try all combinations

Best 3 hospitals from 10 options
 \approx **10^3 combinations**

Practical Strategies for Data Accumulation

~~Idea 1: Try all combinations~~

Idea 2: Hospital-level characteristics
(mortality rates, demographics, etc.)

Practical Strategies for Data Accumulation

~~Idea 1: Try all combinations~~

Idea 2: Hospital-level characteristics

Idea 3: Divergence-based metrics

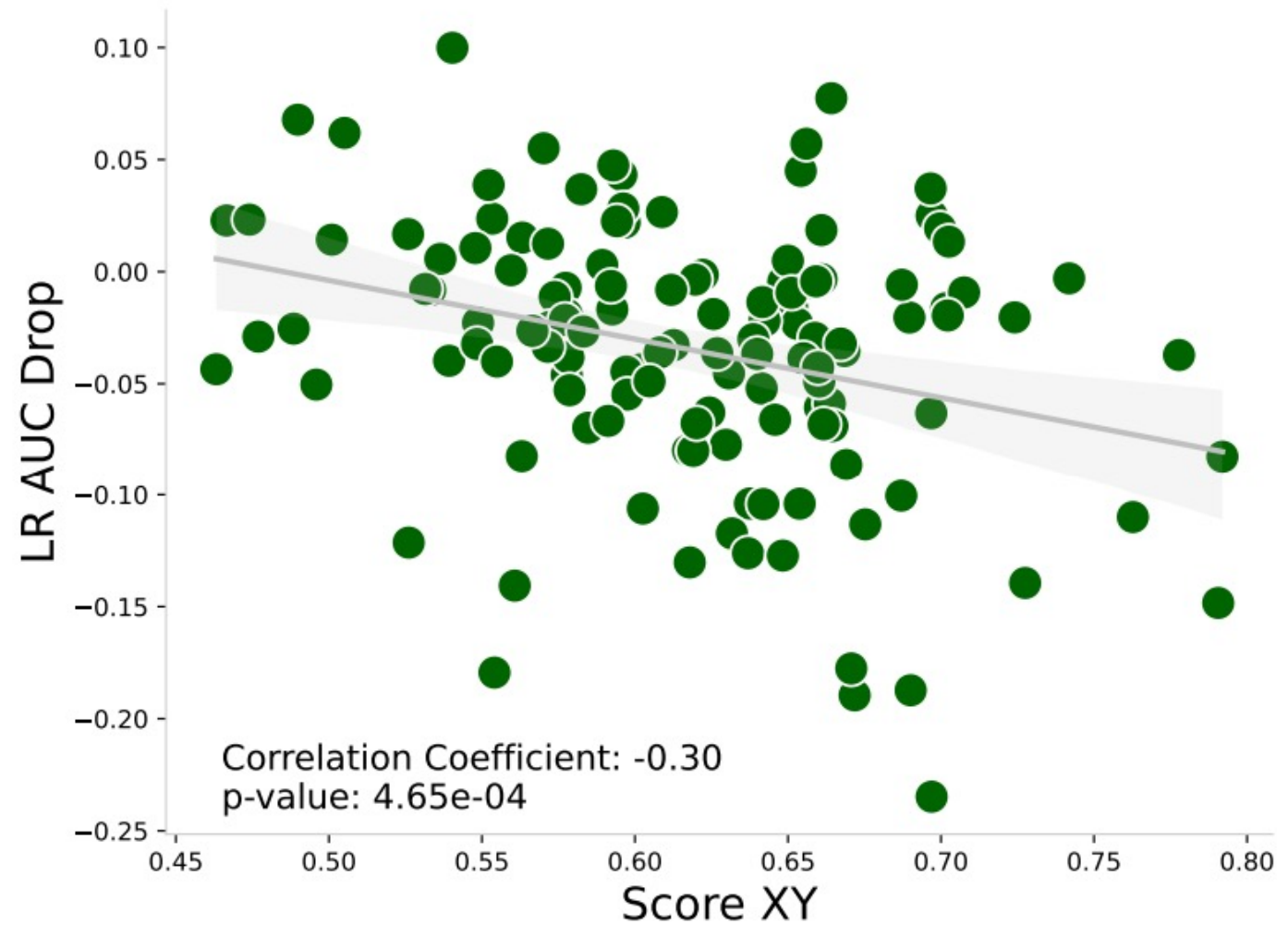
$$\text{KL Ratio } X = E_{x \in P} \left(\log \frac{P(x)}{Q(x)} \right)$$

$$\text{KL Ratio } XY = E_{(x,y) \in P} \left(\log \frac{P(x,y)}{Q(x,y)} \right)$$

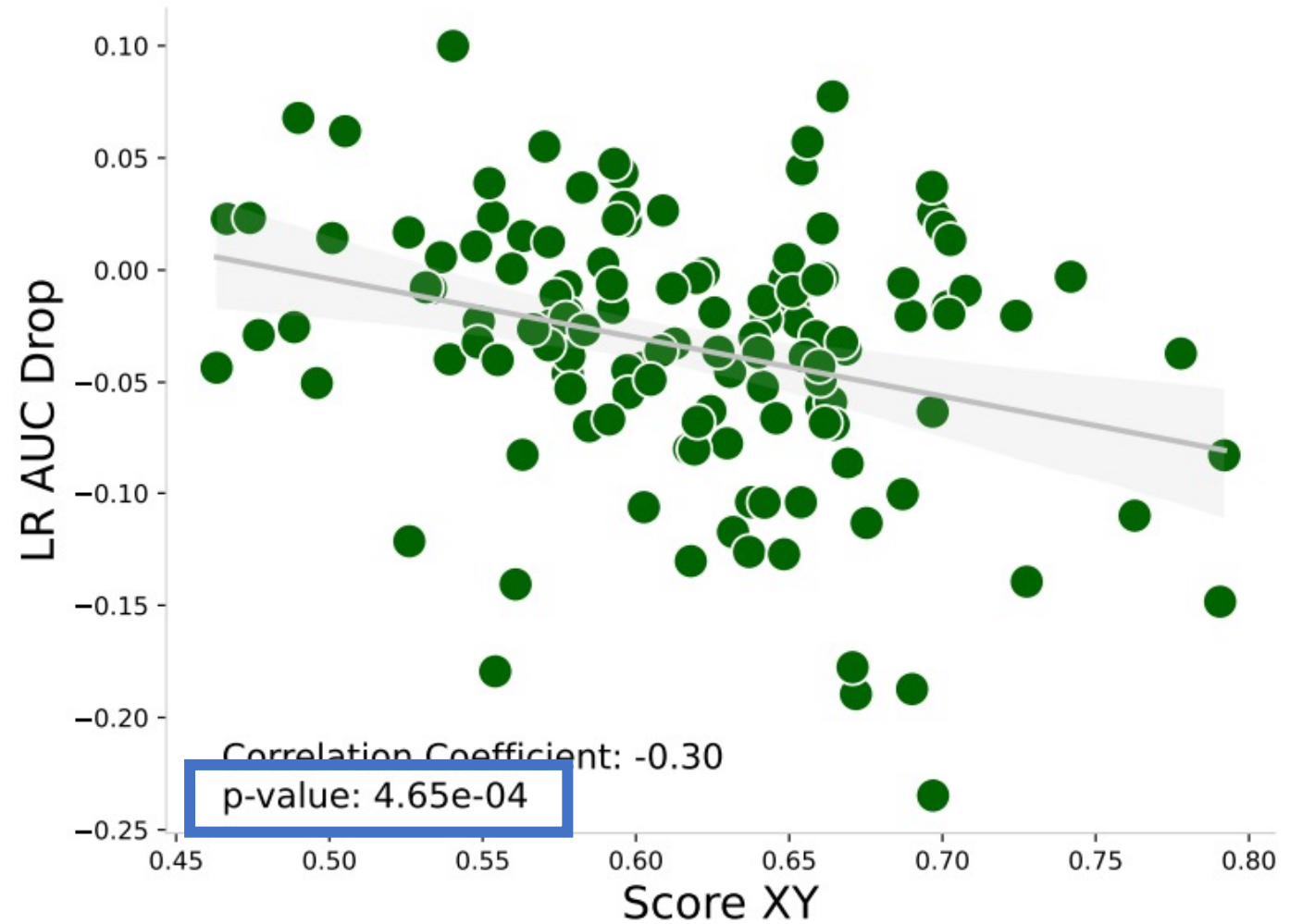
$$\text{Score } X = E_{x \in P} (P(x \in P))$$

$$\text{Score } XY = E_{(x,y) \in P} (P((x,y) \in P))$$

LR AUC drop is **correlated** with ScoreXY, a divergence-based metric

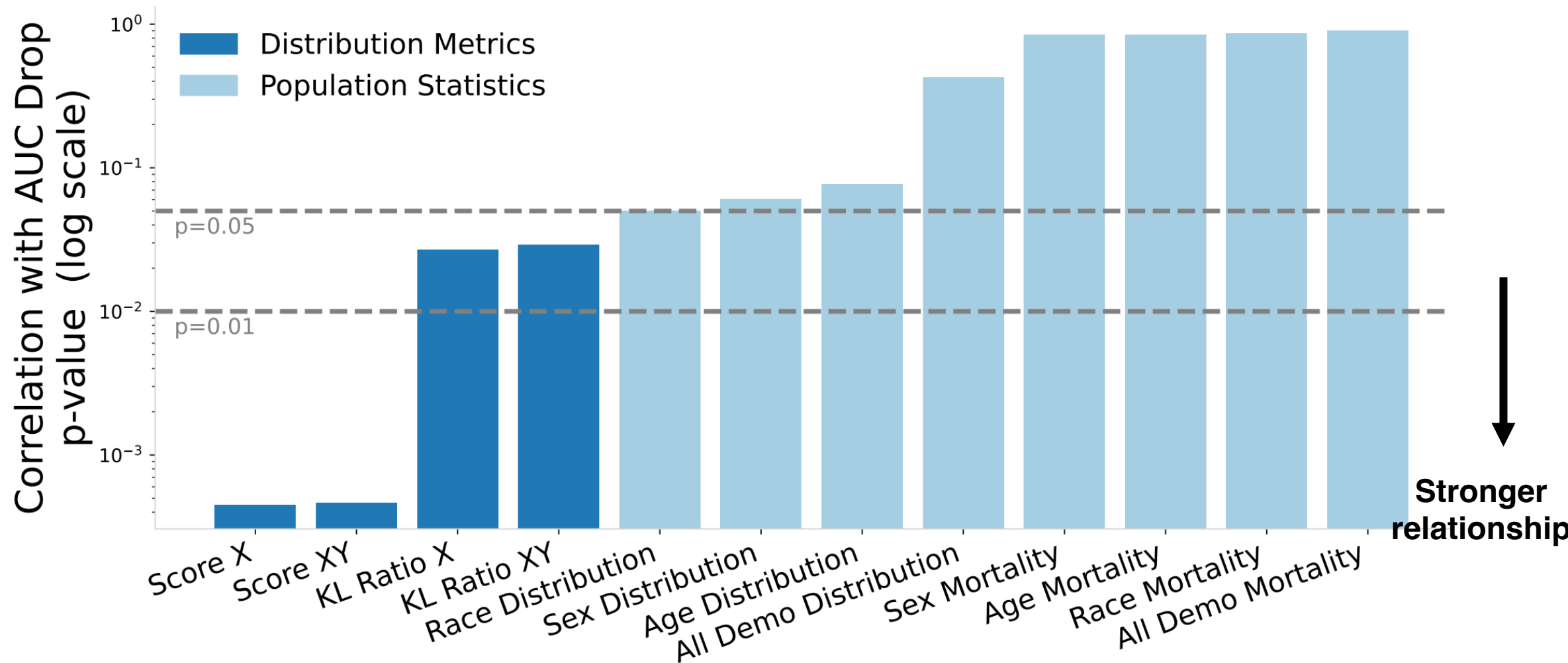


LR AUC drop is **correlated** with ScoreXY, a divergence-based metric

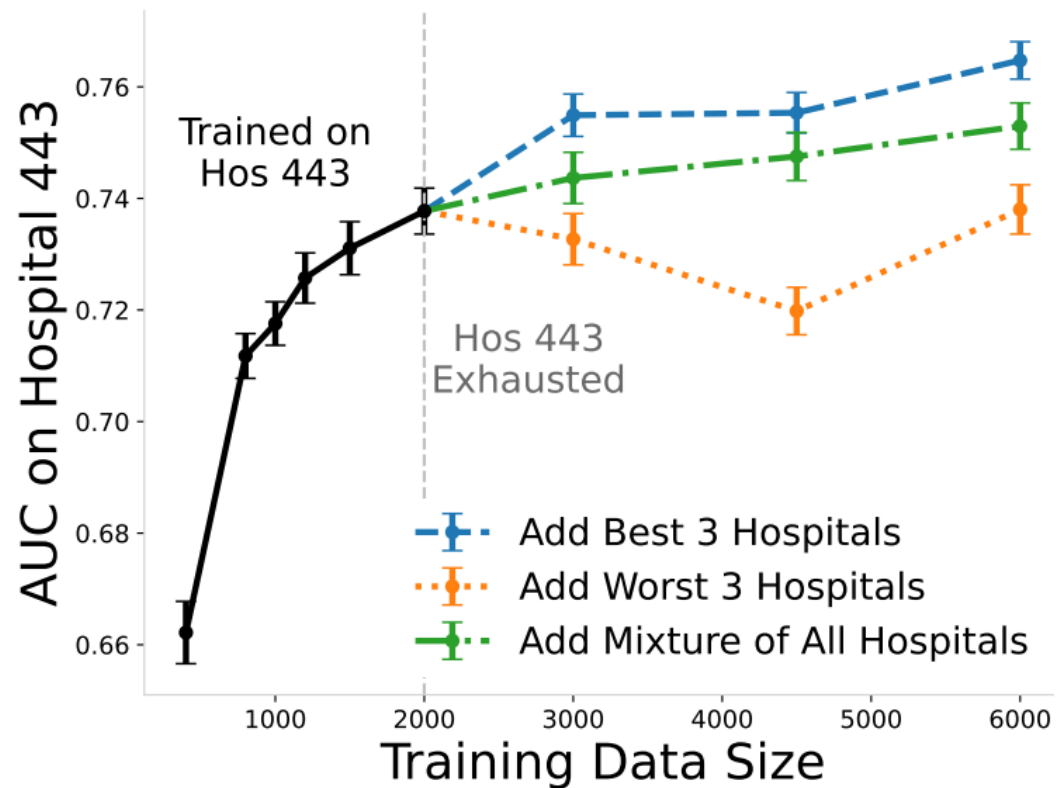




Divergence metrics out-perform other hospital-wide characteristics

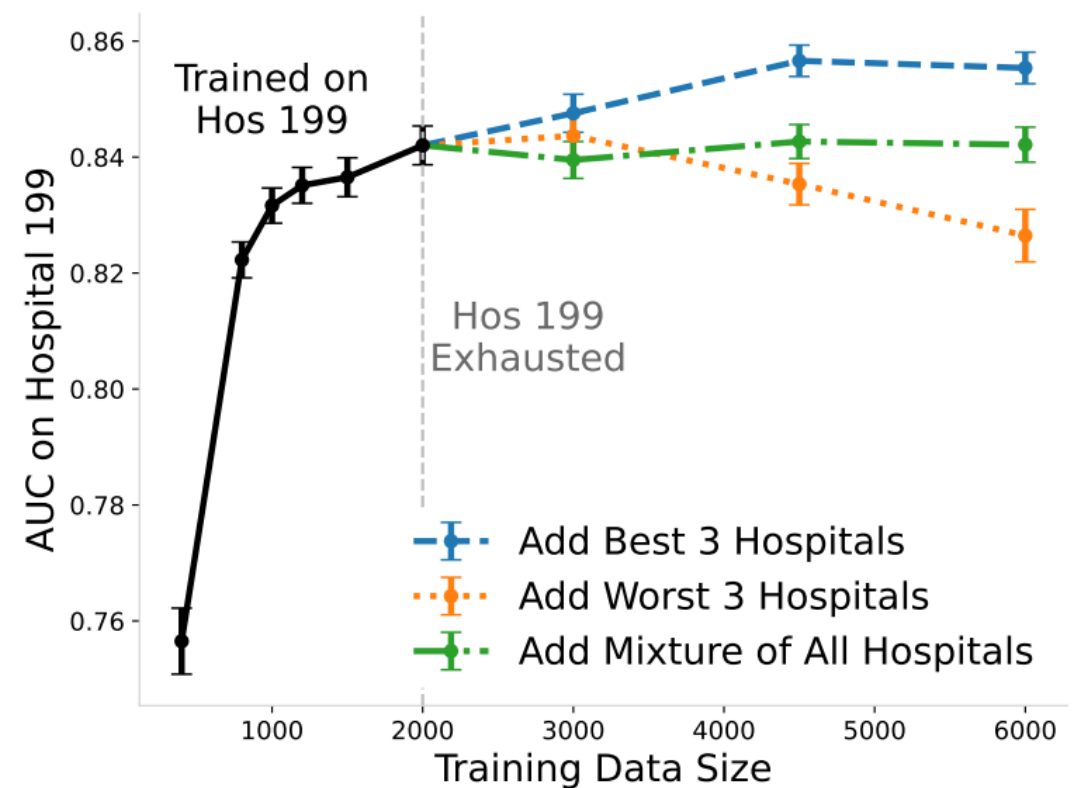
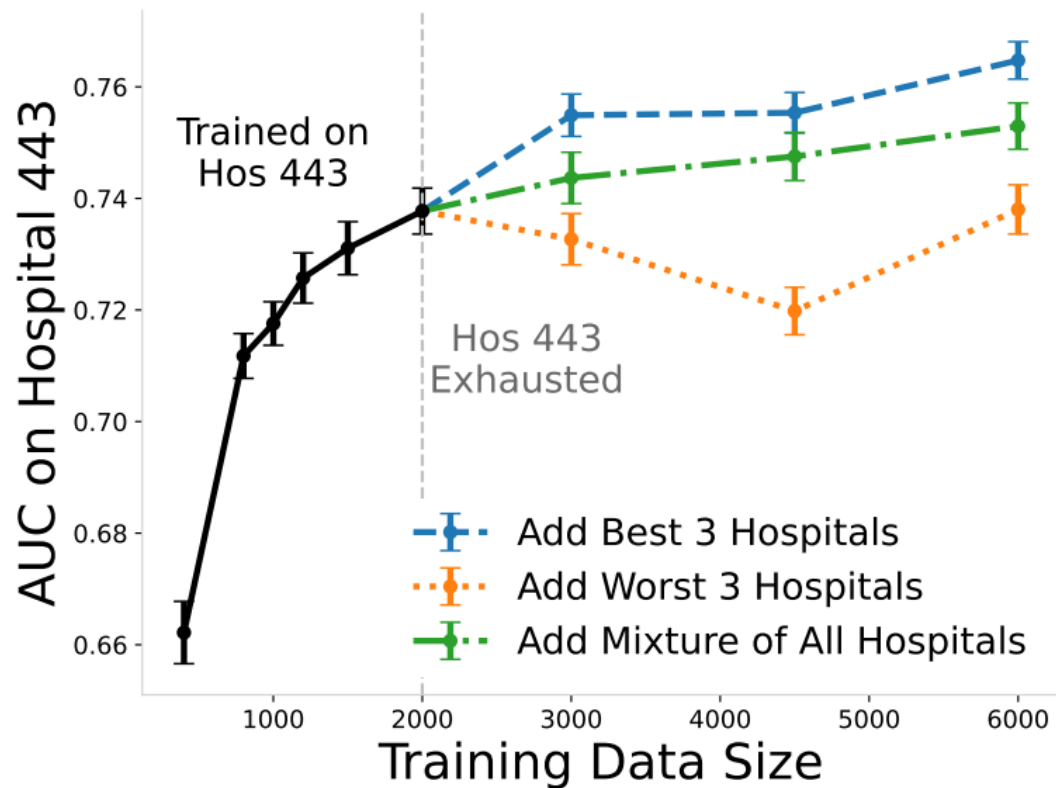


Divergence-based metrics outperform mixture settings



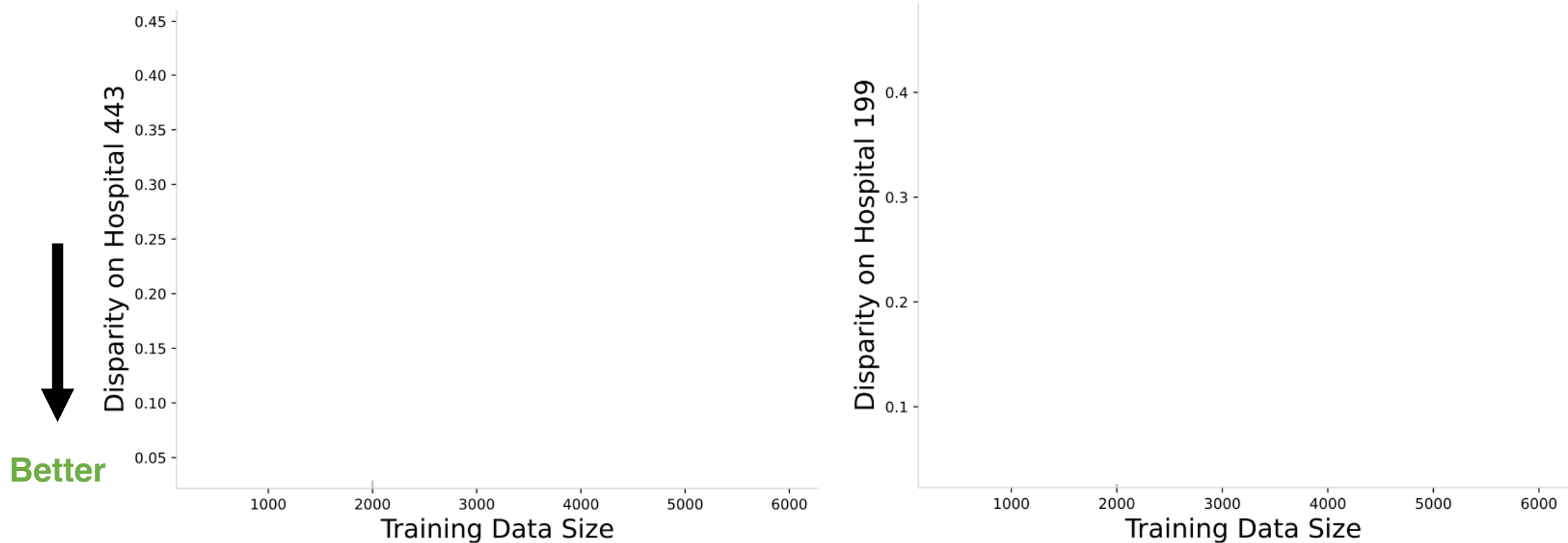
24-hour mortality prediction task on eICU dataset

Divergence-based metrics outperform mixture settings



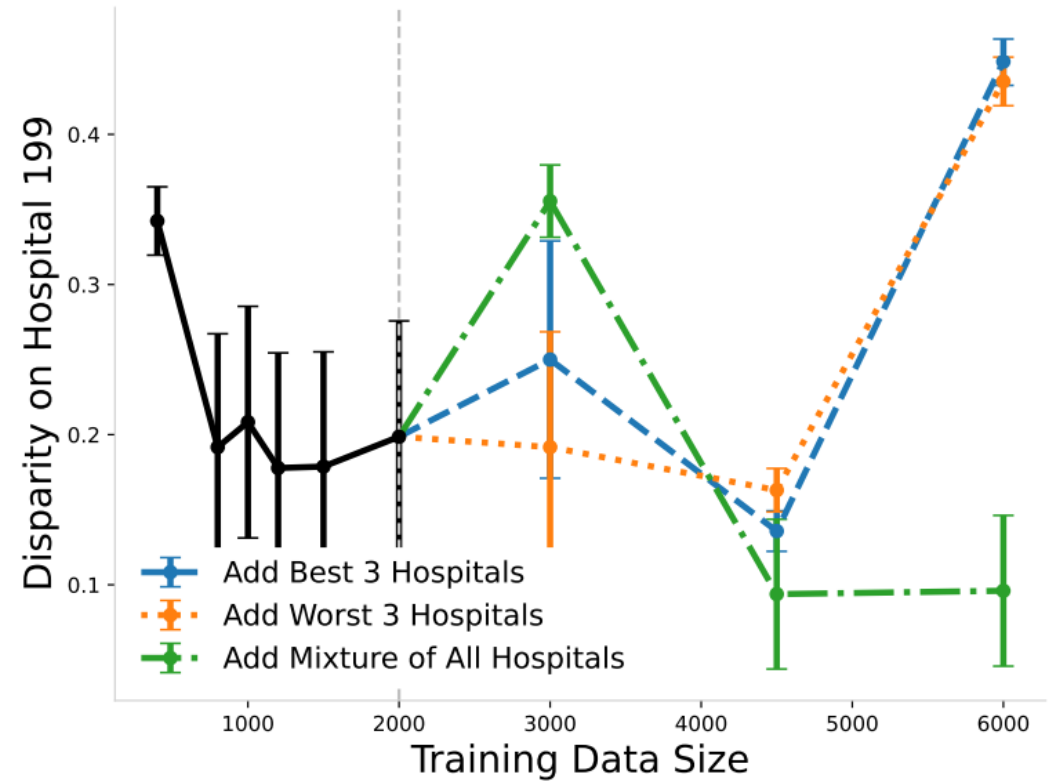
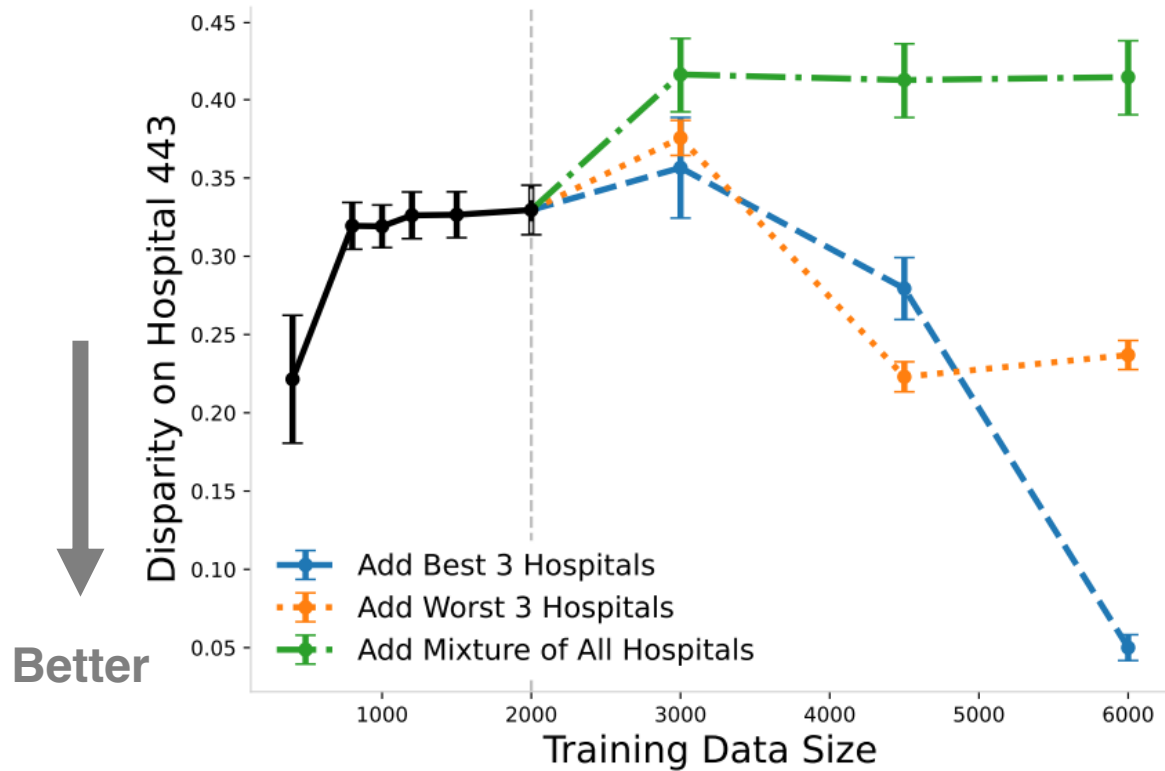
24-hour mortality prediction task on eICU dataset

But! Racial disparity can have mixed results



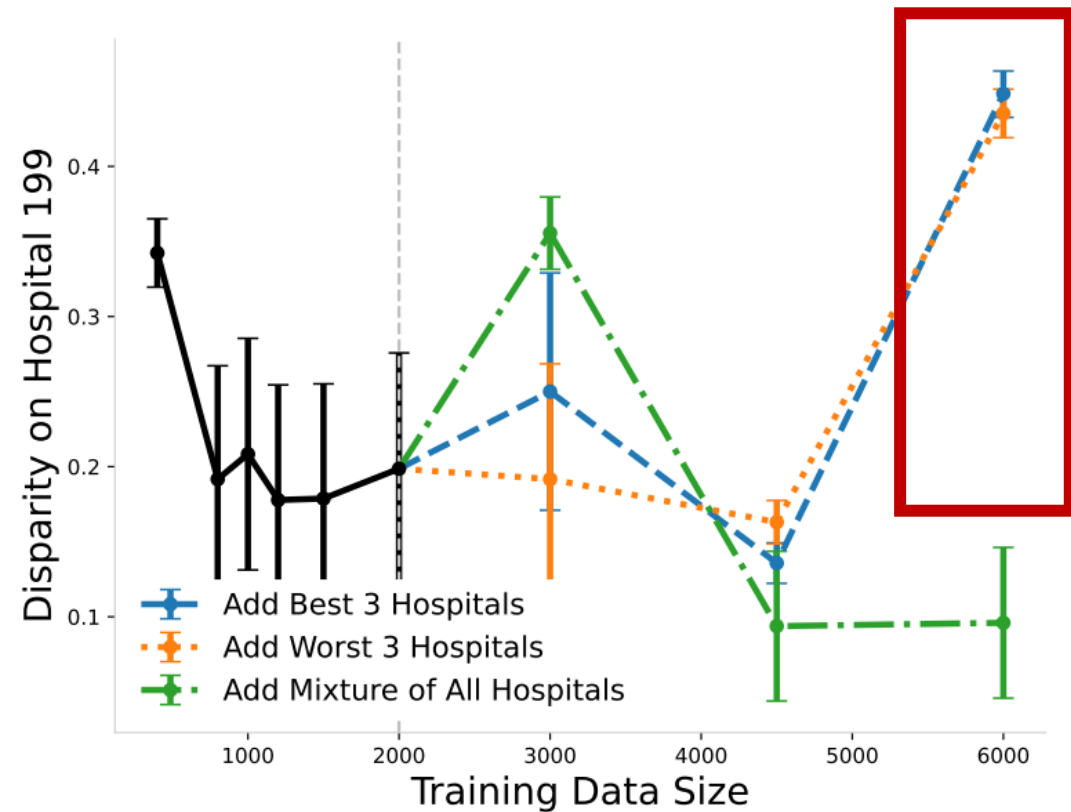
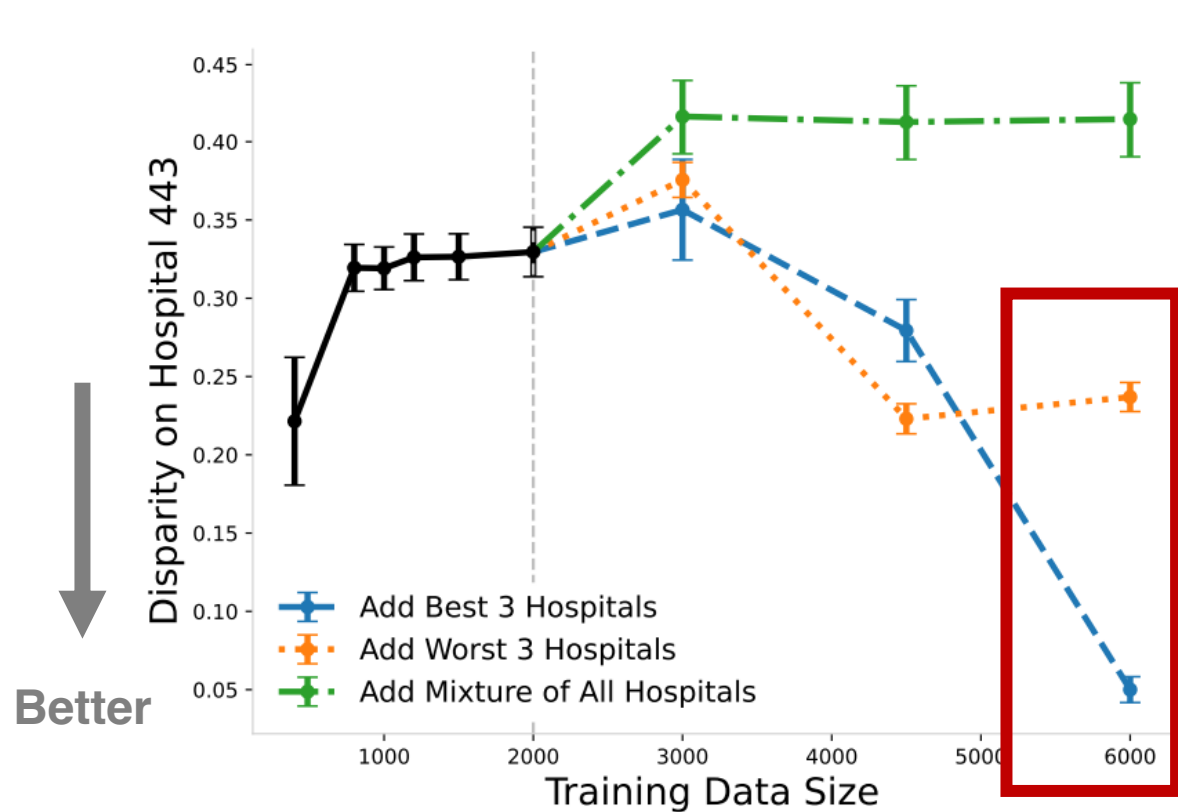
24-hour mortality prediction task on eICU dataset

But! Racial disparity can have mixed results



24-hour mortality prediction task on eICU dataset

But! Racial disparity can have mixed results



24-hour mortality prediction task on eICU dataset

How can we benefit from larger datasets?

1. Data composition is a **data-oriented perspective** to complement existing algorithmic-centered work, e.g., domain adaptation
2. We analyze different data accumulation strategies for health datasets with **multiple sources**
3. We show theoretical and empirical results, including a **divergence-based method** for when to add more data

How can we benefit from larger datasets?

1. Data composition is a **data-oriented perspective** to complement existing algorithmic-centered work, e.g., domain adaptation
2. We analyze different data accumulation strategies for health datasets with **multiple sources**
3. We show theoretical and empirical results, including a **divergence-based method** for when to add more data
4. **Next steps**: fairness, multimodality, impossibility theorems, privacy

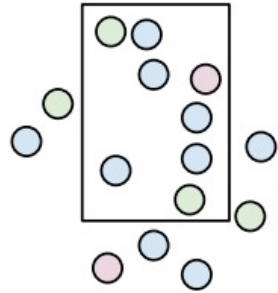
Computational Healthcare for Equity and INclusion

Problem Selection



1. Early detection for intimate partner violence ([PSB 2021](#))
2. Maternal health and high-risk pregnancy ([FaccT 2024](#))

Data Collection



1. Collecting and researching insurance risk scores ([Health Affairs 2022](#))
2. Effect of aggregating data sources ([MLHC 2024](#))

Outcome Definition



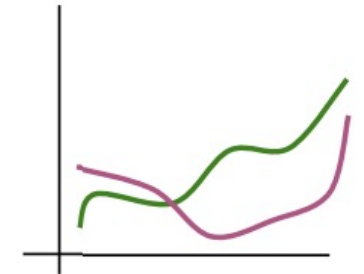
1. Rationales for treatment selection ([Work under review](#))
2. Access to care and how it affects EHR labels ([In progress](#))

Algorithm Development



1. Correcting for patient access to care ([AAAI 2022](#))
2. Developing new LLM bias frameworks ([In progress](#))

Post-Deployment Considerations



1. Bias auditing ([AMA Journal of Ethics 2019](#), [Nature Medicine 2021](#))
2. Mitigating algorithmic bias ([NeurIPS 2018](#))

How can we benefit from larger datasets?

1. We need a **data-oriented perspective**, especially considering **multiple sources**
2. We show theoretical and empirical results, including a **divergence-based method** for when to add more data

