# How Benchmark Prediction from Fewer Data Misses the Mark

Guanhua Zhang\* 1,2, Florian E. Dorner 1,2,3, Moritz Hardt 1,2 1 Max Planck Institute for Intelligent Systems, Tübingen 2 Tübingen AI Center 3 ETH Zurich

#### **Abstract**

Large language model (LLM) evaluation is increasingly costly, prompting interest in methods that speed up evaluation by shrinking benchmark datasets. Benchmark prediction (also called efficient LLM evaluation) aims to select a small subset of evaluation points and predict overall benchmark performance from that subset. In this paper, we systematically assess the strengths and limitations of 11 benchmark prediction methods across 19 diverse benchmarks. First, we identify a highly competitive baseline: Take a random sample and fit a regression model on the sample to predict missing entries. Outperforming most existing methods, this baseline challenges the assumption that careful subset selection is necessary for benchmark prediction. Second, we discover that all existing methods crucially depend on model similarity. They work best when interpolating scores among similar models. The effectiveness of benchmark prediction sharply declines when new models have higher accuracy than previously seen models. In this setting of extrapolation, none of the previous methods consistently beat a simple average over random samples. To improve over the sample average, we introduce a new method inspired by augmented inverse propensity weighting. This method consistently outperforms the random sample average even for extrapolation. However, its performance still relies on model similarity and the gains are modest in general. This shows that benchmark prediction fails just when it is most needed: at the evaluation frontier, where the goal is to evaluate new models of unknown capabilities<sup>†</sup>.

# 1 Introduction

Increasingly, computational cost is a major bottleneck in the evaluation of recent generative models. Growing model size and benchmark task difficulty, as well as the sheer number of available benchmarks all escalate the problem. For example, evaluating a single 176B parameter model on the HELM multi-task benchmark required 4,200 GPU hours [35]; even major companies noted the significant computational burden of evaluation on the BigBench multi-task benchmark [17].

The problem has prompted much recent work on more efficient LLM evaluation. The typical approach is to find a subset of data points to evaluate on, and to predict benchmark performance from these few evaluations. The simplest method is the *random sample mean*: Take a random sample of n evaluation points, and compute the mean of the benchmark metric on the sample. For a metric, like accuracy, with values in the interval [0,1], the sample mean gives an additive approximation up to error  $O(1/\sqrt{n})$ . More sophisticated methods try to improve over this baseline by following a common strategy: cleverly choose a small *core set* of evaluation points, evaluate multiple known models on these points, then fit a model to predict overall benchmark performance from these evaluations.

<sup>\*</sup>Corresponding author: guanhua.zhang@tuebingen.mpg.de

<sup>†</sup>Code is available at https://github.com/socialfoundations/benchmark-prediction.

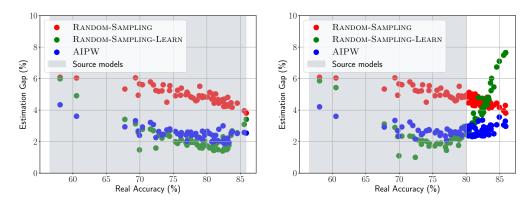


Figure 1: The x-axis denotes the real accuracy in ImageNet while the y-axis denotes the estimation gap (equation 1) for each target model. The gray stands for the accuracy range of source models. Left: source models are randomly sampled across all models. Right: source models are randomly sampled from models with accuracy lower than 80%.

We group existing efforts following this strategy under the term *benchmark prediction*. Previous research has proposed several hypotheses for why benchmark prediction can work: Core sets identify the most informative data points [49], they exploit the dependence between model performance on different data points [60], and they can capture the unobserved abilities of models [43].

The goal of our work is to systematically examine the strengths and weaknesses of benchmark prediction as a solution concept for efficient LLM evaluation.

## 1.1 Our Contributions

We conduct a large-scale, systematic evaluation of 11 state-of-the-art benchmark prediction methods across 19 diverse benchmarks. For each benchmark, we collect detailed performance results for at least 83 models on all data points. We split all models into two groups: source models and target models. For the source models, performance data is available for all data points. In contrast, for the target models, performance information is available only for up to 50 data points. Each method must adhere to this constraint of selecting no more than 50 data points, and the objective is to estimate each target model's mean performance across the full benchmark. To evaluate method effectiveness, we compute the average estimation gap—the absolute difference between the true and estimated full-benchmark performances across all target models.

Many methods work well on similar models, but a simple baseline works best. We first study the *interpolation* regime where the source and target models are random drawn from the same set of all models. Our empirical findings first confirm that in this regime it is possible to reduce the average estimation gap relative to RANDOM-SAMPLING, which simply reports the mean performance on a randomly selected core set. All evaluated methods outperform RANDOM-SAMPLING in over half of the benchmarks. Given that all methods operate on the same number of core-set evaluations, their computational costs are comparable. Thus, a lower *average estimation gap* indicates superior data efficiency—effectively, more informative use of the evaluation budget.

Surprisingly, what works best in the interpolation regime is remarkably simple: after randomly sampling a core set, rather than computing its mean, we fit a regression model to predict the true mean performance. This method, RANDOM-SAMPLING-LEARN, consistently outperforms most other methods and reduces the *average estimation gap* by an average of 37% compared to RANDOM-SAMPLING. This suggests that the manner of core-set selection is relatively unimportant; rather, the key to success is modeling the correlation between core-set and full-benchmark performances.

**Methods fail at the evaluation frontier.** However, our analysis reveals a major limitation: the effectiveness of benchmark prediction methods drops sharply when source and target models are *not* drawn from the same distribution. We call this the *extrapolation* regime. To explore this regime, we conduct an experiment in which we select the top-performing models (according to the full

benchmark) as target models and use the poorer-performing ones as source models. This setup reflects the typical use of benchmarking at the *evaluation frontier* where new models are being released that are likely better than existing ones.

In the extrapolation regime, we show that most benchmark predictions methods fail to outperform the naive RANDOM-SAMPLING baseline. This is illustrated in Figure 1. When source models cover the full range of performances (left), RANDOM-SAMPLING-LEARN more than halves the estimation error. When source models are restricted to the lower-than-80% accuracy (right), RANDOM-SAMPLING-LEARN still outperforms RANDOM-SAMPLING for target models similar to the source distribution, but its predictions substantially degrade for better-performing targets outside the source range.

AIPW is an overlooked exception to the rule. One notable exception is a method inspired by augmented inverse propensity weighting (AIPW)—used in other statistical applications—that we introduce in the context of benchmark prediction. AIPW reliably outperforms RANDOM-SAMPLING both under interpolation and extrapolation. Although it sometimes performs worse than RANDOM-SAMPLING-LEARN when targets resemble sources, it consistently maintains an advantage when they do not, thanks to being a consistent estimator. However, as illustrated in Figure 1 (right), even AIPW sees diminishing improvements as target models' accuracies exceed those of the sources.

Benchmark prediction relies on model similarity. To more systematically examine the generalization of benchmark prediction methods, we calculate the *model similarity* [38], quantifying how closely the predictions of each target model match those of the source models used in training. We observe a strong negative correlation between model similarity and estimation gap: methods that beat RANDOM-SAMPLING tend to do so primarily for targets similar to sources, while accuracy on disimilar models deteriorates. In contrast, RANDOM-SAMPLING exhibits neutral correlation, providing consistent (albeit less accurate) estimates regardless of similarity.

**Main takeaway.** Our findings suggest that while benchmark prediction techniques can be useful in specific scenarios, their reliance on similarity between source and target models poses a risk of misestimating the performance of new models. This underscores the importance of applying these methods with caution, especially for evaluating models that significantly deviate from previous ones.

## 2 Related Work

Evaluating large language models (LLMs) has become increasingly costly as these models grow in size and capabilities [35, 16, 65, 68]. These costs manifest in several ways. First, the collection and annotation of evaluation data can require significant resources [66]. To mitigate these costs, researchers have turned to methods such as using LLMs-as-judges [21, 23] or employing active labeling [32, 31, 10, 12, 70] to generate evaluation data and labels. However, these savings come with drawbacks. For instance, LLM-as-a-judge does not produce reliable evaluation outcomes, as judge models tend to prefer models similar to them, and have other biases [63, 42, 14, 7].

Another significant cost in LLM benchmarking arises from the model inference itself. Generating responses with LLMs can be time-consuming [35, 68, 53], and common inference time scaling techniques [57, 24, 54, 33] may exacerbate this issue. The success of scaling laws [29, 50] in predicting model performance has fueled interest in the development of benchmark prediction techniques [60, 43, 44, 40, 41], which aim to estimate benchmark performance by evaluating LLMs on a limited set of data <sup>2</sup>.

The key idea underpinning benchmark prediction is that not all evaluation examples carry the same amount of information [49]. It is hypothesized that a smaller core set of examples can represent the entire test set, allowing for accurate estimation of overall benchmark performance [60]. This is similar to efficient model training approaches, which aim to identify a subset of training data that enable performance comparable to training on the full dataset [52, 69]. Indeed, a popular benchmark prediction method, k-medoids clustering, is a classical approach to core-set selection for training [15]. However, it is important to recognize that the objectives of training and evaluation differ significantly. While training focuses on minimizing empirical risk and enhancing model performance,

<sup>&</sup>lt;sup>2</sup>Unlike bandit literature [68, 53], which focuses on identifying the best model from a pool, benchmark prediction is more challenging as it seeks to forecast overall benchmark performance for any new model.

evaluation seeks to provide an unbiased estimation of a model's performance to facilitate fair model comparison [40]. Our work challenges the assumption that core-set selection is the key to the success of benchmark prediction by introducing competitive methods that do not rely on core-set selection.

Many existing approaches treat benchmark prediction as a learning problem, aiming to predict a model's overall performance based on its performance on a subset of data [60, 43, 34, 30, 45]. Despite promising results, previous work has highlighted limitations in terms of estimation variance [36]. Going further, we highlight that most benchmark prediction methods rely on model similarity, with estimation performance deteriorating when target models deviate from familiar source models.

## 3 What is Benchmark Prediction?

#### 3.1 Problem Formulation

We define a benchmark as a triplet  $(\mathcal{D}, \mathcal{F}, s)$ . Here  $\mathcal{F}$  refers to the set of models to be evaluated on the benchmark and s represents the evaluation metric. Lastly,  $\mathcal{D}$  represents the benchmark data with  $|\mathcal{D}| = N$  data points. A data point is referred to as  $z \in \mathcal{D}$ , where z = (x, y), x refers to the query and y refers to the ground truth answer.

- s(f,z) refers to the performance of any  $f \in \mathcal{F}$  on any data point  $z \in \mathcal{D}$ . For example,  $s(f,z) = \mathbb{1}[f(x) = y]$  if the benchmark uses standard accuracy as the metric.
- $\bar{s}(f,\mathcal{D}') = \frac{1}{|\mathcal{D}'|} \sum_{z \in \mathcal{D}'} s(f,z)$  represents the average performance of  $f \in \mathcal{F}$  on any  $\mathcal{D}' \subset \mathcal{D}$ .
- $s(f, \mathcal{D}') = \{s(f, z)\}_{z \in \mathcal{D}'}$  represents the vectorized performance of  $f \in \mathcal{F}$  on all data points in  $\mathcal{D}' \subset \mathcal{D}$ , and  $s(\mathcal{F}', z) = \{s(f, z)\}_{f \in \mathcal{F}'}$  represents the vectorized performances of all models in  $\mathcal{F}' \subset \mathcal{F}$  on data point  $z \in \mathcal{D}$ .
- $S(\mathcal{F}', \mathcal{D}') = \{s(f, \mathcal{D}')\}_{f \in \mathcal{F}'} = \{s(\mathcal{F}', z)\}_{z \in \mathcal{D}'}^{\mathsf{T}}$  is the performance matrix of all models in  $\mathcal{F}' \subset \mathcal{F}$  on all data points in  $\mathcal{D}' \subset \mathcal{D}$ .

We refer to  $\mathcal{F}^{(s)} = \{f_1, \dots, f_M\} \subset \mathcal{F}$  as the set of source models, whose performances on every data point of the benchmark  $S(\mathcal{F}^{(s)}, \mathcal{D})$  are known. The rest of the models are the target models  $\mathcal{F}^{(t)} = \mathcal{F} \setminus \mathcal{F}^{(s)}$ , which are only be evaluated on  $n \ll N$  data points to save computational costs.

Benchmark prediction with fewer data aims to estimate  $\bar{s}(f,\mathcal{D})$  for every  $f \in \mathcal{F}^{(t)}$  with only n data points. In practice, benchmark prediction often involves two steps: ① identifying a representative core-set  $\mathcal{C} \subset \mathcal{D}$  with  $|\mathcal{C}| = n$  data points, and ② learning a performance estimator h to estimate the average performance on the full benchmark based on the core-set. Formally, the goal of benchmark prediction is to find  $\mathcal{C}$  and h to minimize the estimation gap over target models,

estimation gap: 
$$\frac{1}{|\mathcal{F}^{(t)}|} \sum_{f \in \mathcal{F}^{(t)}} \left| \bar{s}(f, \mathcal{D}) - h[s(f, \mathcal{C}), S(\mathcal{F}^{(s)}, \mathcal{D})] \right|. \tag{1}$$

For simplicity, in the remainder of the paper, we will denote the estimated performance of target model  $f \in \mathcal{F}^{(t)}$  as h(f), instead of explicitly writing  $h[s(f, \mathcal{C}), S(\mathcal{F}^{(s)}, \mathcal{D})]$ .

# 3.2 Benchmark Prediction Methods

**Previous methods.** In this paper, we examine five widely-used benchmark prediction methods,

- RANDOM-SAMPLING randomly samples a subset as  $\mathcal C$  and directly returns the mean performances as  $h^{\text{RANDOM-SAMPLING}}(f)$ .
- ANCHOR-POINTS-WEIGHTED [60] uses k-medoids clustering to identify  $\mathcal C$  and returns a weighted sum based on the density of each cluster as  $h^{\text{ANCHOR-POINTS-WEIGHTED}}(f)$ .
- ANCHOR-POINTS-PREDICTOR [60] extends ANCHOR-POINTS-WEIGHTED. Instead of directly returning the weighted sum, a linear regression model is learned as  $h^{\text{ANCHOR-POINTS-PREDICTOR}}(f)$ .
- P-IRT [43] extends ANCHOR-POINTS-PREDICTOR by replacing the regression model with an Item Response Theory (IRT) model as  $h^{\text{P-IRT}}(f)$ .
- GP-IRT [43] further generalizes P-IRT by combining its estimation with ANCHOR-POINTS-WEIGHTED as a weighted sum, and use it as  $h^{\text{GP-IRT}}(f)$ .

**New methods.** We introduce six methods that have not yet been applied to benchmark prediction.

- RANDOM-SAMPLING-LEARN randomly samples a subset as  $\mathcal C$  and learns a Ridge regression model g, which predict  $\bar s(f,\mathcal D)$  based on  $s(f,\mathcal C)$ , as  $h^{\text{RANDOM-SAMPLING-LEARN}}(f)$ .
- RANDOM-SEARCH-LEARN performs RANDOM-SAMPLING-LEARN for 10,000 times and selects the run based on cross-validation.
- LASSO trains a Lasso regression model to predict  $\bar{s}(f,\mathcal{D})$  based on  $s(f,\mathcal{D})$  with sparsity constraints on number of non-zero weights lower than n. The learned model is then used as  $h^{\text{LASSO}}(f)$ .
- DOUBLE-OPTIMIZE employs gradient descent to optimize both a subset selection vector, which
  models C, and a linear regression model, h<sup>DOUBLE-OPTIMIZE</sup> [27, 3].
- Principal Component Analysis (PCA) treats benchmark prediction as a matrix completion problem by assuming the performance matrix  $S(\mathcal{F}, \mathcal{D})$  is of low rank. By randomly sampling a subset as  $\mathcal{C}$ , this methods conducts PCA to impute the missing values for target models [59, 6].
- Augmented inverse propensity weighting (AIPW) [48]: Inspired by the application of prediction powered inference [2, 1] to the LLM-as-a-judge setting [5, 14], we apply a more general AIPW estimator to benchmark prediction. We train a Ridge regression model g for every target model f, which predicts the point-wise performance s(f, z) based on  $s(\mathcal{F}^{(s)}, z)$ . Formally,

$$g = \underset{g'}{\operatorname{arg\,min}} \frac{1}{n} \sum_{z \in \mathcal{C}} \left[ g'[\mathbf{s}(\mathcal{F}^{(s)}, z)] - s(f, z) \right]^2. \tag{2}$$

The idea behind the AIPW estimator is to use the predicted performance  $\hat{s}(f,z) = g[s(\mathcal{F}^{(s)},z)]$  as a proxy score to estimate  $\bar{s}(f,\mathcal{D})$  and "debias" that estimator as follows

$$h^{\text{AIPW}}(f) = \bar{s}(f, \mathcal{C}) + \frac{1}{1 + \frac{n}{N-n}} \left( \frac{1}{N-n} \sum_{z \in \mathcal{D} - \mathcal{C}} \hat{s}(f, z) - \frac{1}{n} \sum_{z \in \mathcal{C}} \hat{s}(f, z) \right). \tag{3}$$

Unlike the other learning-based baselines, AIPW is a consistent estimator for  $\bar{s}(f,\mathcal{D})$ [19]. Compared to RANDOM-SAMPLING, it reduces estimator variance by a factor of up to  $\frac{1}{1+\frac{n}{N}}\rho(\hat{s}(f,z),s(f,z))^2$  [14], where  $\rho$  is the Pearson correlation coefficient.

More details of each method are listed in Appendix A.

## 4 Experiments

# 4.1 Experiment Setup

We select a diverse range of benchmarks from the following sources<sup>3</sup>.

- HELM-Lite benchmarks [35]: OpenbookQA [39], GSM8K [9], LegalBench [22], Math [26], MedQA [28], and MMLU [25]. We obtain the per-data point performances of  $|\mathcal{F}|=83$  models from the official leaderboard.
- GLUE benchmarks [61]: MRPC [13], RTE [11, 18, 4], SST-2 [55], MNLI [64], and QNLI [46]. We use the per-data performances of  $|\mathcal{F}| = 87$  models provided by AnchorPoint<sup>4</sup> [60].
- OpenLLM benchmarks [16]: IFEval [67], Math [26], MMLU-Pro [62], Arc-Challenge [8], BBH [58], GPQA [47] and MUSR [56]. We use  $|\mathcal{F}|=448$  models provided by Huggingface <sup>5</sup> and collect their performance scores.
- ImageNet [51]: We collect  $|\mathcal{F}| = 110$  models from Pytorch Hub <sup>6</sup> and evaluate them on ImageNet.

A summary of benchmark statistics is provided in Appendix B.

<sup>&</sup>lt;sup>3</sup>Since P-IRT and GP-IRT requires s(f,z) to be binary, we only use benchmarks with accuracy as metric.

<sup>&</sup>lt;sup>4</sup>The provided score file for QQP is broken so we exclude it.

 $<sup>^5</sup> https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard\#$ 

 $<sup>^6</sup>$ https://pytorch.org/vision/stable/models.html#classification

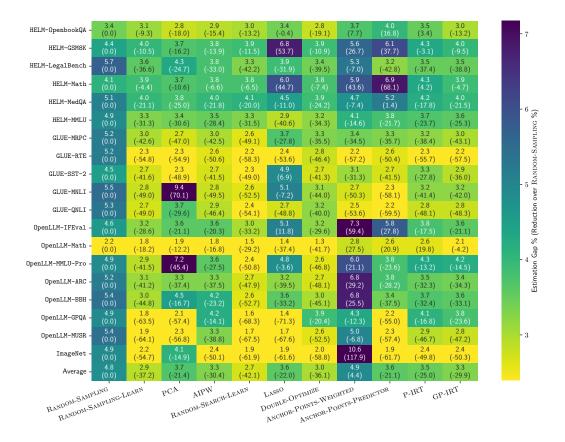


Figure 2: The estimation gaps  $(\downarrow)$  for target models (equation 1) under the interpolation split, where source and target models are identically distributed. Each target is evaluated on n=50 data points. The estimation gap reduction  $(\downarrow)$  over RANDOM-SAMPLING is shown in parentheses. A negative reduction means that the method achieves a lower gap than RANDOM-SAMPLING. % is omitted.

#### 4.2 Estimation Gap Reduction under Interpolation

As done in previous work [60, 43], we examine the effectiveness of benchmark prediction methods under the interpolation model split where source models are identically distributed with target models.

Interpolation model split. For each benchmark, we randomly select 75% of models as source models  $\mathcal{F}^{(s)}$ , for which performance scores across all data points  $S(\mathcal{F}^{(s)},\mathcal{D})$  are available. The remaining 25% of models serve as target models  $\mathcal{F}^{(t)}$  for assessment of benchmark prediction methods. Each target model is evaluated on only n=50 data points unless specified otherwise. Benchmark prediction methods are used to estimate the full benchmark average performance  $\bar{s}(f,\mathcal{D})$  of each target model  $f\in\mathcal{F}^{(t)}$  and evaluated based on the estimation gap from equation 1. Each experiment is repeated over 100 random trials, and we report the average estimation gap across all target models in these trials to ensure robustness. See standard errors in Appendix C.

**Results.** The results are presented in Figure 2. Compared to RANDOM-SAMPLING, all other benchmark methods effectively reduce the estimation gap in over half of the evaluated benchmarks. Notably, nine out of ten methods reduce the estimation gap by more than 20% on average across all benchmarks, as indicated in the last row. This verifies the effectiveness of benchmark prediction methods in the interpolation setting, where source and target models are identically distributed. Interestingly, the top-performing method is the simple baseline, RANDOM-SEARCH-LEARN, which achieves a 42.1% reduction compared to RANDOM-SAMPLING averaged accross all benchmarks. In comparison to the previous state-of-the-art, GP-IRT, which leads to a 29.9% reduction on average, RANDOM-SEARCH-LEARN results in a lower estimation gap in nearly all benchmarks.

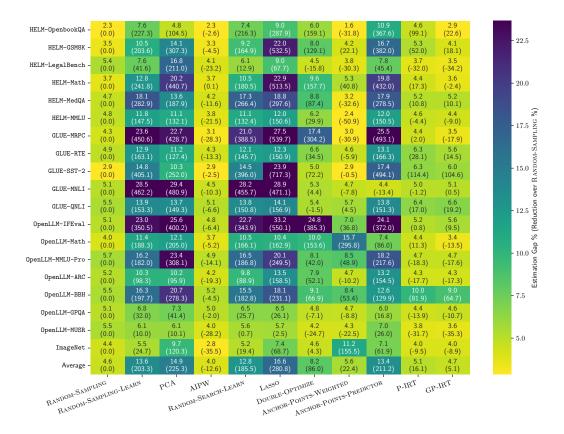


Figure 3: The estimation gaps ( $\downarrow$ ) for target models (equation 1) under extrapolation split, where source models are the lowest-performing 50%, and target models are the top 30%. Each target model is evaluated on n=50 data points. We also report the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted.

On the other hand, the selection of the core-set does not significantly enhance the effectiveness of benchmark prediction. For example, the second best-performing method, RANDOM-SAMPLING-LEARN, also consistently outperforms RANDOM-SAMPLING across all benchmarks, despite the sole difference being the use of a Ridge regression model rather than directly averaging across the core-set. With a 37.2% reduction in estimation gap, it performs comparably to RANDOM-SEARCH-LEARN, despite the latter conducting 10,000 iterations of RANDOM-SAMPLING-LEARN to identify the best subset. Moreover, it surpasses methods like DOUBLE-OPTIMIZE and GP-IRT, which select subsets through optimization or clustering. Another benchmark prediction method, AIPW, which also utilizes a randomly sampled core-set, consistently achieves a lower estimation gap across all benchmarks, yielding results comparable to the state-of-the-art GP-IRT. These findings challenge the prevailing notion that the core of benchmark prediction lies in identifying the most informative or representative subset. Instead, our results suggest that the primary driver of benchmark prediction success is learning to predict the mean, with core-set selection playing a relatively minor role.

# 4.3 Estimation Gap Increase under Extrapolation

We examine the effectiveness of benchmark prediction methods under the extrapolation model split where target models all perform better than source models.

**Extrapolation model split.** Different from the random source-target model split last subsection, we begin by ranking all models for a given benchmark based on their average performance on the full benchmark  $\bar{s}(f,\mathcal{D})$ . The lowest-performing 50% of these models are designated as source models, while the top 30% serve as target models for evaluating benchmark prediction methods. This strategy

reflects real-world model development scenarios, where developers debug and assess improved models based on existing less effective models. The estimation gap as defined in equation 1 is used for measuring the effectiveness of benchmark prediction. We again repeat each experiment 100 times.

**Results.** The results are shown in Figure 3. The average estimation gap for RANDOM-SAMPLING (4.6%) is largely comparable with the interpolation setting (4.8%) as it doesn't rely on source models. However, for all other methods, the estimation gap increases when compared to the interpolation setting. Nearly all benchmark prediction methods that outperform RANDOM-SAMPLING in the interpolation scenario now show diminished performance. Notably, the previous best method RANDOM-SEARCH-LEARN now results in a 185.1% increase in estimation gap than RANDOM-SAMPLING, and performs worse than RANDOM-SAMPLING across all benchmarks. The only method that still outperforms RANDOM-SAMPLING on average is AIPW, beating RANDOM-SAMPLING in 18 out of 19 benchmarks. This is because AIPW, like RANDOM-SAMPLING, is a consistent estimator, but has lower variance than random RANDOM-SAMPLING when its predictor is effective. However, the estimation gap reduction (-12.6%) of AIPW over RANDOM-SAMPLING in the extrapolation setting is also less pronounced than in the interpolation setting (-30.4%).

This stark contrast between interpolation and extrapolation settings underscores the heavy reliance of most benchmark prediction methods on the similarity between source and target models. This is unsurprising, given that many methods approach benchmark prediction as a machine learning problem, which often struggles in out-of-domain scenarios. However, unlike traditional machine learning, which primarily emphasizes in-domain performance, a key objective of benchmarking is to assess and identify new superior models. Therefore, extrapolation is a more prevalent and pertinent setting than interpolation in the context of benchmarking, and the decline in the estimation gap of benchmark prediction methods in this setting calls for more caution.

#### 4.4 Reliance on Model Similarity

In this subsection, we investigate the extent to which benchmark prediction methods rely on the similarity between target and source models.

**Model similarity.** We follow previous works [38, 20] and define the model similarity of target model f to all source models  $\mathcal{F}^{(s)}$  as follows,

$$\mathcal{S}(f, \mathcal{F}^{(s)}, \mathcal{D}) = \frac{1}{M} \sum_{f' \in \mathcal{F}^{(s)}} \frac{c_{obs} - c_{exp}}{1 - c_{exp}}.$$
 (4)

Here,  $c_{exp} = \bar{s}(f,\mathcal{D})\bar{s}(f',\mathcal{D}) + (1-\bar{s}(f,\mathcal{D}))(1-\bar{s}(f',\mathcal{D}))$  measures the chance agreement rate, i.e., the expected probability of  $\{s(f,z)=s(f',z)\}$  if s(f,z) is independent of s(f',z). In contrast,  $c_{obs} = \frac{1}{N}\sum_{z\in\mathcal{D}}\mathbbm{1}[s(f,z)=s(f',z)]$  is the observed agreement rate. For simplicity, we use  $\mathcal{S}(f)$  to denote  $\mathcal{S}(f,\mathcal{F}^{(s)},\mathcal{D})$  in the remainder of the paper.  $\mathcal{S}(f)$  quantifies how similar the performance pattern of the target model f is to all source models  $\mathcal{F}^{(s)}$ , with a higher value indicating greater similarity [20].

We aim to examine the correlation between model similarity and estimation gap. However, we note that the estimation depends on the standard deviation of s(f,z). Since we use accuracy as the metric in our experiment, s(f,z) is Bernoulli with parameter  $p_f = \bar{s}(f,\mathcal{D})$  and standard deviation  $\sigma_f = \sqrt{p_f(1-p_f)}$ . By randomly sampling n data points as  $\mathcal{C}$ , Chebyshev's inequality ensures that

$$|\bar{s}(f,\mathcal{C}) - \bar{s}(f,\mathcal{D})| < \sigma_f / \sqrt{\alpha n}$$
 (5)

with probability at least  $(1-\alpha)$ . In other words, the performance of target models with lower  $\sigma_f$  is easier to estimate with the same amount of data. Thus, the standard deviation of the basic estimation gap could potentially confound the observed correlation between model similarity and estimation gap. Consider the method RANDOM-SAMPLING, whose estimation does not depend on source models. If all target models with low  $\sigma_f$  coincidentally have high  $\mathcal{S}(f)$ , while those with high  $\sigma_f$  have low  $\mathcal{S}(f)$ , then a spurious correlation between estimation gap and model similarity to target models could appear even for RANDOM-SAMPLING. To prevent this, we define the normalized estimation gap as

normalized estimation gap for 
$$f$$
:  $\mathcal{E}(f) = \frac{1}{\sigma_f} |\bar{s}(f, \mathcal{D}) - h(f)|$ . (6)

Then we measure the Pearson correlation between model similarity in equation 4 and the normalized estimation gap in equation 6.

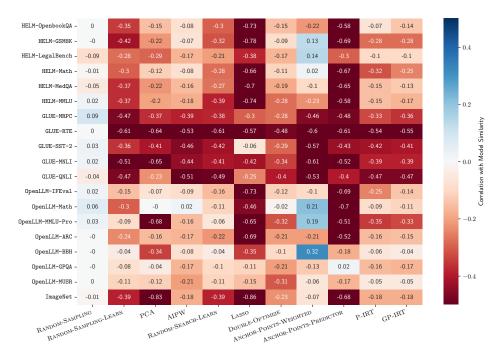


Figure 4: The Pearson correlation between normalized per-model estimation gap (equation 6) and model similarity (equation 4). Negative correlation indicates that target models that are dissimilar to source models tend to have larger estimation gap, and vice versa.

**Results.** The results are shown in Figure 4. A clear negative correlation between model similarity and estimation gap emerges for almost all benchmark prediction methods except for RANDOM-SAMPLING. In particular, the best-performing method under the interpolation model split, RANDOM-SAMPLING-LEARN, exhibits a negative correlation below -0.2 in 13/19 benchmarks. Despite its asymptotic unbiasedness, we also find negative correlations for AIPW. This is perhaps unsurprising: While AIPW is consistent independent of how well its regression model  $g[s(\mathcal{F}^{(s)},z)]$  predicts s(f,z), its variance depends precisely on that prediction quality. If the predictions are good, AIPW improves substantially over RANDOM-SAMPLING, while there is no improvement when predictions are fully uninformative. But intuitively, predicting s(f,z) is harder when f is very different from the models  $\mathcal{F}^{(s)}$  used for training the predictor  $g[s(\mathcal{F}^{(s)},z)]$ .

#### 4.5 Ablation on Core-set Size

We conduct an ablation study on the size of the core-set n. We experiment with  $n \in \{10, 20, 50, 100, 200\}$ , and the summarized results are shown in Table 1 (detailed results can be found in Appendix C). As expected, the estimation gap generally decreases as n increases for most methods. Our previous conclusions remain valid across both settings. With larger core-set sizes, most methods continue to perform better than RANDOM-SAMPLING in the interpolation split but fail to do so in the extrapolation model split. Interestingly, we also find that RANDOM-SAMPLING outperforms all other methods when given twice as much data, even in the interpolation model split.

AIPW remains effective in both settings. However, its advantage over RANDOM-SAMPLING diminishes as n increases. While AIPW reduces the estimation gap by -30.4% in interpolation and -12.6% in extrapolation for n=50, these advantages shrink to -12.4% in interpolation and a mere -2.3% in extrapolation for n=200. This is because the estimator variance reduction factor of AIPW is up to  $\frac{1}{1+\frac{n}{N}}\rho(\hat{s}^(f,z),s(f,z))^2$ . On the other hand, the advantage of AIPW remains significant when the dataset is large and thus  $\frac{n}{N}$  is small. Figure 5 compares AIPW with n=50 to RANDOM-SAMPLING with n=100 data points using ImageNet. AIPW achieves a lower average normalized estimation gap compared to RANDOM-SAMPLING, despite using only half the data. However, the normalized estimation gap for AIPW is biased with respect to model similarity. In contrast, the

Table 1: Ablation study on the core-set size n. We report the estimation gap averaged over all benchmarks. % is neglected for each metric. The lowest estimation gap in each column is highlighted in bold. See detailed results in Appendix C.

	n = 10		Interpolati $n = 50$		n = 200	n = 10		Extrapolat $\mid n = 50$	$ \begin{array}{c c} \text{ion} \\                                    $	n = 200
RANDOM-SAMPLING	11.0	7.7	4.8	3.3	2.1	10.7	7.4	4.6	3.1	2.0
RANDOM-SAMPLING-LEARN	5.4	4.2	2.9	2.1	1.5	17.6	15.8	13.6	12.1	11.1
PCA	6.6	5.2	3.7	2.8	2.1	19.9	17.6	14.9	12.2	9.3
AIPW	8.3	5.4	3.3	2.3	1.8	9.6	6.5	4.0	2.8	2.0
RANDOM-SEARCH-LEARN	4.5	3.7	2.7	2.0	1.4	16.0	14.4	12.8	11.8	11.1
LASSO	7.8	6.1	3.6	2.6	2.2	22.0	19.3	16.6	15.3	14.6
DOUBLE-OPTIMIZE	6.6	4.8	3.0	2.3	1.9	11.3	9.0	8.2	8.0	7.0
ANCHOR-POINTS-WEIGHTED	8.9	6.9	4.9	4.0	3.2	10.4	6.7	5.6	4.7	3.4
ANCHOR-POINTS-PREDICTOR	4.7	4.1	3.6	3.4	4.1	16.2	14.8	13.4	12.4	11.2
P-IRT	7.3	5.9	3.5	2.1	1.3	9.8	8.2	5.1	3.7	3.0
GP-IRT	7.2	5.7	3.3	2.1	1.4	9.7	7.8	4.7	3.4	2.5

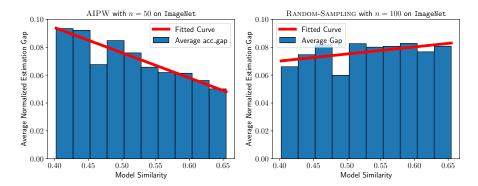


Figure 5: Average normalized estimation gap relative to model similarity for AIPW (n=50) and RANDOM-SAMPLING (n=100) on ImageNet. Each bar represents the target models whose similarity to source models falls within the corresponding range. The normalized estimation gap is defined as shown in equation 6. On average, AIPW outperforms RANDOM-SAMPLING, even with half the data. However, RANDOM-SAMPLING shows better performance when model similarity is low.

normalized estimation gap under RANDOM-SAMPLING remains largely neutral regarding model similarity. Consequently, while AIPW reduces the average, it produces a higher gap for models with low similarity compared to RANDOM-SAMPLING with twice the data.

# 5 Conclusion

In this paper, we study the problem of benchmark prediction from fewer data and examine 11 benchmark prediction methods. Our findings call into question the necessity of meticulous core-set selection and reveal that these methods are most proficient at interpolating scores among similar models. However, except RANDOM-SAMPLING and AIPW, all methods face significant difficulties when predicting target models that differ substantially from those they have encountered before.

We caution against the indiscriminate use of benchmark prediction techniques, as their dependence on model similarity causes most of them to fail precisely when most needed: at the evaluation frontier, where the aim is to assess new models with unknown capabilities. Even in the context of interpolation, no method outperforms RANDOM-SAMPLING, when that simple baseline is given access to twice as much data. Thus, while we recommend to use AIPW as a consistent estimator with lower variance, this suggests that simply raising the sampling budget for RANDOM-SAMPLING can be competitive, especially in settings where predictions of other models for fitting AIPW are costly to obtain.

**Acknowledgement.** We thank Yatong Chen and Jiduan Wu for helpful discussions and feedback on draft versions of this work. Florian Dorner is grateful for financial support from the Max Planck ETH Center for Learning Systems (CLS).

## References

- [1] Anastasios Nikolas Angelopoulos, Stephen Bates, Clara Fannjiang, Michael I. Jordan, and Tijana Zrnic. Prediction-powered inference. *Science*, 382:669 674, 2023.
- [2] Anastasios Nikolas Angelopoulos, John C. Duchi, and Tijana Zrnic. Ppi++: Efficient prediction-powered inference. *ArXiv*, abs/2311.01453, 2023.
- [3] Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *ArXiv*, abs/1308.3432, 2013.
- [4] Luisa Bentivogli, Ido Dagan, Hoa Trang Dang, Danilo Giampiccolo, and Bernardo Magnini. The fifth PASCAL recognizing textual entailment challenge. 2009.
- [5] Pierre Boyeau, Anastasios N Angelopoulos, Nir Yosef, Jitendra Malik, and Michael I Jordan. Autoeval done right: Using synthetic data for model evaluation. arXiv preprint arXiv:2403.07008, 2024.
- [6] Jian-Feng Cai, Emmanuel J Candès, and Zuowei Shen. A singular value thresholding algorithm for matrix completion. *SIAM Journal on optimization*, 20(4):1956–1982, 2010.
- [7] Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. Humans or llms as the judge? a study on judgement biases. *arXiv preprint arXiv:2402.10669*, 2024.
- [8] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*, 2018.
- [9] Karl Cobbe, Vineet Kosaraju, Mo Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *ArXiv*, abs/2110.14168, 2021.
- [10] Ciprian A Corneanu, Sergio Escalera, and Aleix M Martinez. Computing the testing error without a testing set. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2677–2685, 2020.
- [11] Ido Dagan, Oren Glickman, and Bernardo Magnini. The PASCAL recognising textual entailment challenge. In *Machine learning challenges. evaluating predictive uncertainty, visual object classification, and recognising tectual entailment*, pages 177–190. Springer, 2006.
- [12] Weijian Deng and Liang Zheng. Are labels always necessary for classifier accuracy evaluation? In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15069–15078, 2021.
- [13] William B Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the International Workshop on Paraphrasing*, 2005.
- [14] Florian E Dorner, Vivian Yvonne Nastl, and Moritz Hardt. Limits to scalable evaluation at the frontier: Llm as judge won't beat twice the data. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [15] Reza Zanjirani Farahani and Masoud Hekmatfar. Facility location: concepts, models, algorithms and case studies. 2009.
- [16] Clémentine Fourrier, Nathan Habib, Alina Lozovskaya, Konrad Szafer, and Thomas Wolf. Open llm leaderboard v2. https://huggingface.co/spaces/open-llm-leaderboard/open\_llm\_leaderboard, 2024.
- [17] Deep Ganguli, Nicholas Schiefer, Marina Favaro, and Jack Clark. Challenges in evaluating AI systems, 2023.
- [18] Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. The third PASCAL recognizing textual entailment challenge. In *Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing*, pages 1–9. Association for Computational Linguistics, 2007.

- [19] Adam N Glynn and Kevin M Quinn. An introduction to the augmented inverse propensity weighted estimator. *Political analysis*, 18(1):36–56, 2010.
- [20] Shashwat Goel, Joschka Struber, Ilze Amanda Auzina, Karuna K. Chandra, Ponnurangam Kumaraguru, Douwe Kiela, Ameya Prabhu, Matthias Bethge, and Jonas Geiping. Great models think alike and this undermines ai oversight. *ArXiv*, abs/2502.04313, 2025.
- [21] Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Yuanzhuo Wang, and Jian Guo. A survey on llm-as-a-judge. *ArXiv*, abs/2411.15594, 2024.
- [22] Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai Porat, Jason Hegland, Jessica Wu, Joe Nudell, Joel Niklaus, John Nay, Jonathan H. Choi, Kevin Tobia, Margaret Hagan, Megan Ma, Michael Livermore, Nikon Rasumov-Rahe, Nils Holzenberger, Noam Kolt, Peter Henderson, Sean Rehaag, Sharad Goel, Shang Gao, Spencer Williams, Sunny Gandhi, Tom Zur, Varun Iyer, and Zehua Li. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models, 2023.
- [23] Veronika Hackl, Alexandra Elena Müller, Michael Granitzer, and Maximilian Sailer. Is gpt-4 a reliable rater? evaluating consistency in gpt-4 text ratings. *ArXiv*, abs/2308.02575, 2023.
- [24] Moritz Hardt and Yu Sun. Test-time training on nearest neighbors for large language models. arXiv preprint arXiv:2305.18466, 2023.
- [25] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Xiaodong Song, and Jacob Steinhardt. Measuring massive multitask language understanding. ArXiv, abs/2009.03300, 2020.
- [26] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Xiaodong Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. ArXiv, abs/2103.03874, 2021.
- [27] Eric Jang, Shixiang Shane Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *ArXiv*, abs/1611.01144, 2016.
- [28] Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *ArXiv*, abs/2009.13081, 2020.
- [29] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. arXiv preprint arXiv:2001.08361, 2020.
- [30] Alex Kipnis, Konstantinos Voudouris, Luca M. Schulze Buschoff, and Eric Schulz. metabench a sparse benchmark of reasoning and knowledge in large language models. 2024.
- [31] Jannik Kossen, Sebastian Farquhar, Yarin Gal, and Tom Rainforth. Active testing: Sample-efficient model evaluation. In *International Conference on Machine Learning*, 2021.
- [32] Jannik Kossen, Sebastian Farquhar, Yarin Gal, and Tom Rainforth. Active surrogate estimators: An active learning approach to label-efficient model evaluation. *ArXiv*, abs/2202.06881, 2022.
- [33] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474, 2020.
- [34] Yang Li, Jie Ma, Miguel Ballesteros, Yassine Benajiba, and Graham Horwood. Active evaluation acquisition for efficient llm benchmarking. *ArXiv*, abs/2410.05952, 2024.

- [35] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher R'e, Diana Acosta-Navas, Drew A. Hudson, E. Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel J. Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan S. Kim, Neel Guha, Niladri S. Chatterji, O. Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas F. Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. Holistic evaluation of language models. *Annals of the New York Academy of Sciences*, 1525:140 146, 2023.
- [36] Lovish Madaan, Aaditya K. Singh, Rylan Schaeffer, Andrew Poulton, Oluwasanmi Koyejo, Pontus Stenetorp, Sharan Narang, and Dieuwke Hupkes. Quantifying variance in evaluation benchmarks. *ArXiv*, abs/2406.10229, 2024.
- [37] Pranav Mani, Peng Xu, Zachary Chase Lipton, and Michael Oberst. No free lunch: Non-asymptotic analysis of prediction-powered inference. *ArXiv*, abs/2505.20178, 2025.
- [38] Horia Mania, John Miller, Ludwig Schmidt, Moritz Hardt, and Benjamin Recht. Model similarity mitigates test set overuse. *ArXiv*, abs/1905.12580, 2019.
- [39] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391, Brussels, Belgium, October-November 2018. Association for Computational Linguistics.
- [40] David Owen. How predictable is language model benchmark performance? *ArXiv*, abs/2401.04757, 2024.
- [41] Lorenzo Pacchiardi, Konstantinos Voudouris, Ben Slater, Fernando Mart'inez-Plumed, Jos'e Hern'andez-Orallo, Lexin Zhou, and Wout Schellaert. Predictaboard: Benchmarking llm score predictability. *ArXiv*, abs/2502.14445, 2025.
- [42] Arjun Panickssery, Samuel R. Bowman, and Shi Feng. Llm evaluators recognize and favor their own generations. *ArXiv*, abs/2404.13076, 2024.
- [43] Felipe Maia Polo, Lucas Weber, Leshem Choshen, Yuekai Sun, Gongjun Xu, and Mikhail Yurochkin. tinybenchmarks: evaluating llms with fewer examples. *ArXiv*, abs/2402.14992, 2024.
- [44] Felipe Maia Polo, Ronald Xu, Lucas Weber, Mírian Silva, Onkar Bhardwaj, Leshem Choshen, Allysson Flavio Melo de Oliveira, Yuekai Sun, and Mikhail Yurochkin. Efficient multi-prompt evaluation of llms. *arXiv preprint arXiv:2405.17202*, 2024.
- [45] Ameya Prabhu, Vishaal Udandarao, Philip H. S. Torr, Matthias Bethge, Adel Bibi, and Samuel Albanie. Efficient lifelong model evaluation in an era of rapid progress. In *Neural Information Processing Systems*, 2024.
- [46] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of EMNLP*, pages 2383–2392. Association for Computational Linguistics, 2016.
- [47] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark. ArXiv, abs/2311.12022, 2023.
- [48] James M Robins and Andrea Rotnitzky. Semiparametric efficiency in multivariate regression models with missing data. *Journal of the American Statistical Association*, 90(429):122–129, 1995.
- [49] Pedro Rodriguez, Joe Barrow, Alexander Miserlis Hoyle, John P. Lalor, Robin Jia, and Jordan L. Boyd-Graber. Evaluation examples are not equally informative: How should that change nlp leaderboards? In *Annual Meeting of the Association for Computational Linguistics*, 2021.

- [50] Yangjun Ruan, Chris J. Maddison, and Tatsunori B. Hashimoto. Observational scaling laws and the predictability of language model performance. ArXiv, abs/2405.10938, 2024.
- [51] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115:211 252, 2014.
- [52] Noveen Sachdeva and Julian McAuley. Data distillation: A survey. Trans. Mach. Learn. Res., 2023, 2023.
- [53] Chengshuai Shi, Kun Yang, Jing Yang, and Cong Shen. Best arm identification for prompt learning under a limited budget. *arXiv preprint arXiv:2402.09723*, 2024.
- [54] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling Ilm test-time compute optimally can be more effective than scaling model parameters. arXiv preprint arXiv:2408.03314, 2024.
- [55] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of EMNLP*, pages 1631–1642, 2013.
- [56] Zayne Sprague, Xi Ye, Kaj Bostrom, Swarat Chaudhuri, and Greg Durrett. Musr: Testing the limits of chain-of-thought with multistep soft reasoning. *ArXiv*, abs/2310.16049, 2023.
- [57] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning*, pages 9229–9248. PMLR, 2020.
- [58] Mirac Suzgun, Nathan Scales, Nathanael Scharli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V. Le, Ed H. Chi, Denny Zhou, and Jason Wei. Challenging big-bench tasks and whether chain-of-thought can solve them. In *Annual Meeting of the Association for Computational Linguistics*, 2022.
- [59] Roman Vershynin. Four lectures on probabilistic methods for data science. ArXiv, abs/1612.06661, 2016.
- [60] Rajan Vivek, Kawin Ethayarajh, Diyi Yang, and Douwe Kiela. Anchor points: Benchmarking models with much fewer examples. ArXiv, abs/2309.08638, 2023.
- [61] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In BlackboxNLP@EMNLP, 2018.
- [62] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max W.F. Ku, Kai Wang, Alex Zhuang, Rongqi "Richard" Fan, Xiang Yue, and Wenhu Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *ArXiv*, abs/2406.01574, 2024.
- [63] Koki Wataoka, Tsubasa Takahashi, and Ryokan Ri. Self-preference bias in llm-as-a-judge. *ArXiv*, abs/2410.21819, 2024.
- [64] Adina Williams, Nikita Nangia, and Samuel R. Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of NAACL-HLT*, 2018.
- [65] Dingli Yu, Simran Kaur, Arushi Gupta, Jonah Brown-Cohen, Anirudh Goyal, and Sanjeev Arora. Skill-mix: a flexible and expandable family of evaluations for ai models. ArXiv, abs/2310.17567, 2023.
- [66] Xiang Yue, Boshi Wang, Kai Zhang, Ziru Chen, Yu Su, and Huan Sun. Automatic evaluation of attribution by large language models. In *Conference on Empirical Methods in Natural Language Processing*, 2023.

- [67] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models. *ArXiv*, abs/2311.07911, 2023.
- [68] Jin Peng Zhou, Christian K. Belardi, Ruihan Wu, Travis Zhang, Carla P. Gomes, Wen Sun, and Kilian Q. Weinberger. On speeding up language model evaluation. *ArXiv*, abs/2407.06172, 2024.
- [69] Xiao Zhou, Renjie Pi, Weizhong Zhang, Yong Lin, Zonghao Chen, and T. Zhang. Probabilistic bilevel coreset selection. *ArXiv*, abs/2301.09880, 2023.
- [70] Vilém Zouhar, Peng Cui, and Mrinmaya Sachan. How to select datapoints for efficient human evaluation of nlg models?, 2025.

# **NeurIPS Paper Checklist**

## 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: See Section 1.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: See Appendix D.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

## 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: This paper is mainly an empirical work and doesn't provide many new theoretical results.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: See Appendix B and the code.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We release our codes in the supplemental materials.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

# 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See Appendix B.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: See Appendix C.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

#### 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: See Appendix C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Authors have reviewed the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

# 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Guidelines:

- Justification: See Appendix D.
  - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

• The answer NA means that there is no societal impact of the work performed.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.

- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This paper doesn't release any new data or model.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All used models and datasets are well cited in Section 4.

#### Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
  package should be provided. For popular datasets, paperswithcode.com/datasets
  has curated licenses for some datasets. Their licensing guide can help determine the
  license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This paper doesn't provide new assets.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

#### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper doesn't involve crowd-sourcing experiments.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper doesn't involve crowd-sourcing experiments.

# Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

#### 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs. Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

# A Details of Benchmark Prediction Methods

#### A.1 Problem Formulation

We repeat the notation and the problem formulation here for the reader's convenience.

- A benchmark is represented as a triplet  $(\mathcal{D}, \mathcal{F}, s)$ .
- $\mathcal{D}$  represents the benchmark data with  $|\mathcal{D}| = N$  data points. A data point is referred to as  $z \in \mathcal{D}$ , where z = (x, y), x refers to the query and y refers to the ground truth answer.
- $\mathcal{F}$  refers to all potential models that can be evaluated on the benchmark.
- s represents the metric of the benchmark.
  - s(f,z) refers to the performance of any  $f \in \mathcal{F}$  on any data point  $z \in \mathcal{D}$ . For example,  $s(f,z) = \mathbb{1}[f(x) = y]$  if the benchmark uses standard accuracy as the metric.
  - $-\bar{s}(f,\mathcal{D}') = \frac{1}{|\mathcal{D}'|} \sum_{z \in \mathcal{D}'} s(f,z)$  represents the average performance of  $f \in \mathcal{F}$  on any  $\mathcal{D}' \subset \mathcal{D}$ .
  - $s(f, \mathcal{D}') = \{s(f, z)\}_{z \in \mathcal{D}'}$  represents the vectorized performance of  $f \in \mathcal{F}$  on all data points in  $\mathcal{D}' \subset \mathcal{D}$ , and  $s(\mathcal{F}', z) = \{s(f, z)\}_{f \in \mathcal{F}'}$  represents the vectorized performances of all models in  $\mathcal{F}' \subset \mathcal{F}$  on data point  $z \in \mathcal{D}$ .
  - $S(\mathcal{F}', \mathcal{D}') = \{s(f, \mathcal{D}')\}_{f \in \mathcal{F}'} = \{s(\mathcal{F}', z)\}_{z \in \mathcal{D}'}^{\mathsf{T}}$  as the performance matrix of all models in  $\mathcal{F}' \subset \mathcal{F}$  on all data points in  $\mathcal{D}' \subset \mathcal{D}$ .
- $\mathcal{F}^{(s)} = \{f_1, \dots, f_M\} \subset \mathcal{F}$  refers to a set of source models, whose performances on every data point of the benchmark  $S(\mathcal{F}^{(s)}, \mathcal{D})$  are known.
- The rest of the models are referred to as target models  $\mathcal{F}^{(t)} = \mathcal{F} \setminus \mathcal{F}^{(s)}$ , which can only be evaluated on at most  $n \ll N$  data points to save computational costs.

Benchmark prediction with fewer data aims to estimate  $\bar{s}(f,\mathcal{D})$  for every  $f \in \mathcal{F}^{(t)}$  with only n data points. In practice, benchmark prediction often involves two steps: ① identifying a representative core-set  $\mathcal{C} \subset \mathcal{D}$  with  $|\mathcal{C}| = n$  data points, and ② learning a performance estimator h to estimate the average performance on the full benchmark based on the core-set. Formally, the goal of benchmark prediction is to find  $\mathcal{C}$  and h to minimize the estimation gap over target models,

estimation gap: 
$$\frac{1}{|\mathcal{F}^{(t)}|} \sum_{f \in \mathcal{F}^{(t)}} \left| \bar{s}(f, \mathcal{D}) - h[s(f, \mathcal{C}), S(\mathcal{F}^{(s)}, \mathcal{D})] \right|. \tag{7}$$

For simplicity, in the remainder of the paper, we will denote the estimated performance of target model  $f \in \mathcal{F}^{(t)}$  as h(f), instead of explicitly writing  $h[s(f, \mathcal{C}), S(\mathcal{F}^{(s)}, \mathcal{D})]$ .

# A.2 Benchmark Prediction Methods

**Previous methods** In this paper, we examine five widely-used benchmark prediction methods,

• RANDOM-SAMPLING randomly samples a subset as C and directly returns the mean performance,

$$h^{\text{RANDOM-SAMPLING}}(f) = \bar{s}(f, \mathcal{C}).$$
 (8)

If the benchmark metric s is standard accuracy, the gap  $|\bar{s}(f,\mathcal{C}) - \bar{s}(f,\mathcal{D})|$  is bounded by  $\mathcal{O}(\sqrt{1/n})$  with high probability based on Hoeffding's inequality.

• ANCHOR-POINTS-WEIGHTED [60] treats benchmark prediction as a k-medoids clustering problem. The selected medoids are used as C, and a weight vector  $\theta \in \mathbb{R}^n$  is calculated as the normalized cluster size of each medoid. The final estimate for any target model  $f \in \mathcal{F}^{(t)}$  is

$$h^{\text{ANCHOR-POINTS-WEIGHTED}}(f) = \mathbf{s}(f, \mathcal{C})^T \boldsymbol{\theta}.$$
 (9)

• ANCHOR-POINTS-PREDICTOR [60] extends ANCHOR-POINTS-WEIGHTED. Instead of directly returning the weighted sum, a linear regression model g[s(f, C)] is learned to predict s(f, D - C).

$$h^{\text{Anchor-Points-Predictor}}(f) = \bar{g}[s(f, \mathcal{C})]$$
 (10)

where 
$$g = \underset{g'}{\operatorname{arg\,min}} \frac{1}{M} \sum_{f \in \mathcal{F}^{(s)}} \| s(f, \mathcal{D} - \mathcal{C}) - g'[s(f, \mathcal{C})] \|_2^2$$
, (11)

where we note that g[s(f, C)] is a (N - n) dimensional vector and we use  $\bar{g}[s(f, C)]$  as its mean.

• P-IRT [43] extends ANCHOR-POINTS-PREDICTOR by replacing the regression model g in equation 11 with an Item Response Theory (IRT) model. Following the notation for ANCHOR-POINTS-PREDICTOR, we estimate performance for any  $f \in \mathcal{F}^{(t)}$  as follows:

$$h^{\text{P-IRT}}(f) = \frac{N-n}{N} \bar{g}[s(f,\mathcal{C})] + \frac{n}{N} \bar{s}(f,\mathcal{C}). \tag{12}$$

GP-IRT [43] further generalizes P-IRT by combining its estimation with ANCHOR-POINTS-WEIGHTED as a weighted sum,

$$h^{\text{GP-IRT}}(f) = \lambda h^{\text{ANCHOR-POINTS-WEIGHTED}}(f) + (1 - \lambda)h^{\text{P-IRT}}(f), \tag{13}$$

where  $\lambda$  is chosen heuristically to control the error of P-IRT.

**New methods** We introduce six methods that have not yet been applied to benchmark prediction.

• RANDOM-SAMPLING-LEARN randomly samples a subset as C and adopts a Ridge regression model g for estimation as follows,

$$h^{\text{Random-Sampling-Learn}}(f) = g[s(f, C)]$$
 (14)

where 
$$g = \underset{g'}{\operatorname{arg\,min}} \frac{1}{M} \sum_{f \in \mathcal{F}^{(s)}} \left| \bar{s}(f, \mathcal{D}) - g'[s(f, \mathcal{C})] \right|.$$
 (15)

- RANDOM-SEARCH-LEARN performs RANDOM-SAMPLING-LEARN for 10,000 times and selects the best-performing subset as  $\mathcal C$  based on cross-validation. A Ridge regression model g is then trained and used in the same way as RANDOM-SELECTION-LEARN.
- LASSO trains a Lasso regression model with weights  $\theta \in \mathbb{R}^N$  as follows,

$$h^{\text{LASSO}}(f) = s(f, \mathcal{C})^T \boldsymbol{\theta}_{\mathcal{C}}$$
(16)

where 
$$\boldsymbol{\theta} = \underset{\boldsymbol{\theta}'}{\operatorname{arg min}} \frac{1}{n} \sum_{\boldsymbol{r} \in \mathcal{C}} \left[ \boldsymbol{s}(f, \mathcal{D})^{\mathsf{T}} \boldsymbol{\theta}' - \bar{\boldsymbol{s}}(f, \mathcal{D}) \right]^2 + \lambda \|\boldsymbol{\theta}'\|_1,$$
 (17)

where  $\lambda$  is selected so that only n dimensions of  $\theta$  are non-zero and  $\theta_{\mathcal{C}}$  is the non-zero slice of  $\theta$ .

• DOUBLE-OPTIMIZE optimizes both a subset selection vector  $\pi \in \mathbb{R}^N$  and a linear regression model with weights  $\theta \in \mathbb{R}^N$  with gradient descent as follows,

$$h^{\text{DOUBLE-OPTIMIZE}}(f) = [\boldsymbol{s}(f,\mathcal{D}) \cdot \text{TopMask}(\boldsymbol{\pi};n)]^{\text{T}}\boldsymbol{\theta} \tag{18}$$

where 
$$\boldsymbol{\pi}, \boldsymbol{\theta} = \underset{\boldsymbol{\pi}', \boldsymbol{\theta}'}{\operatorname{arg min}} \left\{ [\boldsymbol{s}(f, \mathcal{D}) \cdot \operatorname{TopMask}(\boldsymbol{\pi}'; n)]^{\mathsf{T}} \boldsymbol{\theta}' - \bar{\boldsymbol{s}}(f, \mathcal{D}) \right\}^{2},$$
 (19)

where  $\cdot$  refers to the bitwise multiplication between two vectors, and TopMask $(\pi';n)$  replaces the top n largest values of  $\pi'$  with 1s and the rest with 0s. We directly pass the gradient on TopMask $(\pi';n)$  to  $\pi'$  during optimization following the Straight-Through technique [27, 3].

- Principal Component Analysis (PCA) treats benchmark prediction as a matrix completion problem. This method assumes the performance matrix  $S(\mathcal{F},\mathcal{D})$  is of low rank. By randomly sampling a subset as  $\mathcal{C}$ , this methods conducts PCA to impute the missing values for target models [59, 6]. As a more intuitive view, one could also take the acquired principal components as model capability indicators [50], i.e., the  $(M \times k)$  PCA-transformed scores indicate the k-capabilities of each model, while the  $(k \times N)$  principal components represent the capability requirements for each data point. We select k among  $\{2,5,10,20\}$  through cross-validation. The Pseudo codes are in Algorithm 1.
- Augmented inverse propensity weighting (AIPW) [48]: Inspired by the application of prediction powered inference [2, 1] to the LLM-as-a-judge setting [5, 14], we apply a more general AIPW estimator to benchmark prediction. We train a Ridge regression model g for every target model f, which predicts the point-wise performance s(f, z) based on  $s(\mathcal{F}^{(s)}, z)$ . Formally,

$$g = \arg\min_{g'} \frac{1}{n} \sum_{z \in \mathcal{C}} \left[ g'[s(\mathcal{F}^{(s)}, z)] - s(f, z) \right]^2.$$
 (20)

The idea behind the AIPW estimator is to use the predicted performance  $\hat{s}(f,z) = g[s(\mathcal{F}^{(s)},z)]$  as a proxy score to estimate  $\bar{s}(f,\mathcal{D})$  and "debias" that estimator as follows

$$h^{\text{AIPW}}(f) = \bar{s}(f, \mathcal{C}) + \frac{1}{1 + \frac{n}{N-n}} \left( \frac{1}{N-n} \sum_{z \in \mathcal{D} - \mathcal{C}} \hat{s}(f, z) - \frac{1}{n} \sum_{z \in \mathcal{C}} \hat{s}(f, z) \right). \tag{21}$$

Unlike the other learning-based baselines, AIPW is a consistent estimator for  $\bar{s}(f,\mathcal{D})$ [19]. Compared to RANDOM-SAMPLING, it reduces estimator variance by a factor of up to  $\frac{1}{1+\frac{n}{N}}\rho(\hat{s}^(f,z),s(f,z))^2$  [14], where  $\rho$  is the Pearson correlation coefficient. Recent research [37] shows that AIPW estimator will outperform random sampling if and only if the correlation between  $\hat{s}(f,z)$  and s(f,z) is above a certain level that depends on n.

# **Algorithm 1** PCA Impute Process

- 1: Input: Data matrix with missing values
- 2: **Parameters:** number of components k, max iteration max\_iter, stopping threshold tol
- 3: Output: Imputed data matrix
- 4: Step 1: Initialization
- 5: Compute initial values for missing entries using column means
- 6: Step 2: Iterative Imputation
- 7: **for** iteration  $\leftarrow 1$  to max\_iter **do**
- 8: **PCA Decomposition:**
- 9: Perform  $\overrightarrow{PCA}$  retaining k components
- 10: Transform data to the lower-dimensional space
- 11: Reconstruct the data from the lower-dimensional space
- 12: Evaluate Convergence:
- 13: Compute the norm of differences between imputed and original values at missing entries
- 14: **if** norm < tol **then**
- 15: Break the loop
- 16: **end if**
- 17: Update Imputed Values:
- 18: Replace missing values with reconstructed values
- 19: **end fo**i
- 20:
- 21: return Fully imputed data matrix

# **B** Additional Experiment Setup

We select a diverse range of benchmarks from the following sources<sup>7</sup>.

- HELM-Lite benchmarks [35]:
  - OpenbookQA [39]: N = 500 data points.
  - GSM8K [9]: N = 1000 data points.
  - LegalBench [22]: N=2047 data points.
  - Math [26]: N = 437 data points.
  - MedQA [28]: N = 1000 data points.
  - MMLU [25]: N=567 data points.

We obtain the per-data point performances of  $|\mathcal{F}|=83$  models from the official leaderboard. Note that Helm-Lite often only uses a subset of the original testing set for each benchmark to save compute.

- GLUE benchmarks [61]:
  - MRPC [13]: N = 408 data points.
  - RTE [11, 18, 4]: N = 277 data points.
  - SST-2 [55]: N = 872 data points.
  - MNLI [64]: N = 9815 data points.
  - QNLI [46]: N = 5463 data points.

We use the per-data performances of  $|\mathcal{F}| = 87$  models provided by AnchorPoint<sup>8</sup> [60].

- OpenLLM benchmarks [16]:
  - IFEval [67]: N = 541 data points.
  - Math [26]: N=894 data points. Only level 5 MATH questions are used in OpenLLM.
  - MMLU-Pro [62]: N = 12032 data points.
  - Arc-Challenge [8]: N = 1172 data points.
  - BBH [58]: N = 5761 data points.
  - GPQA [47]:  $N=1192~\mathrm{data}$  points.
  - MUSR [56]: N = 756 data points.

We use  $|\mathcal{F}| = 448$  models provided by Huggingface <sup>9</sup> and collect their performance scores.

• ImageNet [51]: We collect  $|\mathcal{F}|=110$  models from Pytorch Hub  $^{10}$  and evaluate them on ImageNet with N=50,000 data points.

For simplicity, we report the overall average accuracy directly for MMLU, MMLU-Pro, and BBH, rather than the weighted average accuracy computed across sub-tasks. Alternatively, one could apply benchmark predictions separately to each sub-task and then calculate the weighted average accuracy.

<sup>&</sup>lt;sup>7</sup>Since P-IRT and GP-IRT requires s(f, z) to be binary, we only use benchmarks with accuracy as metric.

<sup>&</sup>lt;sup>8</sup>The provided score file for QQP is broken so we exclude it.

<sup>9</sup>https://huggingface.co/spaces/open-llm-leaderboard/open\_llm\_leaderboard#

<sup>10</sup> https://pytorch.org/vision/stable/models.html#classification

Table 2: Training and inference time of each method on ImageNet with N=50000 data points and  $|\mathcal{F}|=110$  models. Training is based on 83 source models, and inference is on 27 target models.

	Training Time (s)	Inference Time (s)
RANDOM-SAMPLING	0.00	0.00
RANDOM-SAMPLING-LEARN	0.02	0.00
PCA	0.59	19.20
AIPW	0.00	0.27
RANDOM-SEARCH-LEARN	81.02	0.00
Lasso	105.58	0.01
DOUBLE-OPTIMIZE	4.88	0.00
ANCHOR-POINTS-WEIGHTED	84.26	0.00
ANCHOR-POINTS-PREDICTOR	197.71	0.26
P-IRT	585.72	0.90
GP-IRT	1750.20	0.89

Table 3: Average estimation gap between the predicted rankings based on the coreset and the actual rankings based on the full benchmark, measured by Kendall's  $\tau$  ( $\uparrow$ ). The results are averaged over all benchmarks.

	Interpolation				Extrapolation					
	n = 10	n = 20	n = 50	n = 100	n = 200	n = 10	n = 20	n = 50	n = 100	n = 200
RANDOM-SAMPLING	0.52	0.61	0.70	0.78	0.84	0.36	0.43	0.53	0.63	0.73
RANDOM-SAMPLING-LEARN	0.57	0.66	0.75	0.81	0.86	0.07	0.12	0.18	0.27	0.36
PCA	0.55	0.63	0.72	0.78	0.83	0.04	0.10	0.21	0.40	0.57
AIPW	0.52	0.62	0.72	0.79	0.84	0.33	0.40	0.51	0.61	0.70
RANDOM-SEARCH-LEARN	0.66	0.70	0.76	0.82	0.86	0.13	0.13	0.20	0.29	0.38
Lasso	0.68	0.71	0.77	0.81	0.82	0.05	0.06	0.12	0.19	0.22
DOUBLE-OPTIMIZE	0.58	0.66	0.76	0.81	0.84	0.31	0.36	0.44	0.50	0.58
ANCHOR-POINTS-WEIGHTED	0.65	0.70	0.76	0.81	0.85	0.37	0.43	0.50	0.60	0.69
ANCHOR-POINTS-PREDICTOR	0.67	0.72	0.77	0.80	0.80	0.21	0.25	0.32	0.38	0.44
P-IRT	0.52	0.58	0.71	0.80	0.87	0.28	0.31	0.42	0.56	0.69
GP-IRT	0.53	0.59	0.72	0.80	0.86	0.28	0.33	0.45	0.59	0.71

# **C** Additinoal Experiment Results

# C.1 Detailed Results

In this paper, we experiment with  $n \in \{10, 20, 50, 100, 200\}$  under both interpolation and extrapolation settings. The detailed results, along with standard errors, are reported in Figures 6, 7, 8, 9, 10, 11, 12, 13, 14, and 15.

#### **C.2** Running time

While some of the benchmark prediction methods could potentially benefit from the use of GPUs, we opted to run all methods without them, as they are sufficiently fast on standard hardware. Table 2 presents the training and inference times for each method on ImageNet. Among the models, GPIRT is the slowest during training because it involves fitting a large Item Response Theory (IRT) model. During inference, PCA is the slowest, as it requires multiple imputations of the entire matrix. Although AIPW needs training a separate regressor for each target model during inference, the regressor is small, making the inference process remain efficient.

# **C.3** Ranking Preservation

We further compare the predicted rankings of target models with the actual rankings based on the full benchmark using Kendall's  $\tau$ . Specifically, we calculate Kendall's  $\tau$  for each random trial and average the results over 100 trials. Our conclusions mostly remain unchanged, with almost all benchmark prediction methods outperforming Random Sampling under interpolation, while none can surpass RANDOM-SAMPLING under extrapolation.

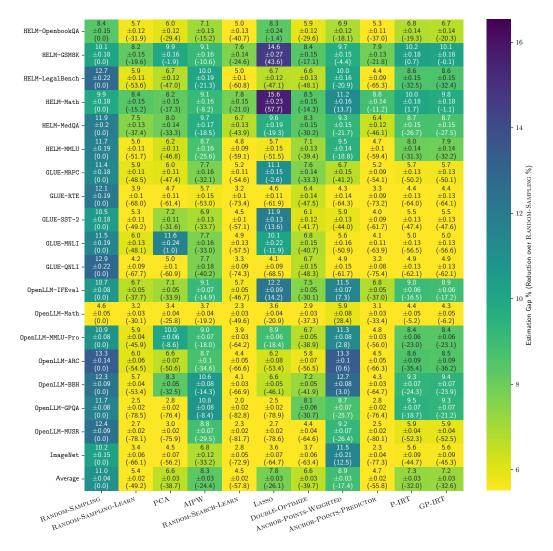


Figure 6: The estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1) under interpolation model split, where source models are identically distributed with target models. Each target model can only be evaluated on n=10 data points. We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

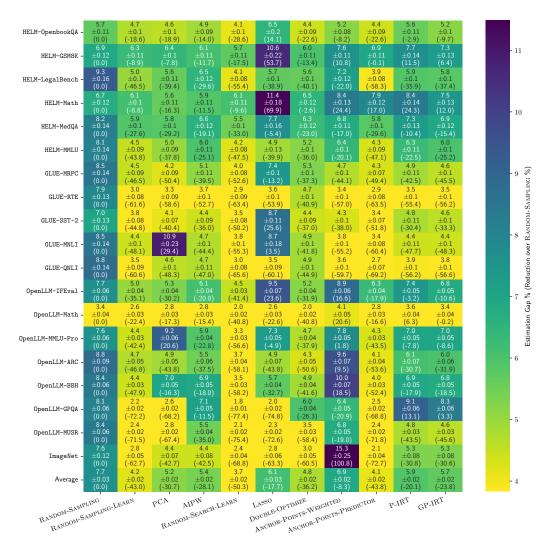


Figure 7: The estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1) under interpolation model split, where source models are identically distributed with target models. Each target model can only be evaluated on n=20 data points. We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

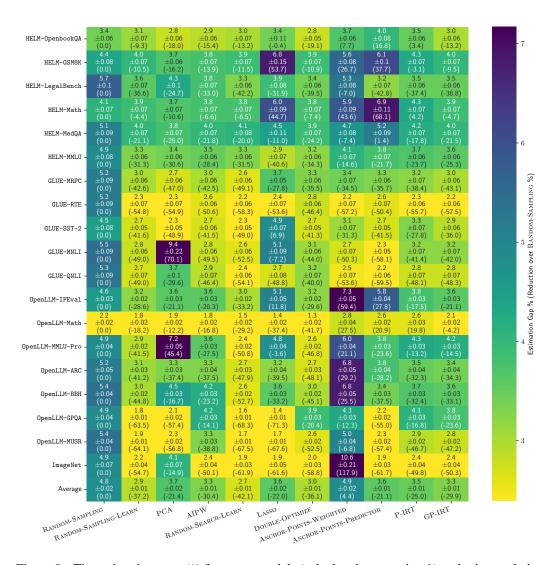


Figure 8: The estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1) under interpolation model split, where source models are identically distributed with target models. Each target model can only be evaluated on n=50 data points. We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

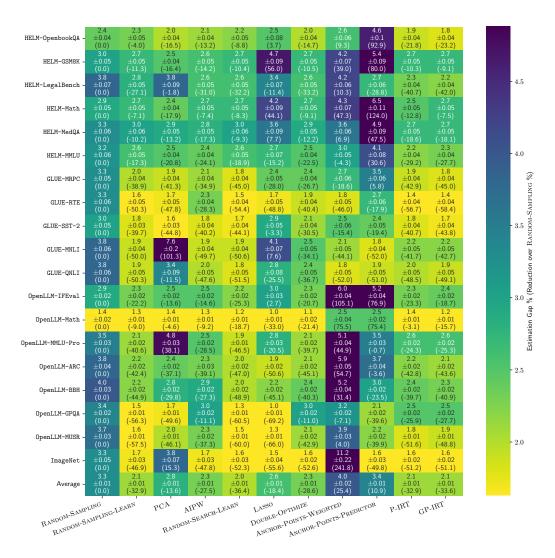


Figure 9: The estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1) under interpolation model split, where source models are identically distributed with target models. Each target model can only be evaluated on n=100 data points. We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

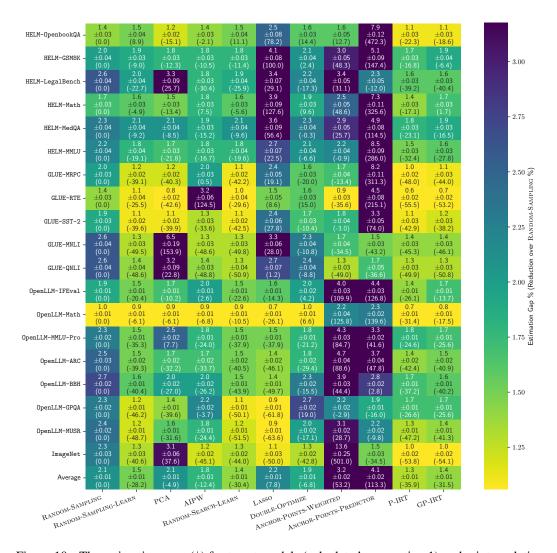


Figure 10: The estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1) under interpolation model split, where source models are identically distributed with target models. Each target model can only be evaluated on n=200 data points. We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

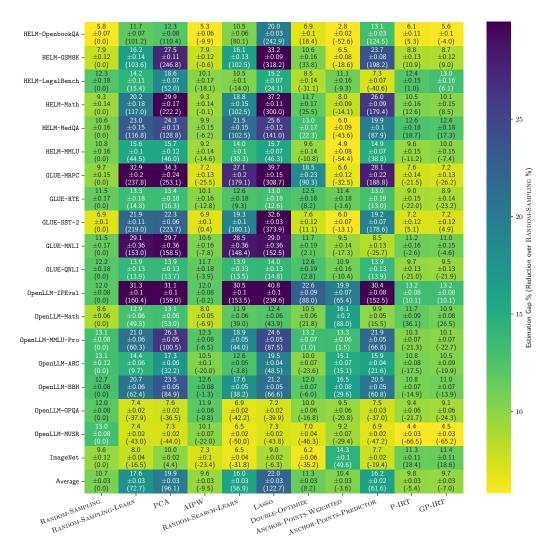


Figure 11: The estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1) under extrapolation model split, where source models are the lowest-performing 50%, and target models are the top 30% based on average performance over the full benchmark. Each target model can only be evaluated on n=10 data points. We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

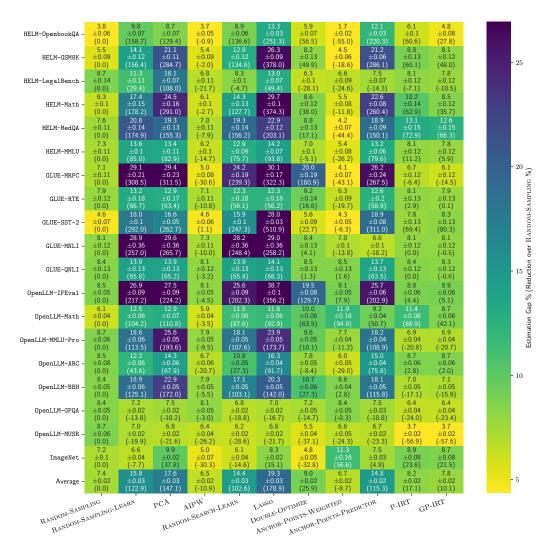


Figure 12: The estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1) under extrapolation model split, where source models are the lowest-performing 50%, and target models are the top 30% based on average performance over the full benchmark. Each target model can only be evaluated on n=20 data points. We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

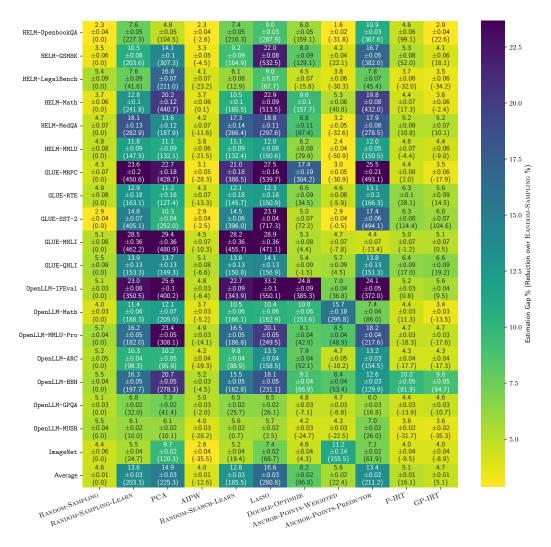


Figure 13: The estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1) under extrapolation model split, where source models are the lowest-performing 50%, and target models are the top 30% based on average performance over the full benchmark. Each target model can only be evaluated on n=50 data points. We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

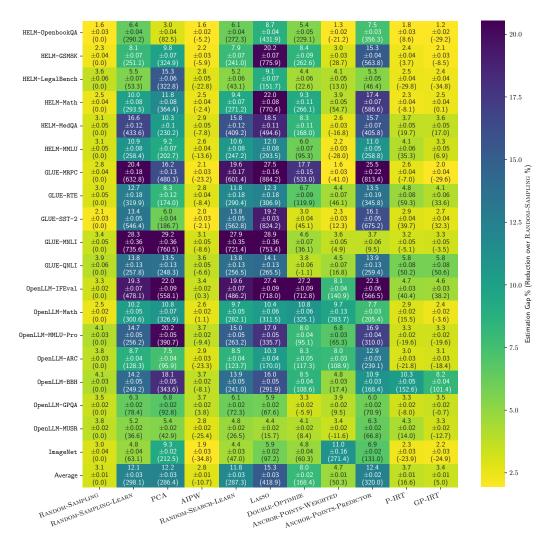


Figure 14: The estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1) under extrapolation model split, where source models are the lowest-performing 50%, and target models are the top 30% based on average performance over the full benchmark. Each target model can only be evaluated on n=100 data points. We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

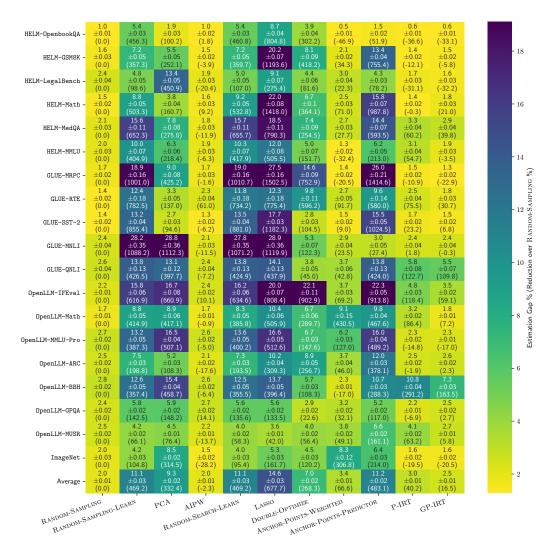


Figure 15: The estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1) under extrapolation model split, where source models are the lowest-performing 50%, and target models are the top 30% based on average performance over the full benchmark. Each target model can only be evaluated on n=200 data points. We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

## C.4 Case Studies

We further investigate two additional experimental settings that deviate from the primary setting in the main paper.

**Fewer source models under interpolation.** Different from the previous interpolation setting that utilized 75% of models as source models, we now use only 10 models as source models for each benchmark and use the rest as target models. All other settings remain unchanged. This setting allows us to assess the effectiveness of benchmark prediction when "training data" from source models is more limited. Results are shown in Figure 16. Consistent with the findings in the paper, most methods still outperform RANDOM-SAMPLING, while RANDOM-SEARCH-LEARN and RANDOM-SAMPLING-LEARN remain to be the best-performing methods.

**Near extrapolation.** We modify the previous extrapolation setting, which used the lowest-performing 50% of models as source models and the top 30% as target models. In this new setting, we designate the top 25% of models as target models and utilize all remaining models as source models. All other settings remain unchanged. This setup enables us to examine whether benchmark prediction methods demonstrate improved performance when the distribution gap between source and target models is reduced. Results are shown in Figure 17. Consistent with the findings in the paper, most methods fail to consistently outperform RANDOM-SAMPLING, except for AIPW.

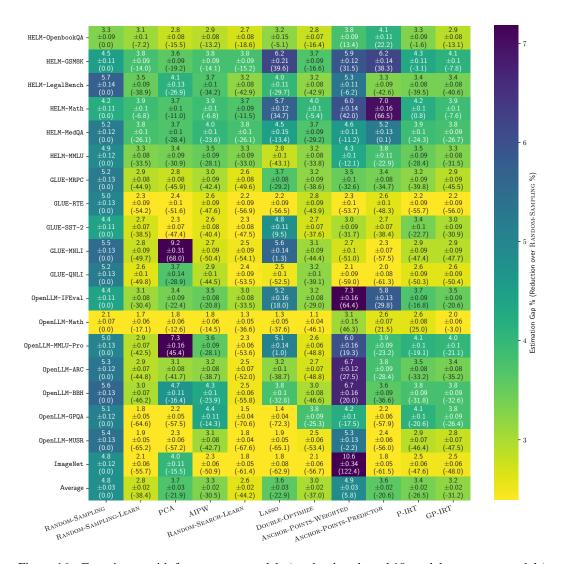


Figure 16: Experiment with fewer source models (randomly selected 10 models as source models) under the interpolation model split. We report the estimation gaps  $(\downarrow)$  for target models (calculated as equation 1). We also report  $\pm$  the standard error of the mean and the estimation gap reduction  $(\downarrow)$  over Random-Sampling in parentheses. A negative reduction implies that the method achieves a lower estimation gap than Random-Sampling. % is omitted. Best viewed in color.

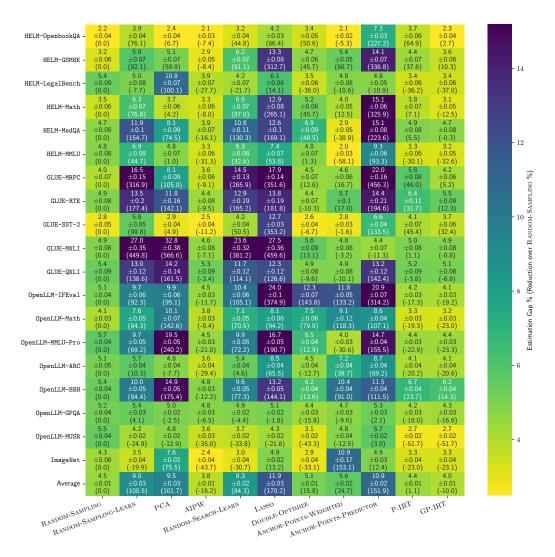


Figure 17: Experiment with the near extrapolation model split by using the top 25% of available models as target models and the remaining bottom 75% models as source models. We report the estimation gaps ( $\downarrow$ ) for target models (calculated as equation 1). We also report  $\pm$  the standard error of the mean and the estimation gap reduction ( $\downarrow$ ) over RANDOM-SAMPLING in parentheses. A negative reduction implies that the method achieves a lower estimation gap than RANDOM-SAMPLING. % is omitted. Best viewed in color.

# **D** Broarder Impacts and Limitations

This paper addresses the benchmark prediction problem in scenarios with limited data. One potential limitation of our study is the relatively small number of models examined. For both the HELM-Lite and GLUE benchmarks, we have collected full benchmark results for fewer than 100 models. Despite conducting 100 random trials for each experiment, including additional and more diverse models could further strengthen the comprehensiveness and robustness of our analysis.

We do not anticipate any direct societal impacts from this work, such as potential malicious or unintended uses, nor do we foresee any significant concerns involving fairness, privacy, or security considerations. Additionally, we have not identified potential harms resulting from the application of this technology.