
How Benchmark Prediction from Fewer Data Misses the Mark

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Evaluating large language models (LLMs) is increasingly costly, motivating meth-
2 ods to speed up evaluation by compressing benchmark datasets. Benchmark
3 prediction aims to select a small subset of evaluation points and predict overall
4 performance from that subset. We systematically assess 11 benchmark prediction
5 methods across 19 benchmarks. First, we identify a strong baseline: take a random
6 sample and fit a regression to predict the missing entries, which outperforms most
7 existing methods and challenges the need for careful subset selection. Second, we
8 show that all methods rely on model similarity: performance degrades markedly
9 when extrapolating to stronger models than those used for training, where few meth-
10 ods beat a simple sample average. We introduce an augmented inverse propensity
11 weighting (AIPW) estimator that consistently improves over the random sample
12 average under both interpolation and extrapolation, though gains remain modest
13 and still depend on similarity. This shows that benchmark prediction fails just
14 when it is most needed: at the evaluation frontier, where the goal is to evaluate new
15 models of unknown capabilities.

16 1 Introduction

17 Computational cost is a major bottleneck in evaluating recent generative models. For example,
18 evaluating a single 176B model on HELM required 4,200 GPU hours [35]; even large organizations
19 report heavy costs on BIG-bench [17]. This has prompted work on efficient LLM evaluation through
20 *benchmark prediction*: finding a subset of data points to evaluate on and predicting benchmark
21 performance from these evaluations. The simplest method is the random sample mean: evaluate n
22 evaluation points and average, which gives an additive approximation up to error $O(1/\sqrt{n})$. Recent
23 work aims to improve this by selecting an informative *core set* and learning mappings from core set
24 to full-benchmark performance. See the discussion of related work in Appendix A.

25 We systematically study the strengths and limits of these methods by evaluating 11 benchmark
26 prediction methods across 19 benchmarks with at least 83 models each. Models are split into
27 source models (full performance data available) and target models (performance data for no more
28 than 50 points). Methods must estimate target models' mean performance using this constraint.
29 Effectiveness is measured by *average estimation gap*-the absolute difference between true and
30 estimated performances.

31 **Many methods work well on similar models, but a simple baseline works best.** In the *interpolation*
32 regime (source and target models from same distribution), a remarkably simple method works
33 best: RANDOM-SAMPLING-LEARN-random sampling followed by regression modeling-reduces
34 the estimation gap by 37% compared to basic random sampling, outperforming most sophisticated
35 methods. This suggests that the manner of core-set selection is relatively unimportant; rather, the key
36 to success is modeling the correlation between core-set and full-benchmark performances.

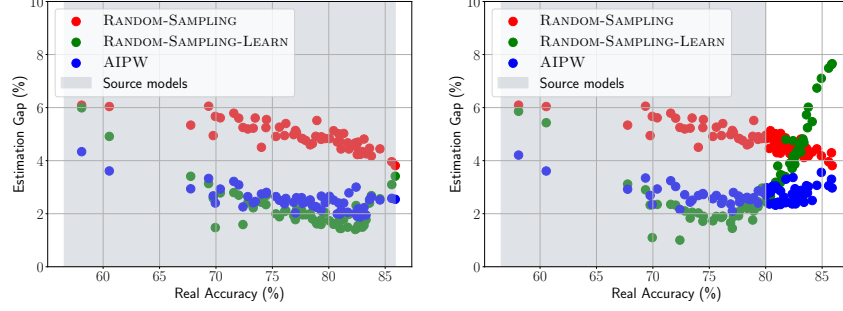


Figure 1: Estimation gap (equation 1) versus real accuracy on ImageNet. Gray shows source model accuracy range. Left: source models randomly sampled across all models. Right: source models sampled from models with lower than 80% accuracy.

Methods fail at the evaluation frontier. In the *extrapolation* regime (target models all better than source models), effectiveness drops sharply. Most methods fail to beat naive random sampling when evaluating new, better models—precisely when efficient evaluation is most needed (Figure 1, right).

AIPW is an overlooked exception to the rule. We introduce augmented inverse propensity weighting (AIPW) to benchmark prediction. Unlike other methods, AIPW consistently outperforms random sampling in both interpolation and extrapolation settings. 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. Our further study reveals that benchmark prediction methods rely heavily on model similarity [38]: methods that beat RANDOM-SAMPLING do so mainly for targets similar to sources, while accuracy on dissimilar models deteriorates. In contrast, RANDOM-SAMPLING exhibits neutral correlation.

2 What is Benchmark Prediction?

Problem formulation. A benchmark is defined as $(\mathcal{D}, \mathcal{F}, s)$ where \mathcal{D} is the dataset with N data points, \mathcal{F} is the model set, and s is the evaluation metric. For any model $f \in \mathcal{F}$ and data point $z = (x, y) \in \mathcal{D}$, we define notation as follows.

- $s(f, z)$ denotes performance of f on point z , for example, $\mathbb{1}[f(x) = y]$ for accuracy.
- $\bar{s}(f, \mathcal{D}') = \frac{1}{|\mathcal{D}'|} \sum_{z \in \mathcal{D}'} s(f, z)$ denotes average performance on subset $\mathcal{D}' \subset \mathcal{D}$.
- $\mathbf{s}(f, \mathcal{D}')$ denotes the vectorized performance of f on $\mathcal{D}' \subset \mathcal{D}$, and $\mathbf{s}(\mathcal{F}', z)$ denotes the vectorized performances of all models in $\mathcal{F}' \subset \mathcal{F}$ on data point z .
- $S(\mathcal{F}', \mathcal{D}')$ denotes the performance matrix for models $\mathcal{F}' \subset \mathcal{F}$ on points $\mathcal{D}' \subset \mathcal{D}$.

Given source models $\mathcal{F}^{(s)} \subset \mathcal{F}$ with known full performance $S(\mathcal{F}^{(s)}, \mathcal{D})$ and target models $\mathcal{F}^{(t)} = \mathcal{F} \setminus \mathcal{F}^{(s)}$ evaluated on only $n \ll N$ points, benchmark prediction aims to estimate $\bar{s}(f, \mathcal{D})$ for each $f \in \mathcal{F}^{(t)}$ by: ① selecting core-set $\mathcal{C} \subset \mathcal{D}$ with $|\mathcal{C}| = n$, and ② learning estimator h to minimize:

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

Previous benchmark prediction methods:

- RANDOM-SAMPLING: pick \mathcal{C} at random; return the mean on \mathcal{C} .
- ANCHOR-POINTS-WEIGHTED [60]: k-medoids to select \mathcal{C} ; return weighted sum by cluster density.
- ANCHOR-POINTS-PREDICTOR [60]: as above, then linear regression from $\mathbf{s}(f, \mathcal{C})$ to $\bar{s}(f, \mathcal{D})$.
- P-IRT [43]: as above, replace regression with the Item Response Theory (IRT) model.
- GP-IRT [43]: combine P-IRT with Anchor-Points-Weighted aggregation.

New methods introduced:

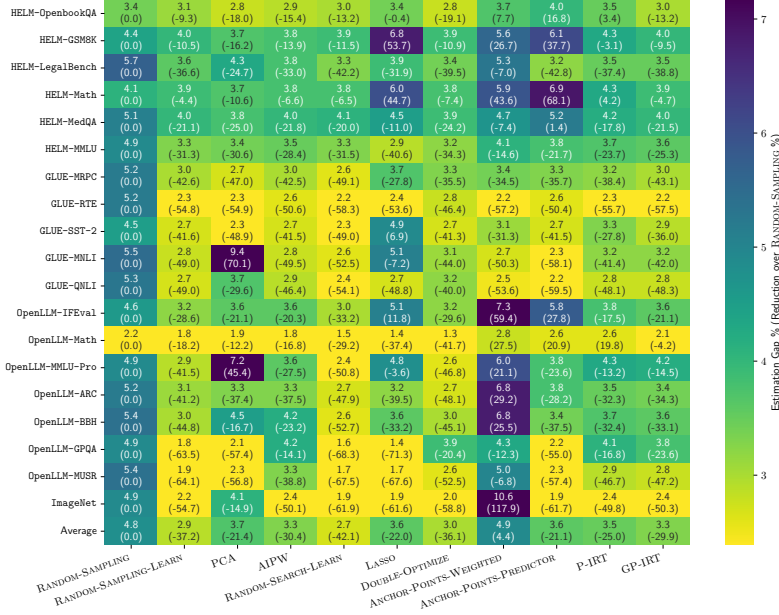


Figure 2: The estimation gaps (\downarrow) for target models (equation 1) under the interpolation split, where source and target models are identically distributed. Each target is evaluated on $n = 50$ data points. The estimation gap reduction (\downarrow) over RANDOM-SAMPLING is shown in parentheses.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

- RANDOM-SAMPLING-LEARN: pick \mathcal{C} at random; Ridge regression from $s(f, \mathcal{C})$ to $\bar{s}(f, \mathcal{D})$.
- RANDOM-SEARCH-LEARN: run RANDOM-SAMPLING-LEARN 10,000 times; select best \mathcal{C} and h .
- LASSO: Lasso regression with sparsity constraint number of non-zero weights $\leq n$.
- DOUBLE-OPTIMIZE: gradient descent for joint core-set and regression optimization.
- PCA: use random \mathcal{C} ; impute target scores with PCA assuming $S(\mathcal{F}, \mathcal{D})$ is low-rank.
- AIPW [48]: train regression g to predict $s(f, z)$ from $s(\mathcal{F}^{(s)}, z)$ for every f . The idea is to use the predicted performance $\hat{s}(f, z) = g[s(\mathcal{F}^{(s)}, z)]$ as a proxy score of $s(f, z)$ and “debias” as follows

$$h^{\text{AIPW}}(f) = \bar{s}(f, \mathcal{C}) + \frac{1}{1 + \frac{n}{N-n}} \left(\frac{1}{N-n} \sum_{z \in \mathcal{D}-\mathcal{C}} \hat{s}(f, z) - \frac{1}{n} \sum_{z \in \mathcal{C}} \hat{s}(f, z) \right). \quad (2)$$

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]. See more details in Appendix B.

3 Experiments

We examine the 11 benchmark prediction methods under both interpolation and extrapolation settings. We select 19 benchmarks from HELM-Lite [35], OpenLM [16], GLUE [61] and ImageNet [51], each with at least 83 models. See more details in Appendix C.

Estimation gap reduction under interpolation. For each benchmark, we randomly select 75% of models as source models $\mathcal{F}^{(s)}$ with full performance scores $S(\mathcal{F}^{(s)}, \mathcal{D})$ available. The remaining 25% serve as target models $\mathcal{F}^{(t)}$, each evaluated on only $n = 50$ data points. We report average estimation gap across all target models in 100 random trials. See standard errors in Appendix D.

Figure 2 shows that compared to RANDOM-SAMPLING, most methods effectively reduce the estimation gap, with nine of ten achieving more than 20% average reduction across benchmarks. The simple baseline RANDOM-SEARCH-LEARN performs best with 42.1% average reduction, outperforming the previous state-of-the-art GP-IRT (29.9% average reduction) on nearly all benchmarks.

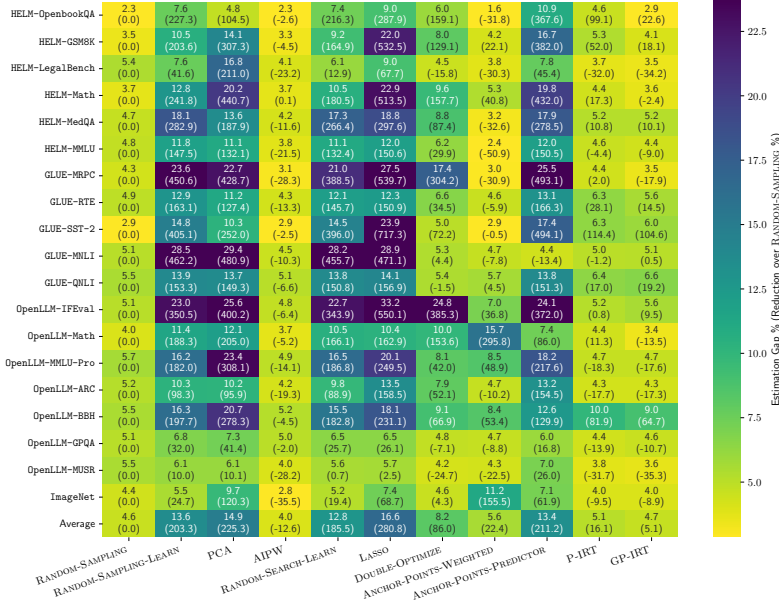


Figure 3: The estimation gaps (\downarrow) for target models (equation 1) under extrapolation split, where source models are the lowest-performing 50%, and target models are the top 30%. Each target model is evaluated on $n = 50$ data points. We also report the estimation gap reduction (\downarrow) over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

Core-set selection does not significantly enhance effectiveness. RANDOM-SAMPLING-LEARN achieves 37.2% average reduction using only Ridge regression on random samples, performing comparably to RANDOM-SEARCH-LEARN despite the latter’s 10,000 optimization iterations. It also surpasses methods with sophisticated subset selection like DOUBLE-OPTIMIZE and GP-IRT. These suggest the primary driver of success is learning to predict the mean rather than core-set selection.

Estimation gap increase under extrapolation. We then examine all 11 methods under *extrapolation*. Models are ranked by full benchmark performance $\bar{s}(f, \mathcal{D})$. The bottom 50% become source models, while the top 30% serve as target models, reflecting real-world scenarios where developers assess improved models based on existing inferior ones.

Figure 3 reveals striking differences from interpolation. While RANDOM-SAMPLING’s estimation gap remains similar (4.6% vs 4.8%), all other methods deteriorate significantly. The previously best RANDOM-SEARCH-LEARN now shows 185.1% increase in estimation gap versus RANDOM-SAMPLING, performing worse across all benchmarks. Only AIPW still outperforms RANDOM-SAMPLING (in 18/19 benchmarks) as it is a consistent estimator, though its advantage shrinks from -30.4% to -12.6% reduction.

This contrast underscores most methods’ heavy reliance on source-target similarity. See a deeper analysis of model similarity and ablation studies in Appendix D. While traditional machine learning emphasizes in-domain performance, benchmarking aims to identify superior new models, making extrapolation more relevant than interpolation. The decline in benchmark prediction effectiveness under extrapolation calls for more caution.

4 Conclusion

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. See more detailed discussion in Appendix E.

References

- [1] Anastasios Nikolas Angelopoulos, Stephen Bates, Clara Fannjiang, Michael I. Jordan, and Tijana Zrnica. Prediction-powered inference. *Science*, 382:669 – 674, 2023.
- [2] Anastasios Nikolas Angelopoulos, John C. Duchi, and Tijana Zrnica. 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 textual 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, M’irian 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 Udandaraao, 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 llm 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, Kavin 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, Ziyang Jiang, Tianle Li, Max W.F. Ku, Kai Wang, Alex Zhuang, Rongqi "Richard" Fan, Xiang Yue, and Wenhui 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, Zirui Chen, Yu Su, and Huan Sun. Automatic evaluation of attribution by large language models. In *Conference on Empirical Methods in Natural Language Processing*, 2023.

- 303 [67] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny
304 Zhou, and Le Hou. Instruction-following evaluation for large language models. *ArXiv*,
305 abs/2311.07911, 2023.
- 306 [68] Jin Peng Zhou, Christian K. Belardi, Ruihan Wu, Travis Zhang, Carla P. Gomes, Wen Sun, and
307 Kilian Q. Weinberger. On speeding up language model evaluation. *ArXiv*, abs/2407.06172,
308 2024.
- 309 [69] Xiao Zhou, Renjie Pi, Weizhong Zhang, Yong Lin, Zonghao Chen, and T. Zhang. Probabilistic
310 bilevel coreset selection. *ArXiv*, abs/2301.09880, 2023.
- 311 [70] Vilém Zouhar, Peng Cui, and Mrinmaya Sachan. How to select datapoints for efficient human
312 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 F.

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 cross-referenced.
- 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 C 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 C.

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 D.

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 D.

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]

Justification: See Appendix F.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- 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.
- 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 3.

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]

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

A 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 data¹.

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, 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.

¹Unlike 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.

B Details of Benchmark Prediction Methods

B.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 $\mathcal{D}' \subset \mathcal{D}$.
 - $\mathbf{s}(f, \mathcal{D}') = \{s(f, z)\}_{z \in \mathcal{D}'}$ represents the vectorized performance of $f \in \mathcal{F}$ on all data points in $\mathcal{D}' \subset \mathcal{D}$, and $\mathbf{s}(\mathcal{F}', z) = \{s(f, z)\}_{f \in \mathcal{F}'}$ represents the vectorized performances of all models in $\mathcal{F}' \subset \mathcal{F}$ on data point $z \in \mathcal{D}$.
 - $S(\mathcal{F}', \mathcal{D}') = \{\mathbf{s}(f, \mathcal{D}')\}_{f \in \mathcal{F}'} = \{\mathbf{s}(\mathcal{F}', z)\}_{z \in \mathcal{D}'}^T$ 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,

$$\text{estimation gap: } \frac{1}{|\mathcal{F}^{(t)}|} \sum_{f \in \mathcal{F}^{(t)}} |\bar{s}(f, \mathcal{D}) - h[\mathbf{s}(f, \mathcal{C}), S(\mathcal{F}^{(s)}, \mathcal{D})]|. \quad (3)$$

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 $h[\mathbf{s}(f, \mathcal{C}), S(\mathcal{F}^{(s)}, \mathcal{D})]$.

B.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 performance,

$$h^{\text{RANDOM-SAMPLING}}(f) = \bar{s}(f, \mathcal{C}). \quad (4)$$

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 \mathcal{C} , and a weight vector $\boldsymbol{\theta} \in \mathbb{R}^n$ is calculated as the normalized cluster size of each medoid. The final estimate for any target model $f \in \mathcal{F}^{(t)}$ is

$$h^{\text{ANCHOR-POINTS-WEIGHTED}}(f) = \mathbf{s}(f, \mathcal{C})^T \boldsymbol{\theta}. \quad (5)$$

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

$$h^{\text{ANCHOR-POINTS-PREDICTOR}}(f) = \bar{\mathbf{g}}[\mathbf{s}(f, \mathcal{C})] \quad (6)$$

$$\text{where } \mathbf{g} = \arg \min_{\mathbf{g}'} \frac{1}{M} \sum_{f \in \mathcal{F}^{(s)}} \|\mathbf{s}(f, \mathcal{D} - \mathcal{C}) - \mathbf{g}'[\mathbf{s}(f, \mathcal{C})]\|_2^2, \quad (7)$$

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

- 700 • P-IRT [43] extends ANCHOR-POINTS-PREDICTOR by replacing the regression model g in
 701 equation 7 with an Item Response Theory (IRT) model. Following the notation for ANCHOR-
 702 POINTS-PREDICTOR, we estimate performance for any $f \in \mathcal{F}^{(t)}$ as follows:

$$h^{\text{P-IRT}}(f) = \frac{N-n}{N} \bar{g}[s(f, \mathcal{C})] + \frac{n}{N} \bar{s}(f, \mathcal{C}). \quad (8)$$

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

$$h^{\text{GP-IRT}}(f) = \lambda h^{\text{ANCHOR-POINTS-WEIGHTED}}(f) + (1 - \lambda) h^{\text{P-IRT}}(f), \quad (9)$$

705 where λ is chosen heuristically to control the error of P-IRT.

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

- 707 • RANDOM-SAMPLING-LEARN randomly samples a subset as \mathcal{C} and adopts a Ridge regression
 708 model g for estimation as follows,

$$h^{\text{RANDOM-SAMPLING-LEARN}}(f) = g[s(f, \mathcal{C})] \quad (10)$$

$$\text{where } g = \arg \min_{g'} \frac{1}{M} \sum_{f \in \mathcal{F}^{(s)}} |\bar{s}(f, \mathcal{D}) - g'[s(f, \mathcal{C})]|. \quad (11)$$

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

$$h^{\text{LASSO}}(f) = s(f, \mathcal{C})^T \theta_{\mathcal{C}} \quad (12)$$

$$\text{where } \theta = \arg \min_{\theta'} \frac{1}{n} \sum_{z \in \mathcal{C}} [s(f, \mathcal{D})^T \theta' - \bar{s}(f, \mathcal{D})]^2 + \lambda \|\theta'\|_1, \quad (13)$$

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

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

$$h^{\text{DOUBLE-OPTIMIZE}}(f) = [s(f, \mathcal{D}) \cdot \text{TopMask}(\pi; n)]^T \theta \quad (14)$$

$$\text{where } \pi, \theta = \arg \min_{\pi', \theta'} \{[s(f, \mathcal{D}) \cdot \text{TopMask}(\pi'; n)]^T \theta' - \bar{s}(f, \mathcal{D})\}^2, \quad (15)$$

716 where \cdot refers to the bitwise multiplication between two vectors, and $\text{TopMask}(\pi'; n)$ replaces
 717 the top n largest values of π' with 1s and the rest with 0s. We directly pass the gradient on
 718 $\text{TopMask}(\pi'; n)$ to π' during optimization following the Straight-Through technique [27, 3].

- 719 • Principal Component Analysis (PCA) treats benchmark prediction as a matrix completion problem.
 720 This method assumes the performance matrix $S(\mathcal{F}, \mathcal{D})$ is of low rank. By randomly sampling a
 721 subset as \mathcal{C} , this methods conducts PCA to impute the missing values for target models [59, 6]. As
 722 a more intuitive view, one could also take the acquired principal components as model capability
 723 indicators [50], *i.e.*, the $(M \times k)$ PCA-transformed scores indicate the k -capabilities of each
 724 model, while the $(k \times N)$ principal components represent the capability requirements for each
 725 data point. We select k among $\{2, 5, 10, 20\}$ through cross-validation. The Pseudo codes are in
 726 Algorithm 1.

- 727 • Augmented inverse propensity weighting (AIPW) [48]: Inspired by the application of prediction
 728 powered inference [2, 1] to the LLM-as-a-judge setting [5, 14], we apply a more general AIPW
 729 estimator to benchmark prediction. We train a Ridge regression model g for every target model f ,
 730 which predicts the point-wise performance $s(f, z)$ based on $s(\mathcal{F}^{(s)}, z)$. Formally,

$$g = \arg \min_{g'} \frac{1}{n} \sum_{z \in \mathcal{C}} [g'[s(\mathcal{F}^{(s)}, z)] - s(f, z)]^2. \quad (16)$$

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

$$h^{\text{AIPW}}(f) = \bar{s}(f, \mathcal{C}) + \frac{1}{1 + \frac{n}{N-n}} \left(\frac{1}{N-n} \sum_{z \in \mathcal{D}-\mathcal{C}} \hat{s}(f, z) - \frac{1}{n} \sum_{z \in \mathcal{C}} \hat{s}(f, z) \right). \quad (17)$$

733 Unlike the other learning-based baselines, AIPW is a consistent estimator for $\bar{s}(f, \mathcal{D})$ [19].
 734 Compared to RANDOM-SAMPLING, it reduces estimator variance by a factor of up to
 735 $\frac{1}{1 + \frac{n}{N}} \rho(\hat{s}(f, z), s(f, z))^2$ [14], where ρ is the Pearson correlation coefficient. Recent research [37]
 736 shows that AIPW estimator will outperform random sampling if and only if the correlation between
 737 $\hat{s}(f, z)$ 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 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  $< \text{tol}$  then
15:     Break the loop
16:   end if
17:   Update Imputed Values:
18:     Replace missing values with reconstructed values
19: end for
20:
21: return Fully imputed data matrix

```

C Additional Experiment Setup

We select a diverse range of benchmarks from the following sources².

- 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 AnchorPoint³ [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]: $N = 1192$ data points.
- MUSR [56]: $N = 756$ data points.

We use $|\mathcal{F}| = 448$ models provided by Huggingface⁴ and collect their performance scores.

- ImageNet [51]: We collect $|\mathcal{F}| = 110$ models from Pytorch Hub⁵ 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.

²Since P-IRT and GP-IRT requires $s(f, z)$ to be binary, we only use benchmarks with accuracy as metric.

³The provided score file for QQP is broken so we exclude it.

⁴https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard#

⁵<https://pytorch.org/vision/stable/models.html#classification>

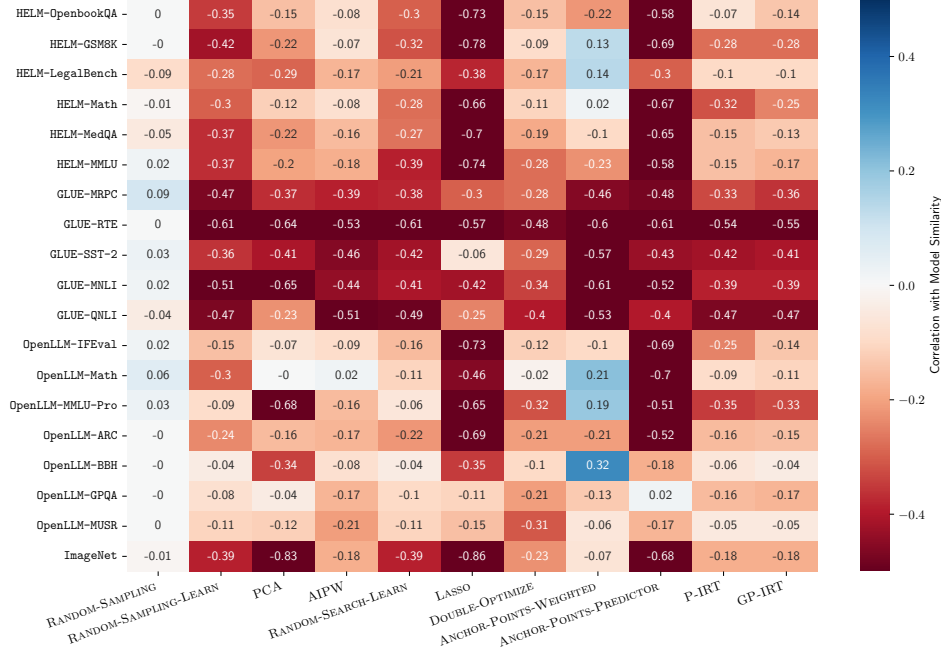


Figure 4: The Pearson correlation between normalized per-model estimation gap (equation 20) and model similarity (equation 18). Negative correlation indicates that target models that are dissimilar to source models tend to have larger estimation gap, and vice versa.

D Additinoal Experiment Results

D.1 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,

$$\mathcal{S}(f, \mathcal{F}^{(s)}, \mathcal{D}) = \frac{1}{M} \sum_{f' \in \mathcal{F}^{(s)}} \frac{c_{obs} - c_{exp}}{1 - c_{exp}}. \quad (18)$$

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, $c_{obs} = \frac{1}{N} \sum_{z \in \mathcal{D}} \mathbb{1}[s(f, z) = s(f', z)]$ 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

$$|\bar{s}(f, \mathcal{C}) - \bar{s}(f, \mathcal{D})| < \sigma_f / \sqrt{\alpha n} \quad (19)$$

with probability at least $(1 - \alpha)$. In other words, the performance of target models with lower σ_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 σ_f coincidentally have high $\mathcal{S}(f)$, while those with high σ_f have low

792 $\mathcal{S}(f)$, then a spurious correlation between estimation gap and model similarity to target models could
 793 appear even for RANDOM-SAMPLING. To prevent this, we define the normalized estimation gap as

$$\text{normalized estimation gap for } f: \quad \mathcal{E}(f) = \frac{1}{\sigma_f} |\bar{s}(f, \mathcal{D}) - h(f)|. \quad (20)$$

794 Then we measure the Pearson correlation between model similarity in equation 18 and the normalized
 795 estimation gap in equation 20.

796 **Results.** The results are shown in Figure 4. A clear negative correlation between model similarity
 797 and estimation gap emerges for almost all benchmark prediction methods except for RANDOM-
 798 SAMPLING. In particular, the best-performing method under the interpolation model split, RANDOM-
 799 SAMPLING-LEARN, exhibits a negative correlation below -0.2 in 13/19 benchmarks. Despite its
 800 asymptotic unbiasedness, we also find negative correlations for AIPW. This is perhaps unsurprising:
 801 While AIPW is consistent independent of how well its regression model $g[s(\mathcal{F}^{(s)}, z)]$ predicts
 802 $s(f, z)$, its variance depends precisely on that prediction quality. If the predictions are good, AIPW
 803 improves substantially over RANDOM-SAMPLING, while there is no improvement when predictions
 804 are fully uninformative. But intuitively, predicting $s(f, z)$ is harder when f is very different from the
 805 models $\mathcal{F}^{(s)}$ used for training the predictor $g[s(\mathcal{F}^{(s)}, z)]$.

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.

	Interpolation					Extrapolation				
	$n = 10$	$n = 20$	$n = 50$	$n = 100$	$n = 200$	$n = 10$	$n = 20$	$n = 50$	$n = 100$	$n = 200$
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

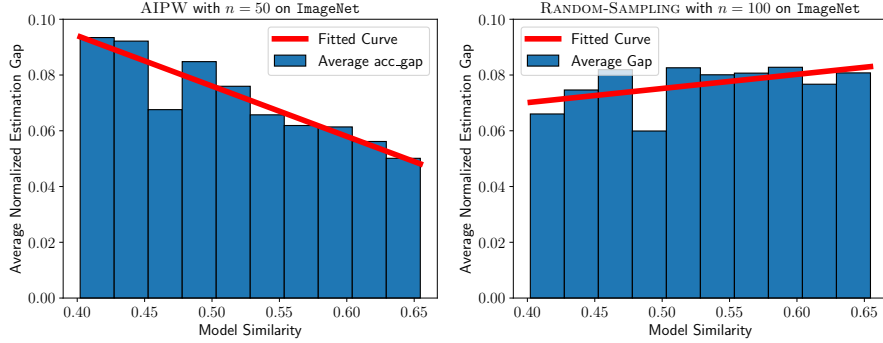


Figure 5: Average normalized estimation gap relative to model similarity for AIPW ($n=50$) and RANDOM-SAMPLING ($n=100$) on ImageNet. Each bar represents the target models whose similarity to source models falls within the corresponding range. The normalized estimation gap is defined as shown in equation 20. On average, AIPW outperforms RANDOM-SAMPLING, even with half the data. However, RANDOM-SAMPLING shows better performance when model similarity is low.

806 D.2 Ablation on Core-set Size

807 We conduct an ablation study on the size of the core-set n . We experiment with $n \in$
808 $\{10, 20, 50, 100, 200\}$, and the summarized results are shown in Table 1 (detailed results can be
809 found in Figures 6, 7, 8, 9, 10, 11, 12, 13, 14, and 15. As expected, the estimation gap generally de-
810 creases as n increases for most methods. Our previous conclusions remain valid across both settings.
811 With larger core-set sizes, most methods continue to perform better than RANDOM-SAMPLING in the
812 interpolation split but fail to do so in the extrapolation model split. Interestingly, we also find that
813 RANDOM-SAMPLING outperforms all other methods when given twice as much data, even in the
814 interpolation model split.

815 AIPW remains effective in both settings. However, its advantage over RANDOM-SAMPLING
816 diminishes as n increases. While AIPW reduces the estimation gap by -30.4% in interpolation and
817 -12.6% in extrapolation for $n = 50$, these advantages shrink to -12.4% in interpolation and a mere
818 -2.3% in extrapolation for $n = 200$. This is because the estimator variance reduction factor of AIPW
819 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
820 when the dataset is large and thus $\frac{n}{N}$ is small. Figure 5 compares AIPW with $n = 50$ to RANDOM-
821 SAMPLING with $n = 100$ data points using ImageNet. AIPW achieves a lower average normalized
822 estimation gap compared to RANDOM-SAMPLING, despite using only half the data. However, the
823 normalized estimation gap for AIPW is biased with respect to model similarity. In contrast, the
824 normalized estimation gap under RANDOM-SAMPLING remains largely neutral regarding model
825 similarity. Consequently, while AIPW reduces the average, it produces a higher gap for models with
826 low similarity compared to RANDOM-SAMPLING with twice the data.

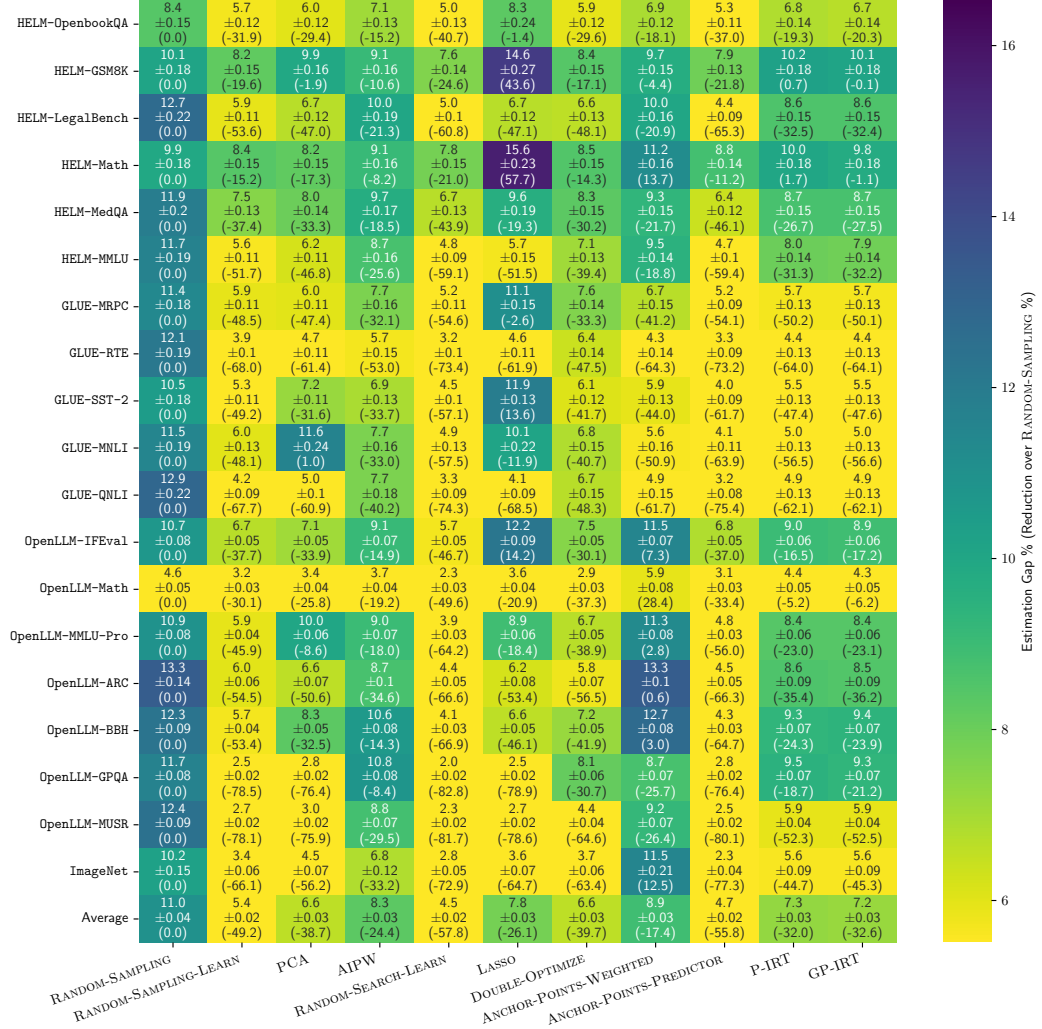


Figure 6: The estimation gaps (\downarrow) for target models (calculated as equation 1) under interpolation model split, where source models are identically distributed with target models. Each target model can only be evaluated on $n = 10$ data points. We also report \pm the standard error of the mean 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.

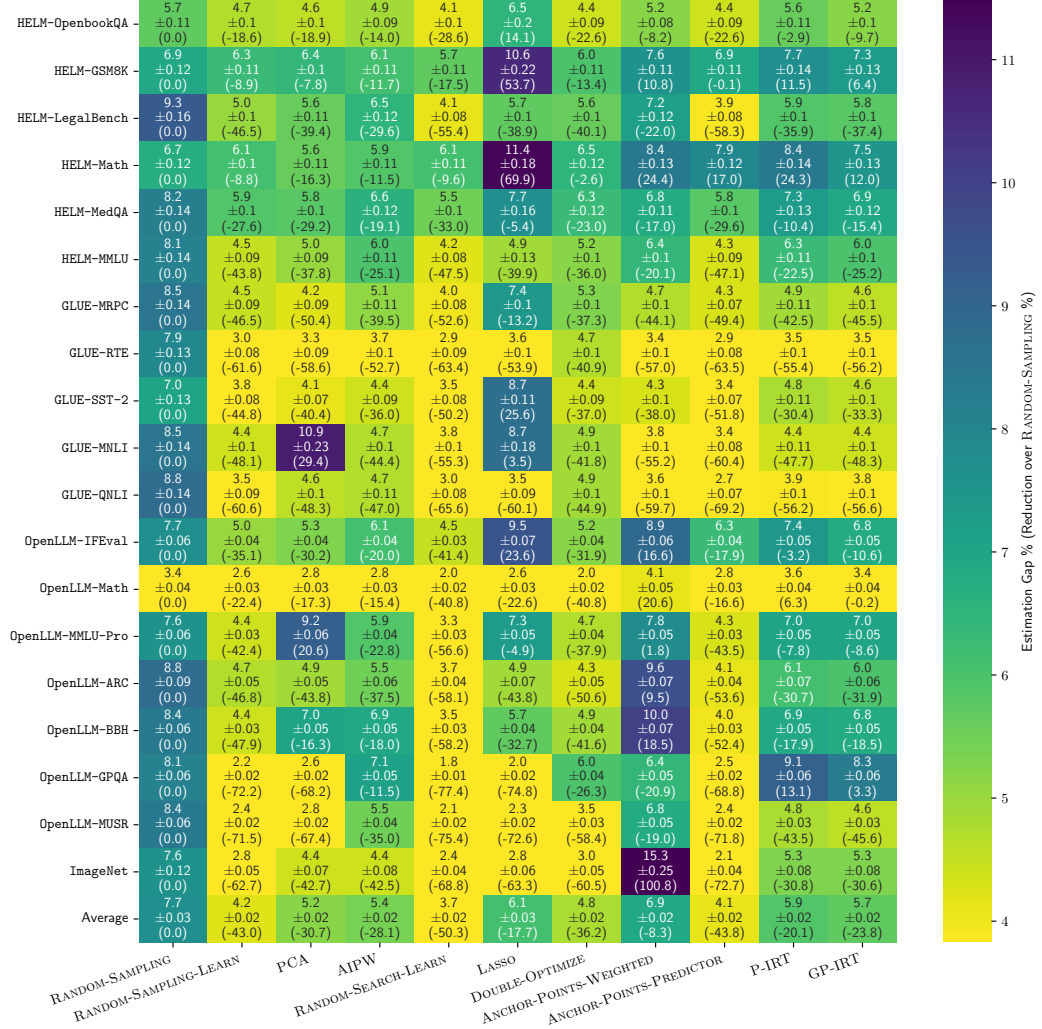


Figure 7: The estimation gaps (\downarrow) for target models (calculated as equation 1) under interpolation model split, where source models are identically distributed with target models. Each target model can only be evaluated on $n = 20$ data points. We also report \pm the standard error of the mean 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.

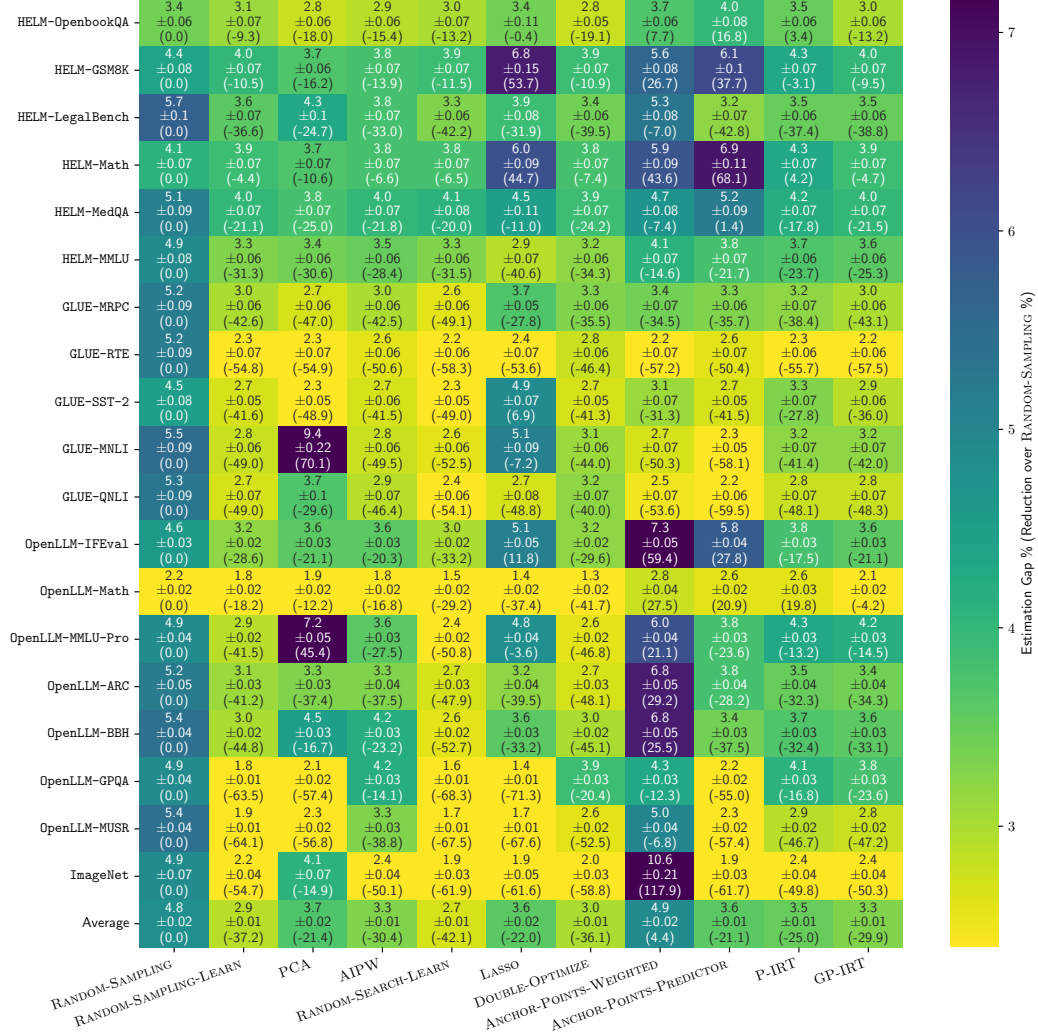


Figure 8: The estimation gaps (\downarrow) for target models (calculated as equation 1) under interpolation model split, where source models are identically distributed with target models. Each target model can only be evaluated on $n = 50$ data points. We also report \pm the standard error of the mean 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.

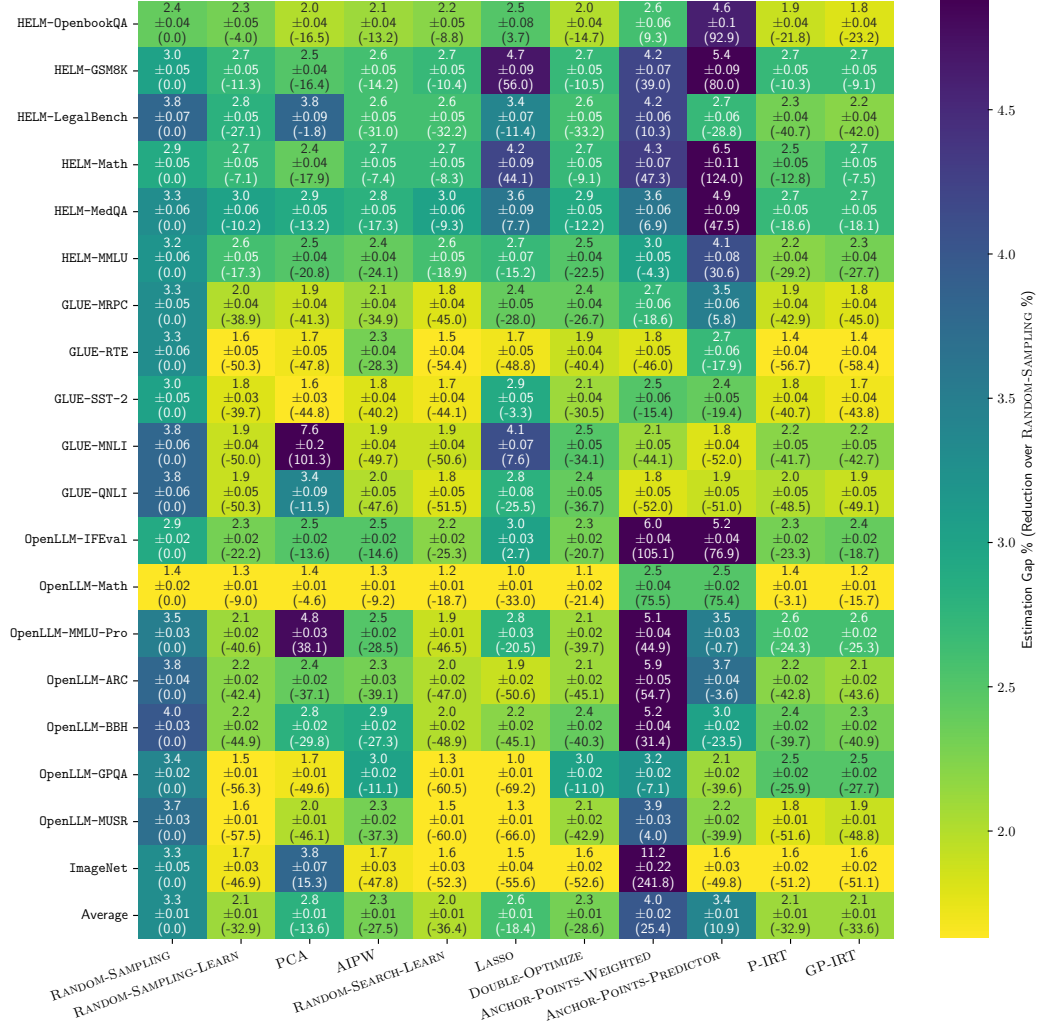


Figure 9: The estimation gaps (\downarrow) for target models (calculated as equation 1) under interpolation model split, where source models are identically distributed with target models. Each target model can only be evaluated on $n = 100$ data points. We also report \pm the standard error of the mean 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.

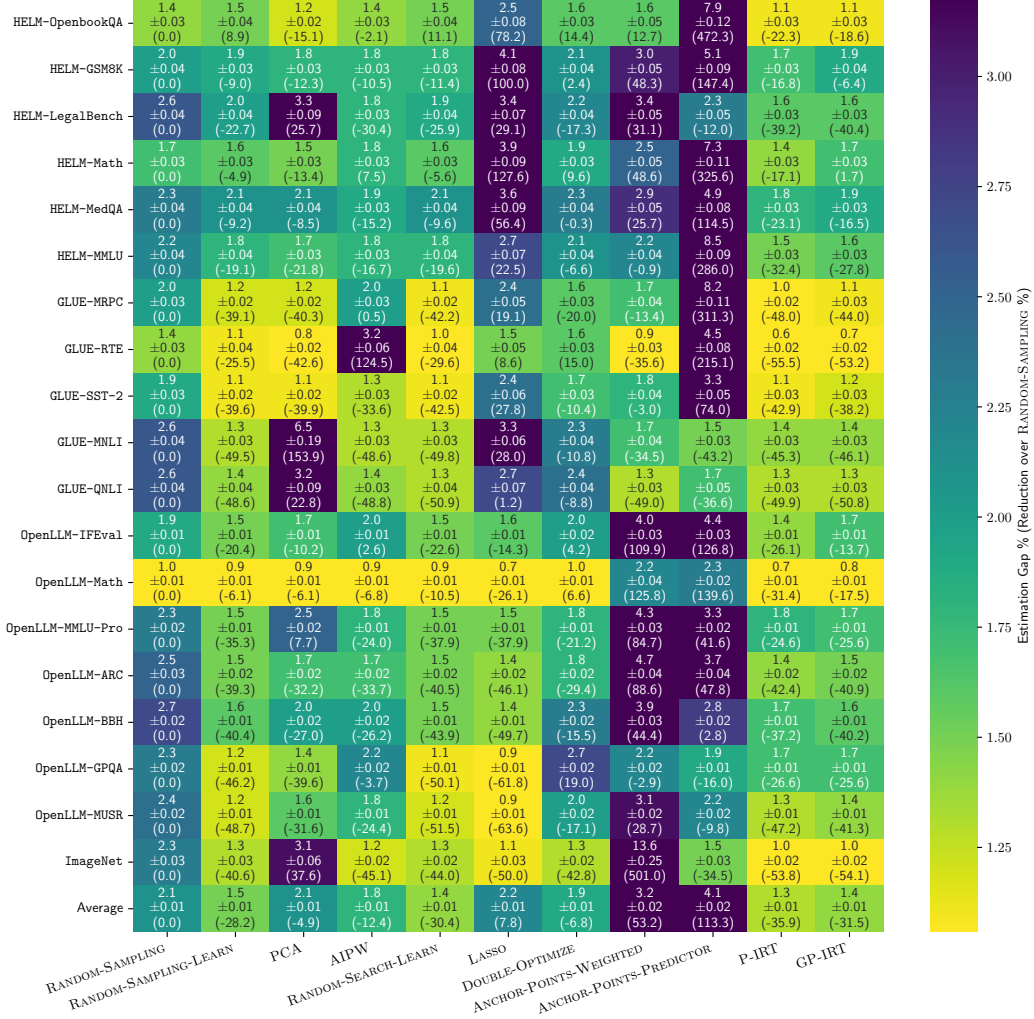


Figure 10: The estimation gaps (↓) for target models (calculated as equation 1) under interpolation model split, where source models are identically distributed with target models. Each target model can only be evaluated on $n = 200$ data points. We also report \pm the standard error of the mean and the estimation gap reduction (↓) 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.

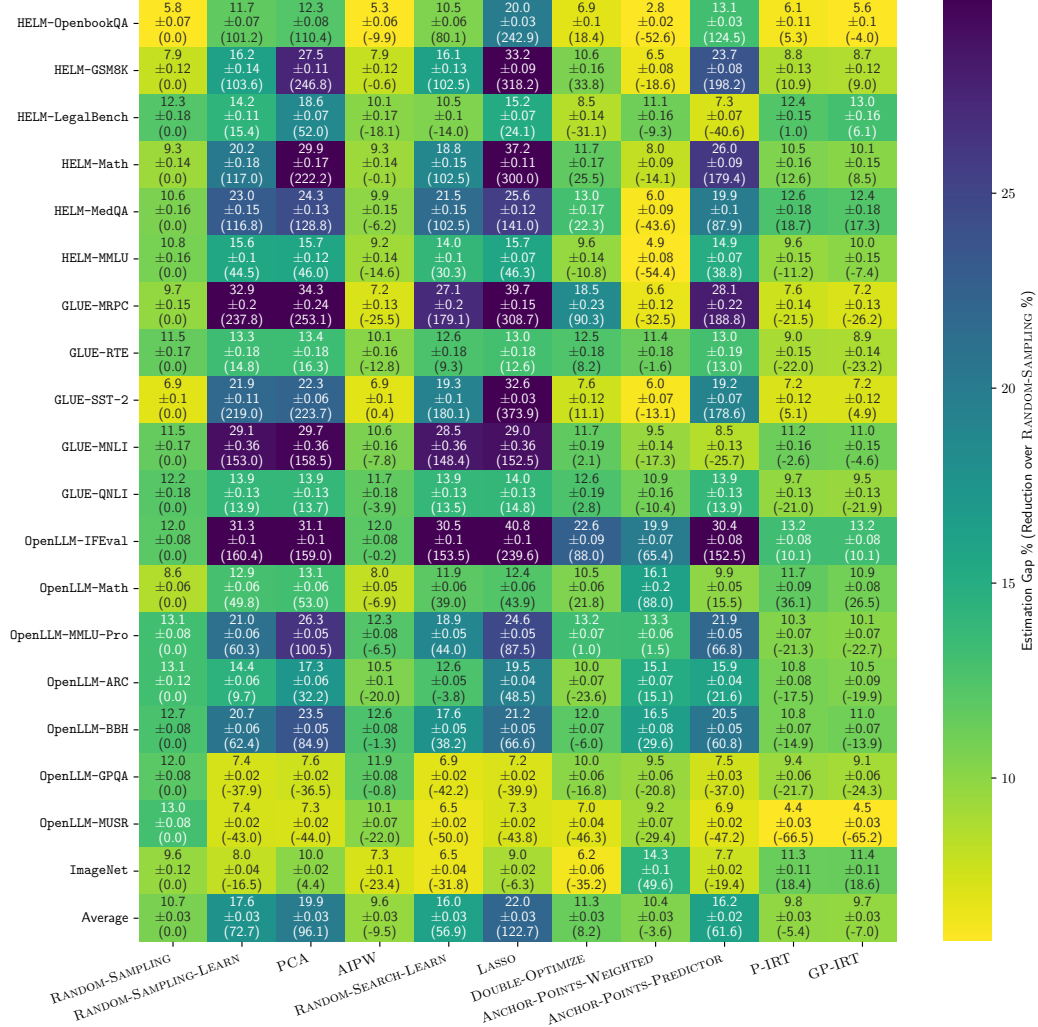


Figure 11: The estimation gaps (↓) for target models (calculated as equation 1) under extrapolation model split, where source models are the lowest-performing 50%, and target models are the top 30% based on average performance over the full benchmark. Each target model can only be evaluated on $n = 10$ data points. We also report \pm the standard error of the mean and the estimation gap reduction (↓) 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.

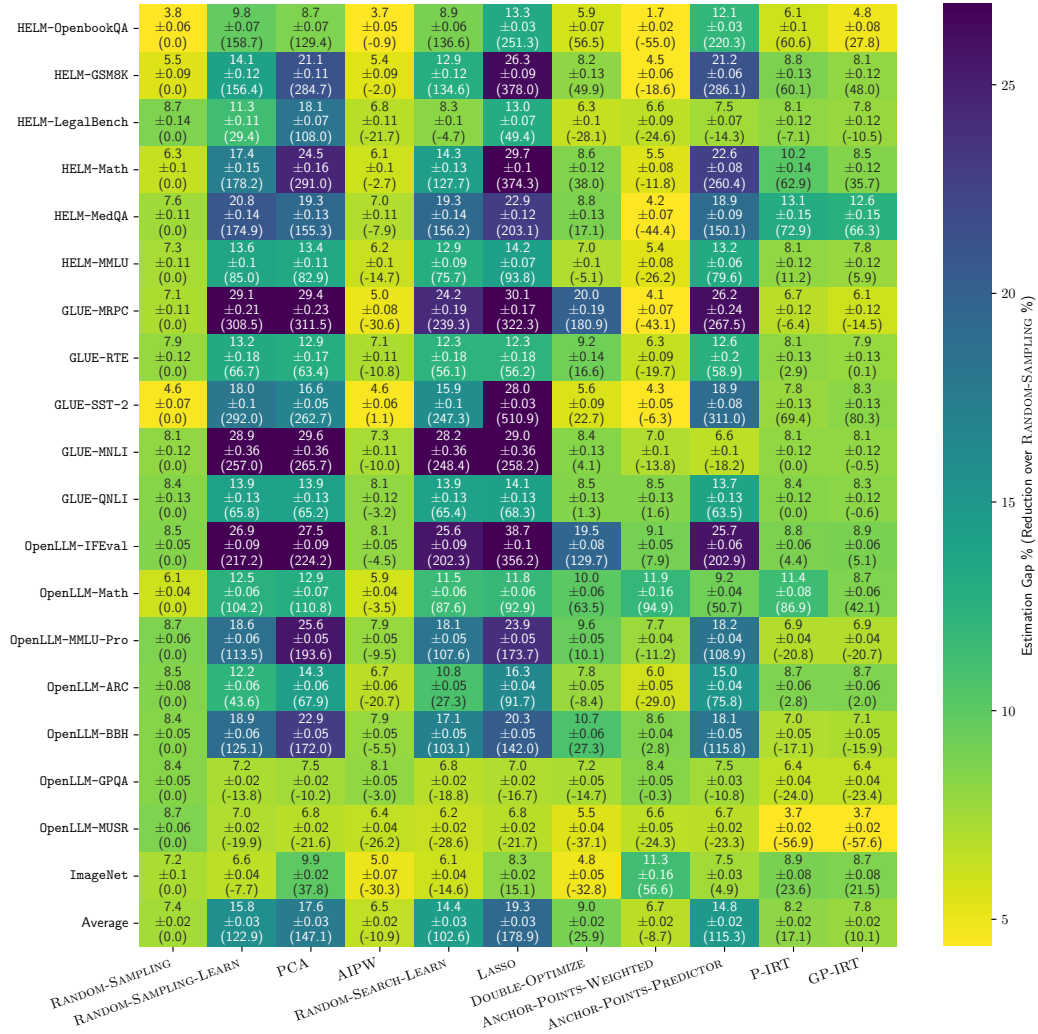


Figure 12: The estimation gaps (↓) for target models (calculated as equation 1) under extrapolation model split, where source models are the lowest-performing 50%, and target models are the top 30% based on average performance over the full benchmark. Each target model can only be evaluated on $n = 20$ data points. We also report \pm the standard error of the mean and the estimation gap reduction (↓) 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.

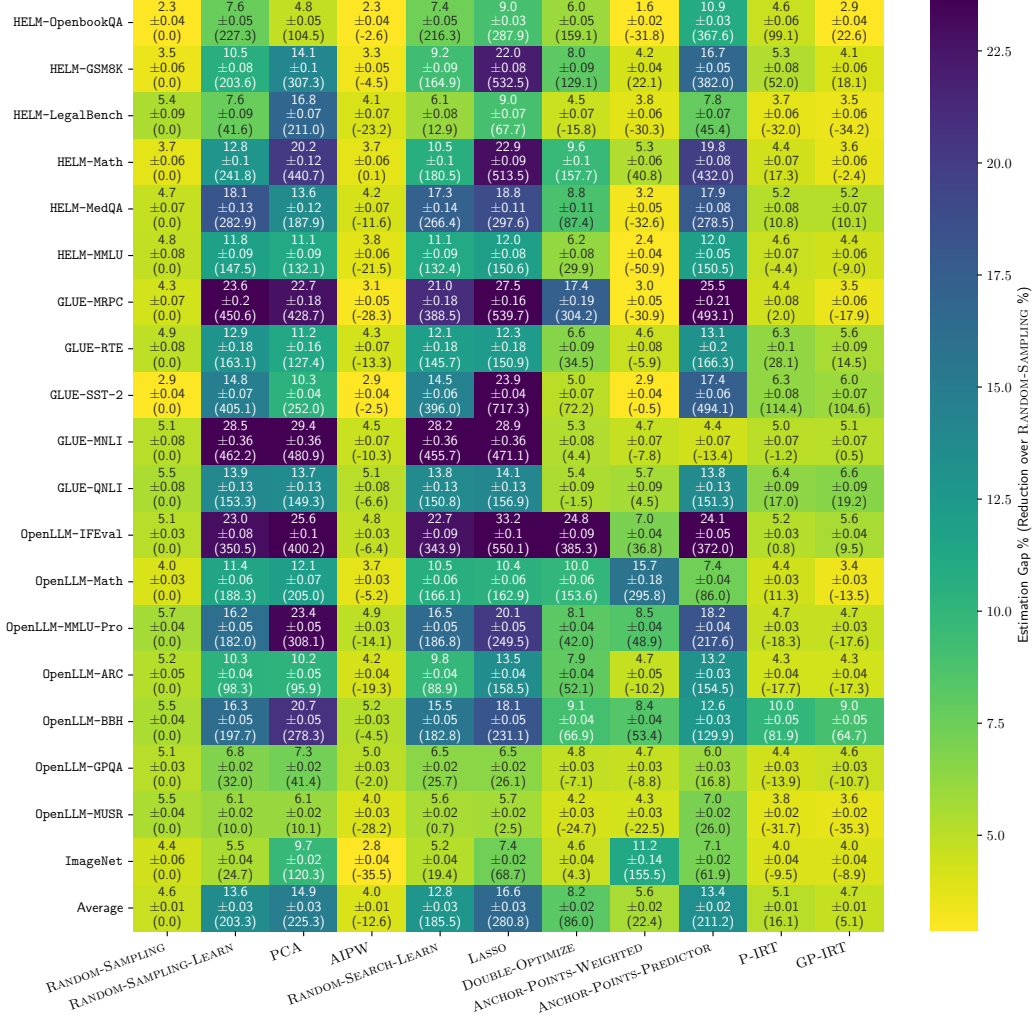


Figure 13: The estimation gaps (↓) for target models (calculated as equation 1) under extrapolation model split, where source models are the lowest-performing 50%, and target models are the top 30% based on average performance over the full benchmark. Each target model can only be evaluated on $n = 50$ data points. We also report \pm the standard error of the mean and the estimation gap reduction (↓) 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.

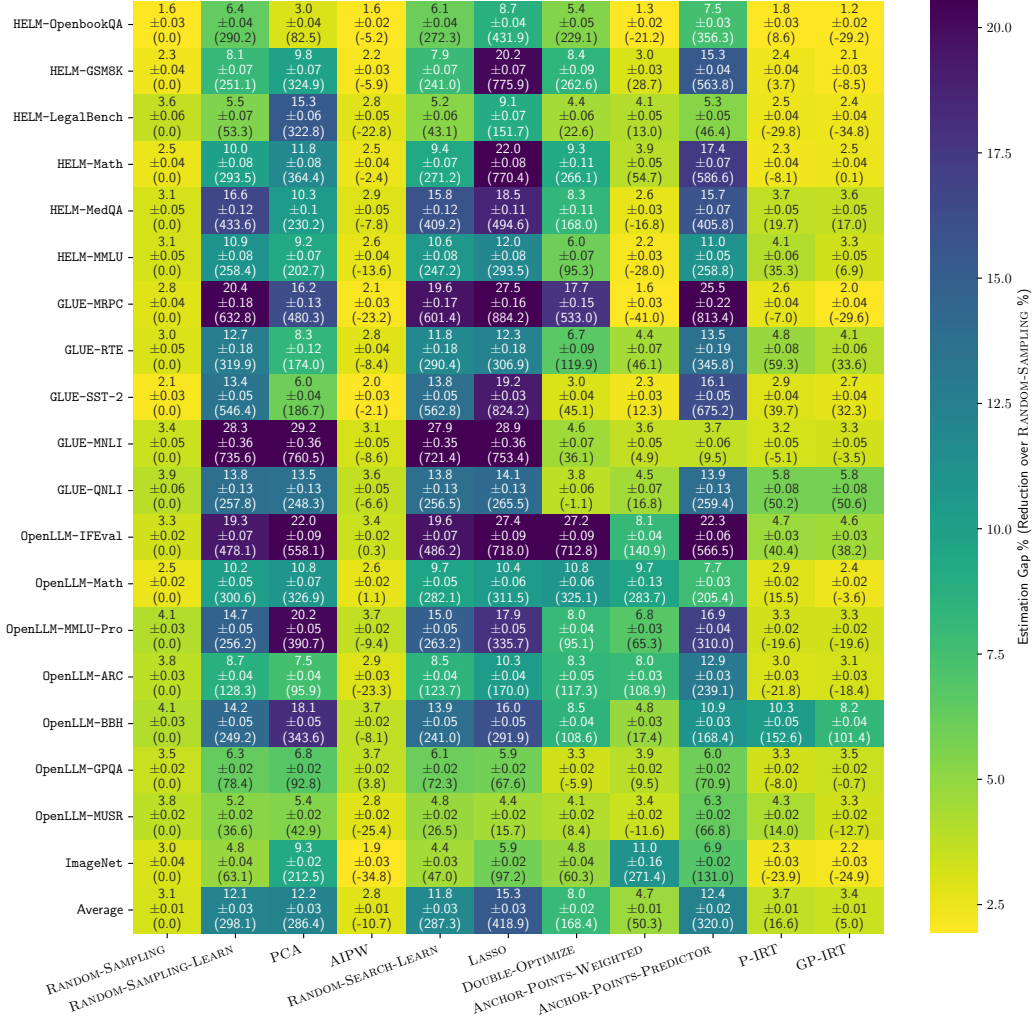


Figure 14: The estimation gaps (↓) for target models (calculated as equation 1) under extrapolation model split, where source models are the lowest-performing 50%, and target models are the top 30% based on average performance over the full benchmark. Each target model can only be evaluated on $n = 100$ data points. We also report \pm the standard error of the mean and the estimation gap reduction (↓) 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.

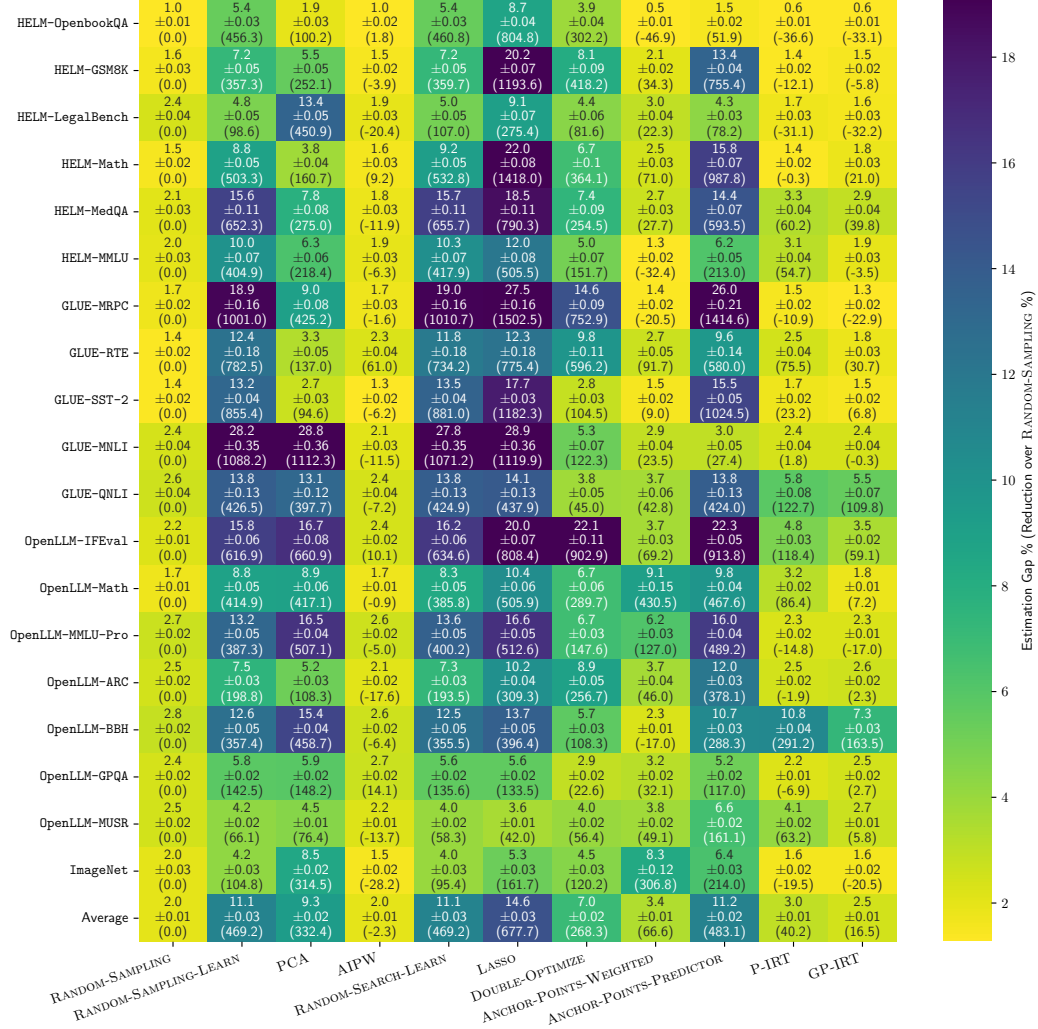


Figure 15: The estimation gaps (\downarrow) for target models (calculated as equation 1) under extrapolation model split, where source models are the lowest-performing 50%, and target models are the top 30% based on average performance over the full benchmark. Each target model can only be evaluated 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.

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 τ (\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

827 D.3 Running time

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

835 D.4 Ranking Preservation

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

841 D.5 Case Studies

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

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

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

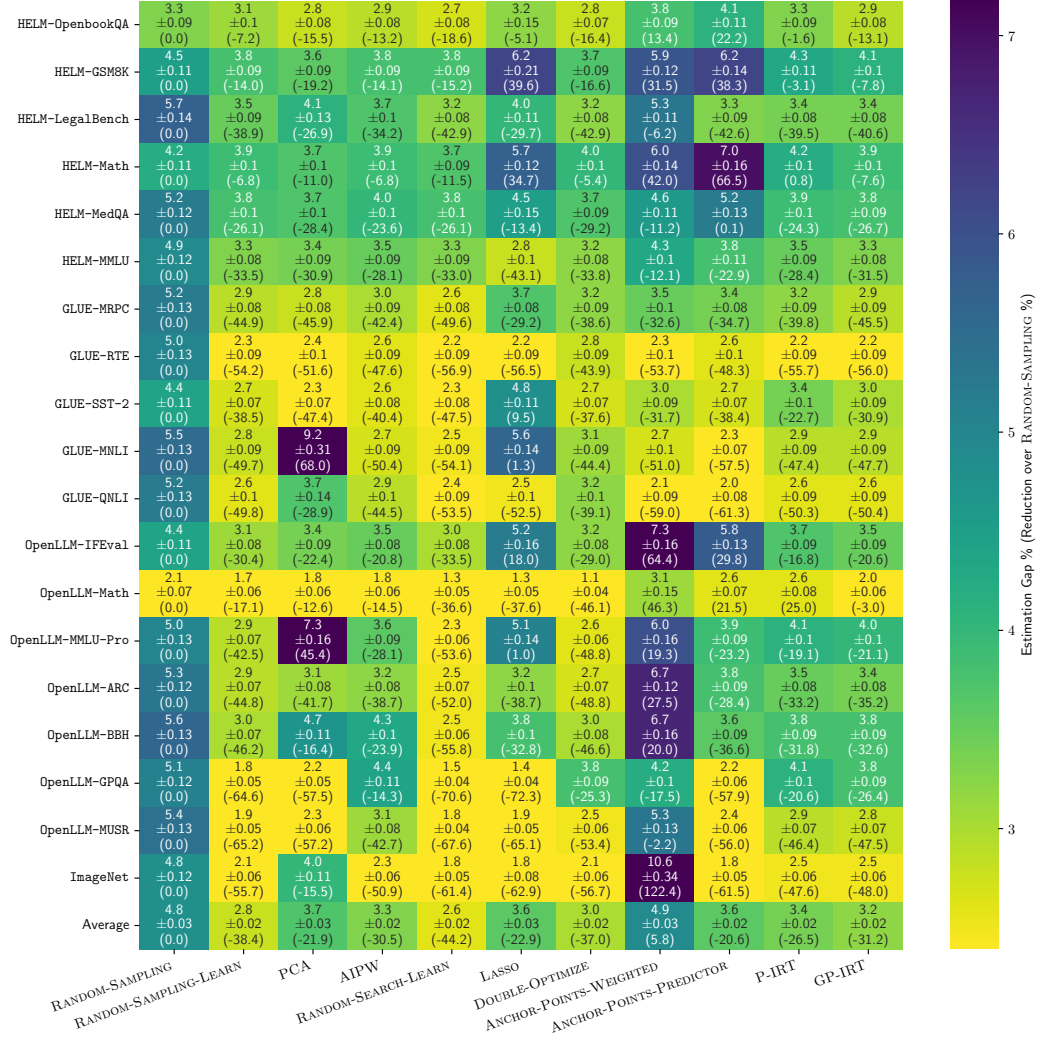


Figure 16: Experiment with fewer source models (randomly selected 10 models as source models) under the interpolation model split. We report the estimation gaps (\downarrow) for target models (calculated as equation 1). We also report \pm the standard error of the mean and the estimation gap reduction (\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.

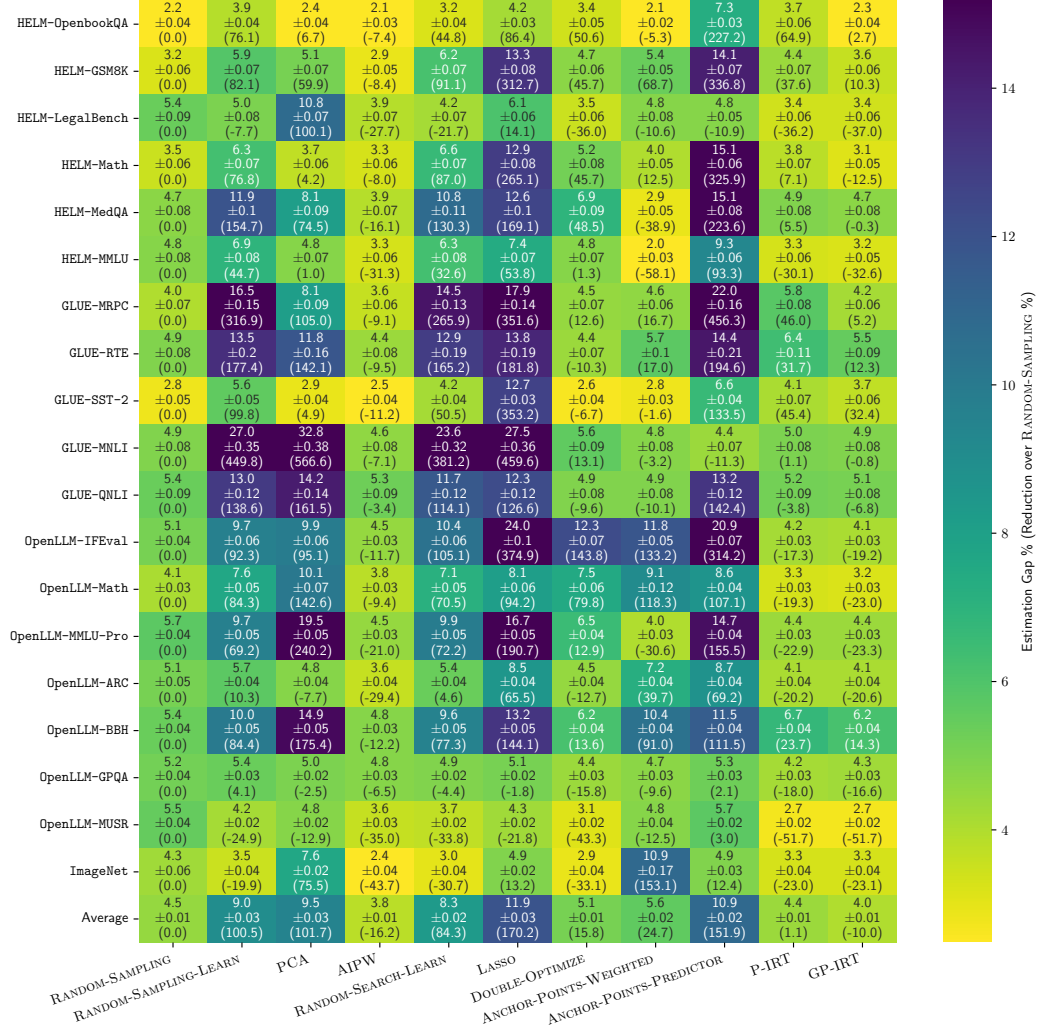


Figure 17: Experiment with the near extrapolation model split by using the top 25% of available models as target models and the remaining bottom 75% models as source models. We report the estimation gaps (\downarrow) for target models (calculated as equation 1). We also report \pm the standard error 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.

858 **E Detailed Conclusion**

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

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

871 **F Broader Impacts and Limitations**

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

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