
SimBA: Simplifying Benchmark Analysis Using Performance Matrices Alone

Nishant Subramani^{*1} Alfredo Gomez^{*1} Mona T. Diab¹

Abstract

Modern language models are evaluated on large benchmarks, which are difficult to make sense of, especially for model selection. Looking at the raw evaluation numbers themselves using a model-centric lens, we propose **SimBA**, a three phase framework to **Simplify Benchmark Analysis**. The three phases of **SimBA** are: *stalk*, where we conduct dataset & model comparisons, *prowl*, where we discover a representative subset, and *pounce*, where we use the representative subset to predict performance on a held-out set of models. Applying **SimBA** to three popular LM benchmarks: HELM, MMLU, and BigBench-Lite reveals that across all three benchmarks, datasets and models relate strongly to one another (*stalk*). We develop an representative set discovery algorithm which covers a benchmark using raw evaluation scores alone. Using our algorithm, we find that with 6.25% (1/16), 1.7% (1/58), and 28.4% (21/74) of the datasets for HELM, MMLU, and BigBenchLite respectively, we achieve coverage levels of at least 95% (*prowl*). Additionally, using just these representative subsets, we can both preserve model ranks and predict performance on a held-out set of models with *near zero* mean-squared error (*pounce*). Taken together, **SimBA** can help model developers improve efficiency during model training and dataset creators validate whether their newly created dataset differs from existing datasets in a benchmark.

1. Introduction

The rapid expansion of language model (LM) benchmarks has resulted in an overabundance of evaluation datasets. However, the relationships among these datasets remain poorly understood. Current evaluation methods primarily

focus on overall model win rates or simple aggregate measures, which fail to provide fine-grained insights into dataset characteristics and model performance trends (Liang et al., 2023b). One approach that the community has taken to mitigate this problem is to look at instance-level predictions and construct coresets, where each individual dataset in a benchmark is subsampled using various heuristics (Rodriguez et al., 2021; Perlitz et al., 2024; Zouhar et al., 2025). This resulting subset is used as a proxy for the entire benchmark. Instance-level sampling has many drawbacks: integration into an already existing evaluation framework is hard, weak statistical signal hinders generalization, and collecting instance-level predictions across many models may be computationally infeasible.

Our work looks to uncover a more structured and simplified understanding of benchmarks through an analysis of the datasets and models directly from the performance matrix *without* collecting any instance-level predictions. Our framework to **Simplify Benchmark Analysis** is called **SimBA** and has three phases:

1. **Stalk**: Analyzing relationships between datasets and measure how models relate to one another across a benchmark.
2. **Prowl**: Discovering a representative subset of datasets from a benchmark that maintains model order.
3. **Pounce**: Predicting model performances using the representative set based on performance patterns.

Using our three phase approach (Figure 1), we analyze HELM (Liang et al., 2023b), MMLU (Hendrycks et al., 2020), and BigBenchLite (Srivastava et al., 2022) and find that both datasets and models correlate well with one another (§2). Motivated by this, we find that:

1. We can identify representative subsets S_{HELM} , S_{MMLU} , and S_{BBL} with just 6.25% (1/16), 1.7% (1/58), and 28.4% (21/74) of datasets respectively that achieve over 95% coverage (§3).
2. Our representative subsets preserve model ranks and can predict performance on a held-out set of models with *near zero* error (§4).

Taken together, our three phase analysis can be used directly by LM practitioners and dataset developers alike to improve efficiency and efficacy.

^{*}Equal contribution ¹Carnegie Mellon University, Language Technologies Institute. Correspondence to: Nishant Subramani <nishant2@cs.cmu.edu>, Alfredo Gomez <alfredo3@cs.cmu.edu>.

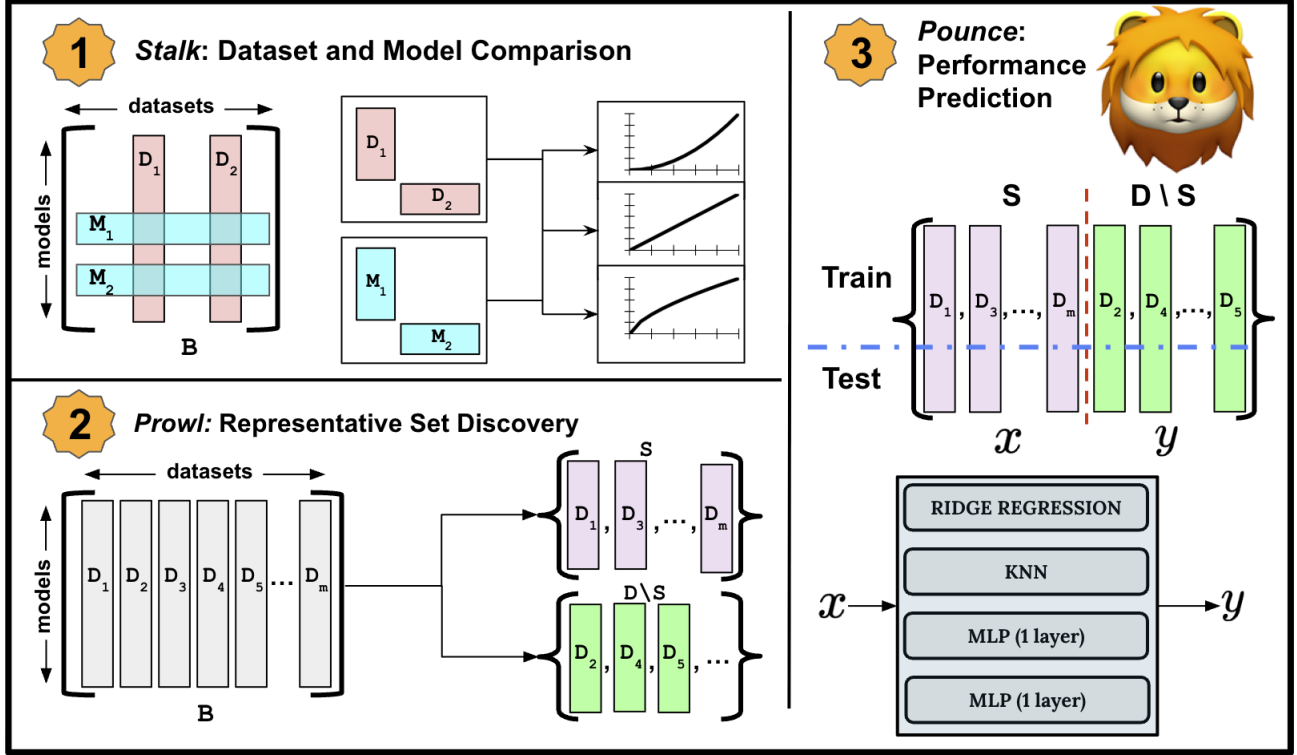


Figure 1: An overview of SimBA, our three phase analysis framework. Its three phases are: *stalk*: dataset & model comparison, *prowl*: representative dataset discovery, and *pounce*: performance prediction.

2. Stalk: Dataset & Model Comparison

A benchmark is represented as a matrix $B \in \mathbb{R}^{m \times d}$, where m is the number of models and d is the number of datasets. B can have missing values. Different datasets evaluate different metrics, which often are scaled differently (e.g. classification vs. generation), so we normalize B . For every dataset D_i with random chance performance x_{random} , we modify every observation x_1, \dots, x_m to be:

$$x_j = \max(0, \frac{x_j - x_{random}}{1 - x_{random}}) \quad (1)$$

Note these new x_j values correspond to percent above random chance. We use this normalization for all analysis.¹ Moreover, all values are further normalized to be within the interval $[0, 1]$.

2.1. Dataset Comparison

Datasets are the essence of a benchmark and the first step in our recipe is to compare datasets. For every pair of datasets D_i and D_j , we compare their performance numbers to identify how they relate to each other using one of four relationships: LINEAR, EXPONENTIAL, POWER-LAW, or

NONE. To measure this, we learn a multi-variate linear regressor f with parameters $W_{reg} \in \mathbb{R}^{m \times 1}$ and $b \in \mathbb{R}$ to optimize the objective:

$$B[:, j] = W_{reg}^T \cdot B[:, i] + b \quad (2)$$

Here, $M[:, i]$ and $M[:, j]$ are the performance numbers across all models for datasets D_i and D_j . We then obtain the R^2 value to quantify the amount of variation explained by the regressor. This naturally works to establish the strength of a linear relationship. To measure an EXPONENTIAL or POWER-LAW relationship, we apply a transformation to the measurements of one dataset before learning the linear regressor.² Since we are predicting one dataset from another, we learn a total of 6 regressors: two per type, one where D_j is predicted from D_i and vice-versa. We choose the one with the largest R^2 value. If that $R^2 < 0.5$ between two datasets, we classify its relationship as NONE.³

²This step resembles applying a nonlinear kernel (e.g., radial-basis function kernel) to an SVM to learn a nonlinear relationship.

³ $R^2 = 0.5$ indicates that only 50% of the variation is explained by the regressor, which we deem is too low to confidently claim a specific relationship. This threshold is a tunable hyperparameter.

¹Srivastava et al. (2022) use the same normalization for Big-Bench for some of their analyses.

2.2. Model Comparison

We compare models across a benchmark in the same way as datasets: for every pair of models M_1, M_2 , we follow the procedure in §2.1 to classify whether the relationship is one of four types: LINEAR, EXPONENTIAL, POWER-LAW, or NONE. Crucially, we only use performance numbers across the benchmark to make this classification.⁴ Identical to dataset comparison, we learn 6 regressors and choose the relationship with the largest R^2 value. We classify the relationship between two models to be NONE if the best regressor results in an $R^2 < 0.5$, just as in §2.1.

3. Prowl: Representative Dataset Discovery

Equipped with an understanding of a benchmark, our next step is to identify a target. More specifically, we want to discover a representative subset of a benchmark S from the performance matrix $B \in \mathbb{R}^{m \times d}$ alone, where m is the number of models and d is the total number of datasets.

3.1. Dataset Similarity

To discover a representative subset, we need a method to determine whether a dataset is redundant.⁵ We compare datasets based on model performance patterns alone using 9 similarity measures: three correlations (PEARSON, SPEARMAN, KENDALL-TAU) and six similarities (COSINE, MANHATTAN, EUCLIDEAN, MINKOWSKI (P=3), WASSERSTEIN, and JENSEN-SHANNON).⁶ We compute similarities between all pairs of datasets, resulting in a matrix $C_{\text{SIM}} \in \mathbb{R}^{d \times d}$, where each entry is a similarity score based on a similarity measure SIM.

3.2. Discovering a Representative Subset

First, we define a proxy metric for the coverage of a candidate representative set S . The PROXY_COVERAGE (δ) of S under a similarity measure SIM is computed as:

$$\delta(S, \text{SIM}) = \frac{\sum_{i \in D} \lambda_i}{|D|} \quad (3)$$

where λ_i is defined as:

$$\lambda_i = \begin{cases} 1, & \text{if } i \in S \\ \max_{j \in S} C_{\text{SIM}}[i, j], & \text{otherwise} \end{cases} \quad (4)$$

⁴We explicitly do not use any information about the model architecture, size, or family.

⁵If there exists a set of datasets $\{D_1, \dots, D_k\}$ that are redundant, we need to keep only one in our representative set.

⁶All similarities are computed as: $\text{SIM} = 1 - \text{DISTANCE}$ or $\text{SIM} = \exp^{-\text{DISTANCE}}$ based on whether DISTANCE is bounded. See Appendix A for more details.

Additionally, we need to define the COVERAGE_GAIN (Ψ) of a set S when a dataset D_i is added under a similarity measure SIM. This is used as a heuristic to measure how much PROXY_COVERAGE (δ) is gained when adding a new dataset D_i to S :

$$\Psi(S, D_i) = \delta(\{S, D_i\}, \text{SIM}) - \delta(S, \text{SIM}) \quad (5)$$

Equipped with these measures, we propose Alg. 1 to discover representative subset. While our implementation supports beam search for representative dataset discovery, empirical evaluation showed that larger beam widths (5, 10, 20) provided no meaningful improvement over the greedy approach across HELM, MMLU, and BigBenchLite.⁷

Algorithm 1 Representative Dataset Discovery

input $C_{\text{SIM}} \in \mathbb{R}^{d \times d}, \gamma, \mathcal{D} = \{D_1, \dots, D_d\}$

output S

- 1: $S \leftarrow \emptyset$
 - 2: **while** $\delta(S) < \gamma \wedge |S| \leq d$ **do**
 - 3: $D^* \leftarrow \arg \max_{D_i \in \mathcal{D} \setminus S} \Psi(S, D_i)$
 - 4: $S \leftarrow S \cup \{D^*\}$
 - 5: **end while**
 - 6: **return** S
-

Baselines: We also evaluate the following baselines as random and simple baselines in dataset selection have been shown to be strong (Diddee & Ippolito, 2025). RANDOM is a baseline where S is populated with a random dataset iteratively without replacement. We average across 1000 random runs in our experiments. GREEDY MINIMUM is a baseline where S is populated with the dataset with the lowest average performance across all models. GREEDY MAXIMUM is a baseline where S is populated with the dataset with the highest average performance across all models. For both greedy baselines, S is also populated iteratively.

3.3. Representative Subset Evaluation

Using Algorithm 1, we discover a representative subset S that has n datasets. To measure how well S covers the original benchmark B , we first find the mean win rate (MWR) of each model as compared to the other models based entirely on S :

$$\text{MWR}(B') = \frac{\sum_{\mu \leq m, d \leq n} \mathbb{I}[B'[\mu, d] > B'[\mu', d]]}{m - 1} \quad (6)$$

Remember that $B' \in \mathbb{R}^{m \times n}$, where m is the number of models and n is the number of datasets in S . B' is matrix formed by collecting all performance numbers for all models for the datasets in S , so $\text{MWR}(B') \in \mathbb{R}^m$. We then

⁷ γ is a PROXY_COVERAGE threshold set by the user.

compute the PEARSON-CORRELATION between $\text{MWR}(B')$ and the mean-win-rate obtained on the full benchmark, $\text{MWR}(B) \in \mathbb{R}^m$, to obtain coverage η :

$$\eta(B') = \text{PEARSON}(\text{MWR}(B), \text{MWR}(B')) \quad (7)$$

Identifying S is a greedy process involving a similarity function SIM . Since S iteratively increases from having one to two to many datasets, we require an algorithm to identify how good a similarity function SIM is in discovering a strong representative subset at every size. Additionally, we also require a way to compare two similarity functions SIM_1 and SIM_2 . Taking inspiration from the receiver-operating characteristic curve and taking the area under it (AUC-ROC) (Marcum, 1960), we develop our own signed AUC measure called *subset coverage AUC* (SCAUC).⁸ To compute SCAUC, we iteratively build S one dataset at a time until we get to the full benchmark ($S = B$). Starting with the first dataset chosen ($|S| = 1$), we compute $\eta(S)$ using equation 7.⁹ We then construct a curve using the d coverage values (η) and compute the signed area under that curve.¹⁰ To compare two similarity functions, we measure their SCAUC values and choose the one with a higher value.

4. Pounce: Performance Prediction

Equipped with the results of the first two phases of our analysis pipeline, we predict the performances directly using a representative subset. Since the representative subset S is just a subset of the full benchmark B , we evaluate whether different regression based approaches can predict $D \setminus S$ (i.e., the subset of D that is not in S) from S alone.

One general approach for matrix prediction is Singular Value Decomposition (SVD) (Candes & Recht, 2009). However, SVD works only on a partially observed matrix, where rows (models) and columns (datasets) have at least one observation. In a realistic setting, we want to use our subset S to predict the other dataset performances $D \setminus S$ on entirely new models, so SVD would not be immediately applicable. As a result, we focus on three regressor types: RIDGE REGRESSION, KNN REGRESSION, and MLP REGRESSION.

Regularized Linear Regression RIDGE REGRESSION uses a linear function with L2 regularization penalty between a model’s performance scores on S and its performance on $D \setminus S$. This regularization helps prevent overfitting when training on small representative subsets (Hastie et al., 2009; Hoerl & Kennard, 1970).

⁸Subramani et al. (2025) also develop a signed AUC measure to measure the tool-calling utility of LLMs.

⁹Note: $\eta(\emptyset)$ is undefined and $\eta(B) = 1$.

¹⁰This is signed area because correlations can be negative.

KNN Regression KNN REGRESSION estimates the performance score y for a dataset by averaging the performance scores of its k nearest neighbors in feature space,

$$y = \frac{1}{k} \sum_{i \in \mathcal{N}_k(X)} y_i \quad (8)$$

where $\mathcal{N}_k(X)$ denotes the set of the k closest datasets to X using a chosen distance metric (e.g., EUCLIDEAN) (Altman, 1992). We use $k=5$ neighbors, or the size of the training set if smaller than 5. KNN does not impose a functional form, allowing it to capture non-linear relationships.

MLP Regression A Multi-Layer Perceptron (MLP) is a feedforward neural network that models non-linear relationships using multiple layers (Rosenblatt, 1958; Ye et al., 2023). We experiment with two MLP architectures a single hidden layer with 12 neurons and a two-layer architecture with 12 neurons in each hidden layer. We include MLP REGRESSION because it can capture complex non-linear relationships between dataset features and model performance.

4.1. Performance Prediction Evaluation

To evaluate how well we can predict performance, mean squared error (MSE) is a natural choice. For a given representative set S , we train a regressor to predict performance on the remaining datasets in the benchmark ($D \setminus S$).¹¹ We compute the MSE of the regressor on the held-out test set on ($D \setminus S$).

In our experiments, we build S sequentially, by greedily adding one dataset at a time according to Algorithm 1. As a result, we can measure MSE at each point for $|S| = 1, \dots, |S| = |D| - 1$.¹² This traces an MSE curve. Using a similar approach to tracing the area under the coverage curve like in §3, we can compute the area under the MSE curve, which we term AUC-MSE.¹³ Note that high values of AUC-MSE indicate high error because this curve is an error curve *not* a performance curve like other AUC curves.

5. Experiments

Benchmarks For our analysis, we look at three benchmarks: HELM (Liang et al., 2023a), MMLU (Hendrycks et al., 2020), and BigBenchLite (Srivastava et al., 2022). We look at the core scenarios of HELM (17 datasets, 29 models), MMLU (58 datasets, 79 models), and BigBenchLite (74 datasets, 45 models), splitting each benchmark into

¹¹ \mathcal{D} is the set of datasets in the benchmark $\mathcal{D} = \{D_1, \dots, D_d\}$.

¹²When $|S| > |D| - 1$, ($|D \setminus S| = 0$).

¹³AUC-MSE is not signed because the minimum error one can get is 0.0, so every point on this curve is in the top right quadrant of a cartesian coordinate plane ($x, y > 0$).

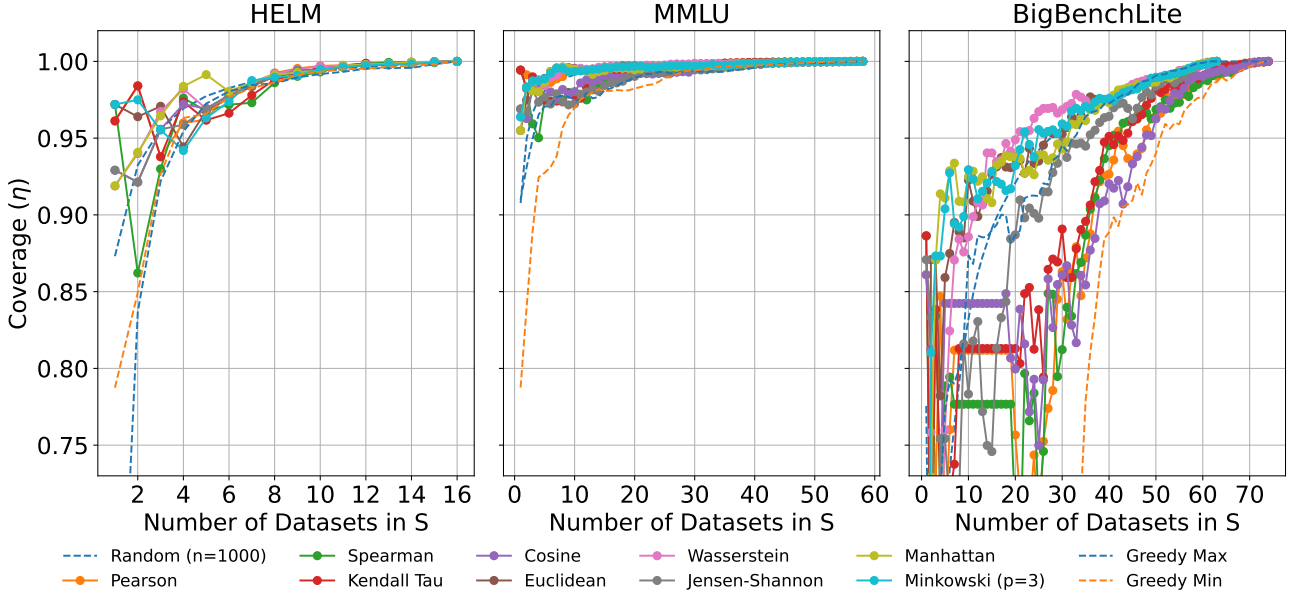


Figure 2: Here, we measure the coverage (η) across HELM, MMLU, and BigBenchLite as our representative subset S grows. We report performance for all three baselines and all nine similarity measures discussed in §3.

training and test sets to validate our analysis. All datasets are included in both splits, but models are separated across training and test sets randomly with 80% of models in training and 20% of the models in test. This means that HELM has 23 models in train and 6 models in test, MMLU has 63 models in train and 16 in test, and BigBenchLite has 36 models in train and 9 in test. For each benchmark, we go through the 3 analysis phases from Figure 1, discuss those results in §6, and expand the analysis to look at robustness in §7.

6. Results

6.1. Stalk: Dataset & Model Comparison

Dataset Comparison: We measure how datasets relate to one another using the methodology in §2.1. Table 1 shows that for HELM, only 5% of the pairs of datasets have a LINEAR relationship, while 37.5% lack any relationship, indicating minor redundancy among datasets.

For MMLU, 4.1% of dataset pairs have a LINEAR relationship and only 7.7% lack a relationship, suggesting that most datasets are predictive of one another through non-linear relationships (EXPONENTIAL: 39.3%, POWER-LAW: 49.0%). For BigBenchLite, 27.9% of dataset pairs have a LINEAR relationship while 61.6% show no relationship at all, indicating that BigBenchLite has the most diverse and independent datasets among the three benchmarks. Taken together, we suspect that finding a small representative dataset would be most difficult for BigBenchLite.

Relationships	HELM	MMLU	BigBenchLite
LINEAR	6	67	754
EXPONENTIAL	25	649	139
POWER-LAW	44	810	143
NONE	45	127	1665
Total	120	1653	2701

Table 1: Dataset comparison results with counts between each pair of datasets and their classifications as LINEAR, EXPONENTIAL, POWER-LAW, and NONE based on the highest R^2 value. See §2.1 for details.

Relationships	HELM	MMLU	BigBenchLite
LINEAR	20	290	75
EXPONENTIAL	72	1258	41
POWER-LAW	294	1302	549
NONE	20	231	325
Total	406	3081	990

Table 2: Model comparison results with counts between each pair of datasets and their classifications as LINEAR, EXPONENTIAL, POWER-LAW, and NONE based on the highest R^2 value. See §2.2 for details.

Similarity Methods	HELM		MMLU		BigBenchLite	
	SCAUC (\uparrow)	$ S^* (\downarrow)$	SCAUC (\uparrow)	$ S^* (\downarrow)$	SCAUC (\uparrow)	$ S^* (\downarrow)$
RANDOM (n = 1000)	0.980	2.5	0.994	2.3	0.928	25.2
GREEDY MINIMUM	0.967	4	0.980	8	0.607	51
GREEDY MAXIMUM	0.957	4	0.988	4	0.933	33
PEARSON	0.984	1	0.997	1	0.886	42
SPEARMAN	0.975	1	0.992	1	0.884	41
KENDALL-TAU	0.983	1	0.994	1	0.899	40
COSINE	0.981	3	0.992	1	0.896	48
MANHATTAN	0.985	3	0.996	1	0.945	32
EUCLIDEAN	0.984	1	0.996	1	0.943	27
MINKOWSKI (P=3)	0.983	1	0.996	1	0.950	22
WASSERSTEIN	0.983	3	0.997	1	0.943	21
JENSEN-SHANNON	0.980	3	0.991	1	0.903	36

Table 3: Performance of our three baselines and seven similarity measures on the identification of a representative subset task for HELM, MMLU, and BigBenchLite. SCAUC is the area under the coverage curves present in Figure 2. S^* is the smallest subset that achieves $\eta(S^*) = 0.95$. **Bold** indicates the best performing system for each metric.

Model Comparison: We measure how models relate to one another via the method in §2.2. Table 2 shows that for HELM, only 4.9% of model pairs have LINEAR relationships, with the majority showing POWER-LAW relationships (72.4%). MMLU demonstrates more complex model relationships with 9.4% LINEAR, 40.8% EXPONENTIAL, and 42.3% POWER-LAW relationships. BigBenchLite shows 7.6% LINEAR relationships and 55.5% POWER-LAW relationships. Notably, all three benchmarks have relatively few model pairs with no discernible relationship (HELM: 4.9%, MMLU: 7.5%, BigBenchLite: 32.8%)

6.2. Prowl: Coverage Analysis

Our goal is to determine the minimum number of datasets needed to achieve a specific coverage level η on a specific benchmark. Since we experiment with nine similarity functions, we want to identify the best similarity measure for this task. To do this, we first compare every pair of datasets D_i, D_j using each of the similarity functions defined in §3. This results in a similarity matrix $C_{\text{SIM}} \in \mathbb{R}^{d \times d}$ for each similarity function.¹⁴ On each similarity matrix C_{SIM} , we apply our coverage algorithm using its respective similarity function and construct S . We measure coverage(η) using equation 7 at each phase of S until $S = \mathcal{D}$. This traces a coverage curve and we measure the area under this coverage curve and report the SCAUC value in Table 3. We also report the size of the *smallest* representative subset S^* that achieves a coverage $\eta \geq 0.95$.

We find that we can achieve coverage levels of at least 95%

¹⁴ C_{SIM} could be an upper (or lower) triangular matrix because our similarity functions are symmetric.

with just 6.25% (1/16), 1.7% (1/58), and 28.4% (21/74) of the datasets for HELM, MMLU, and BigBenchLite respectively. This represents a substantial efficiency gain: particularly for MMLU and HELM, where a single well-chosen dataset can effectively represent nearly the entire benchmark for model ranking purposes. Additionally, Table 3 shows that the choice of similarity measure has varying effects across benchmarks. For HELM and MMLU, most similarity measures perform similarly well, with several achieving the optimal single-dataset representative subset. However, for BigBenchLite, there is more variation in performance, with WASSERSTEIN achieving the best results ($|S^*| = 21$) and several measures like PEARSON and SPEARMAN requiring substantially more datasets ($|S^*| = 42$ and 41 respectively). Finally, using Algorithm 1 with most similarity measures outperforms all baselines on average across all three benchmarks, though the improvement is less pronounced for BigBenchLite due to its more diverse dataset composition. See Table 6 for details on the proportion of times a system outperforms the random baseline across the 1000 random runs.

6.3. Pounce: Performance Prediction

Our goal in this phase is to validate that representative datasets enable accurate performance prediction. Having identified efficient representative subsets in phase II (prowl), we now assess whether performance on these subsets can predict performance on the remaining datasets. We first split our models into training (80%) and test (20%) sets. Using only the training models, we identify representative dataset sets at 80% coverage and train four performance predictors: RIDGE REGRESSION, KNN REGRESSION, and MLP

Regressor	AUC-MSE (\downarrow)		
	HELM	MMLU	BigBenchLite
RIDGE	0.005	0.002	0.002
KNN	0.004	0.002	0.002
MLP (1 layer)	0.081	0.110	0.006
MLP (2 layer)	0.031	0.081	0.006

Table 4: Performance as measured by AUC-MSE across HELM, MMLU, and BigBenchLite. We use MINKOWSKI ($P=3$) as the similarity measure SIM to identify representative subsets S for our four regressors: RIDGE, KNN, MLP (1 layer), and MLP (2 layer). AUC-MSE is the area under the curves in the top row of Figure 3. **Bold** indicates the best performing method for each benchmark.

REGRESSION with one and two layers. We then evaluate these predictors on the test set of models, measuring their ability to predict scores (MSE) for the remaining datasets based solely on performance patterns observed in the representative subset.

As shown in Table 4, both RIDGE REGRESSION and KNN REGRESSION perform exceptionally well across all three benchmarks, achieving *near zero* AUC-MSE values. RIDGE REGRESSION achieves (0.005, 0.002, 0.002) and KNN REGRESSION achieves (0.004, 0.002, 0.002) for HELM, MMLU, and BigBenchLite respectively. Meanwhile, Figure 3 shows that prediction error generally decreases as the representative subset size increases, but with diminishing returns. The MLP models consistently underperform, with the single-layer MLP showing particularly poor results on HELM and MMLU. Additionally, in Figure 3, we find that KNN REGRESSION achieves negligible error with just *one* dataset on both HELM and MMLU.

For MMLU, all four regressors maintain consistently low error rates across different subset sizes, with RIDGE and KNN maintaining the lowest error rates throughout. This complements the analysis presented in §2.1, where MMLU showed the highest proportion of inter-dataset relationships, making it the most predictable benchmark.

7. Analysis

7.1. Robustness Analysis

Our evaluation framework relies on point estimates of the performances of models on individual datasets. As such, robustness to noise is a critical consideration when evaluating different approaches of performance prediction. We evaluate the robustness of performance prediction by perturbing the training set of the benchmark matrix B_{train} as $B'_{train} = B_{train} + \mathcal{N}(\mu, \sigma^2)$ where μ is mean and σ is the noise level parameter. For all noise perturbations we use

	Regressor	AUC-MSE (\downarrow)		
		HELM	MMLU	BigBenchLite
0.05	RIDGE	0.010	0.004	0.004
	KNN	0.017	0.012	0.013
	MLP (1 layer)	0.081	0.079	0.007
0.1	MLP (2 layer)	0.031	0.068	0.007
	RIDGE	0.013	0.007	0.006
	KNN	0.027	0.021	0.013
0.1	MLP (1 layer)	0.087	0.078	0.009
	MLP (2 layer)	0.033	0.067	0.008

Table 5: Performance as measured by AUC-MSE across HELM, MMLU, and BigBenchLite when adding noise. We use MINKOWSKI ($P=3$) as the similarity measure SIM to identify representative subsets S for our four regressors: RIDGE, KNN, MLP (1 layer), and MLP (2 layer). AUC-MSE is the area under the curves in Figure 3. **Bold** indicates the best performing method for each benchmark.

$\mu = 0$ and $\sigma^2 = 0.05$ or $\sigma^2 = 0.1$. This evaluation framework enables us to quantify the stability of our methods under varying noise conditions.

As shown in Figure 3, generally, error rates for all methods increase with greater noise. However, we observe that both RIDGE and KNN REGRESSION achieve low error rates under both noise conditions, across all three benchmarks ($\text{AUC-MSE} < 0.03$). RIDGE consistently maintains AUC-MSE values of 0.013 for all benchmarks and noise conditions. Predictably, the MLP systems are the least robust, but maintain similar error rates to the no noise version. In other words, the MLP systems just are not great at performance prediction.

8. Related Work

The NLP community increasingly evaluates on large benchmarks (Srivastava et al., 2022; Li et al., 2023; Liang et al., 2023b). Some approaches attempt to make evaluation more efficient by doing instance-level reduction (Perlitz et al., 2023; Vivek et al., 2023; Polo et al., 2024). We differ in that we focus on aggregate metrics and do not need access to any instance-level information for efficiency and performance prediction. Ye et al. (2023) also perform performance prediction using simple regressors on BigBench. However, they require features of the model and datasets for accurate prediction, unlike us.

From the field of psychometrics and measurement theory, there exists the idea of convergent and divergent validity (Campbell & Fiske, 1959). Convergent validity suggests that metrics measuring similar underlying constructs should correlate highly with each other. Conversely, discriminant

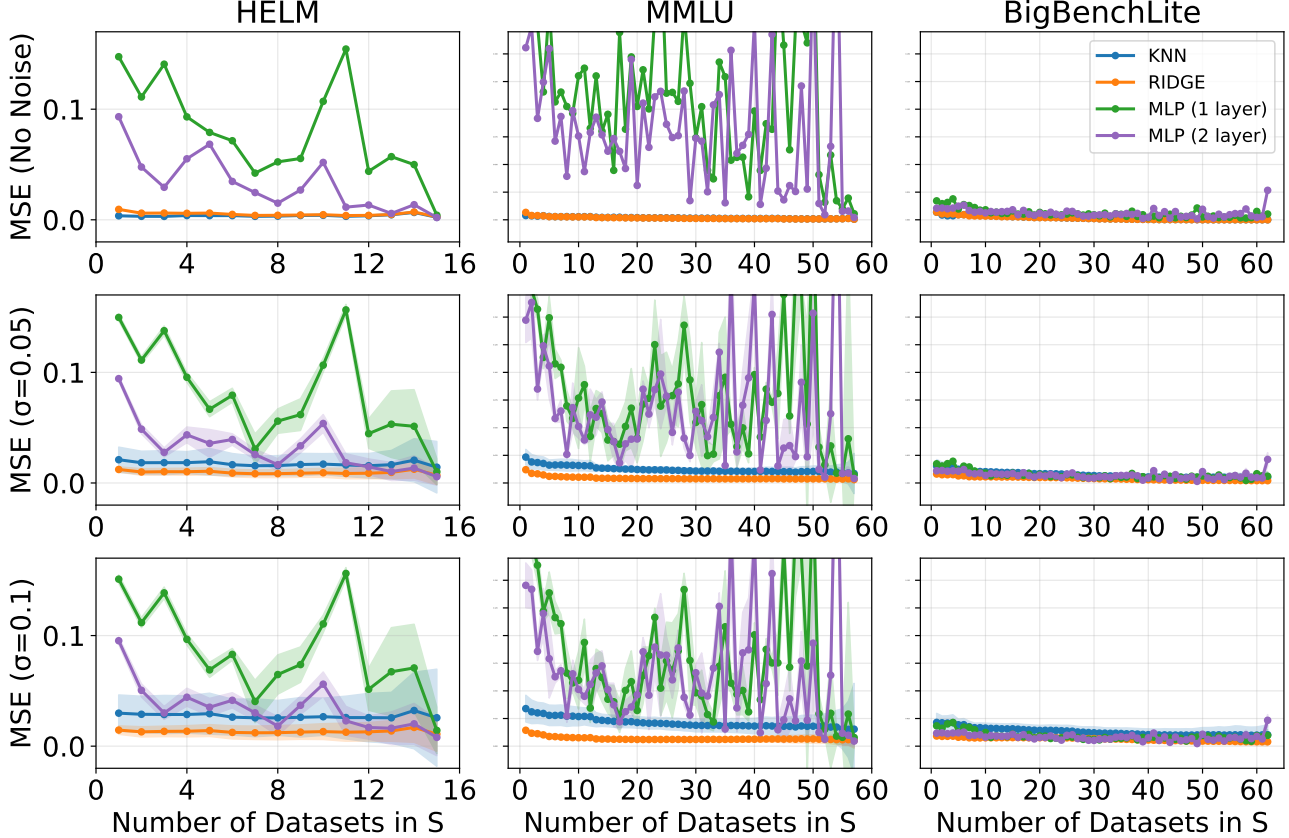


Figure 3: We measure the mean squared error (MSE) on a held-out test set of models of regressors trained on a representative subset S using MINKOWSKI ($P=3$) as SIM. This is measured for HELM, MMLU, and BigBenchLite for the three regressors we experiment with in §4. Additionally, we repeat this experiment after add two magnitudes of noise according to §7.1 ($\sigma = 0.05$ and $\sigma = 0.1$) to the training set of B . Lower scores are better.

validity indicates that metrics capturing fundamentally different aspects should show minimal correlation. Xiao et al. (2023) propose MetricEval, a framework motivated by measurement theory to conceptualize and evaluate the reliable and valid of natural language generation metrics.

9. Conclusion

We propose a three phase approach to **Simplify Benchmark Analysis** called **SimBA**: *stalk* (dataset & model comparison), *prowl* (representative set discovery), and *pounce* (performance prediction). Using our approach, we analyze the HELM, MMLU, and BigBenchLite benchmarks. Our analysis shows that models and datasets alike correlate well with one another (*stalk*). Additionally, using Algorithm 1, we can identify representative subsets with 1, 1, and 21 datasets respectively that achieve a greater than 95% coverage (*prowl*). These representative sets preserve the original model ranks on the benchmark and can be used to predict performance on held-out models with *negligi-*

ble error (*pounce*). Furthermore, **SimBA** can be used by LM practitioners and dataset developers directly to reduce evaluation costs and validate dataset uniqueness.

10. Limitations

Dataset & Model Comparison Our analysis assumes that the relationships between datasets are sufficiently stable over time. As models continue to improve and scale, the nature of these relationships may evolve, potentially requiring periodic reassessment of representative subsets. The approach provides a snapshot analysis based on current model performance matrices, but doesn’t account for how these relationships might change with fundamentally new architectures or training paradigms.

Moreover, the presented analyses requires having a sufficient number of models evaluated on the benchmarks. If the available model lack in diversity in underlying architectures, training data or training methodologies, the identified relationships may not generalize to future models.

Identifying the Representative Dataset Our method is a greedy approach that iteratively chooses datasets such that PROXY_COVERAGE increases. A different, albeit more expensive, approach could exhaustively identify the best combination of datasets using a better search algorithm. We experimented with adapting to beam search, but found no significant improvement, perhaps because PROXY_COVERAGE is correlated, but can be slightly disconnected with COVERAGE depending on the similarity function used.

Performance Prediction Although KNN performs reasonably well, its AUC-MSE values are consistently 2-3 times higher than RIDGE, while the MLP models perform the worst with AUC-MSE values 5-10 times higher than RIDGE, possibly due to overfitting on the limited training data. Models trained with better regularization on more data would have greater stability and be less prone to overfitting.

Overall Risks Our evaluation framework offers a three stage approach to better understand a benchmark. Although representative sets get high coverage, there could be cases when the representative set gets high coverage by chance. In this case, it would be risky to make major decisions that affect users based on a small sample of data.

11. Ethical Considerations

Since **SimBA** does not involve training generative models, the primary ethical concerns center on potential misuse of our framework’s insights and the risk of overconfidence in representative subset evaluations.

Our finding that small representative subsets can achieve high coverage creates opportunities for manipulation. Model developers could strategically evaluate only on datasets where their models perform well, then use **SimBA** to claim coverage over the entire benchmark without actually testing on challenging datasets. This could mislead the research community and downstream users about true model capabilities. Similarly, the choice of similarity measure presents another avenue for selective reporting, as our analysis shows that different similarity functions (PEARSON, MINKOWSKI ($P=3$), WASSERSTEIN, etc.) can yield different representative subsets and coverage results.

To aid in mitigating these concerns, we recommend transparent reporting of representative set selection methodologies, evaluation across multiple correlation methods, and validation on diverse datasets. **SimBA** does not aim to replace comprehensive evaluation, especially for high-stakes deployments. Rather, it serves as a supplementary tool for understanding benchmark structure and improving evaluation efficiency in appropriate contexts.

Acknowledgments

We thank Iz Beltagy, Pradeep Dasigi, Dirk Groeneveld, and Alexis Ross for feedback on very early scoping of this work. Additionally, we thank Harshita Diddee, Athiya Deviyani, and members of MD’s R3Lit lab for helpful discussions and feedback on later versions of the work.

References

- Altman, N. S. An introduction to kernel and nearest-neighbor nonparametric regression. *The American Statistician*, 46(3):175–185, 1992.
- Campbell, D. T. and Fiske, D. W. Convergent and discriminant validation by the multitrait-multimethod matrix. *Psychological bulletin*, 56(2):81, 1959.
- Candes, E. J. and Recht, B. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, 2009.
- Diddee, H. and Ippolito, D. Chasing random: Instruction selection strategies fail to generalize. In Chiruzzo, L., Ritter, A., and Wang, L. (eds.), *Findings of the Association for Computational Linguistics: NAACL 2025*, pp. 1943–1957, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-195-7. URL <https://aclanthology.org/2025.findings-naacl.103/>.
- Hastie, T., Tibshirani, R., and Friedman, J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2nd edition, 2009.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D. X., and Steinhardt, J. Measuring massive multitask language understanding. *ArXiv*, abs/2009.03300, 2020. URL <https://api.semanticscholar.org/CorpusID:221516475>.
- Hoerl, A. E. and Kennard, R. W. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- Li, X., Zhang, T., Dubois, Y., Taori, R., Gulrajani, I., Guestrin, C., Liang, P., and Hashimoto, T. B. Alpaca-eval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval, 5 2023.
- Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., Kumar, A., Newman, B., Yuan, B., Yan, B., Zhang, C., Cosgrove, C., Manning, C. D., R’e, C., Acosta-Navas, D., Hudson, D. A., Zelikman, E., Durmus, E., Ladhak, F., Rong, F., Ren, H., Yao, H., Wang, J., Santhanam,

- K., Orr, L. J., Zheng, L., Yuksekgonul, M., Suzgun, M., Kim, N. S., Guha, N., Chatterji, N. S., Khattab, O., Henderson, P., Huang, Q., Chi, R., Xie, S. M., Santurkar, S., Ganguli, S., Hashimoto, T., Icard, T. F., Zhang, T., Chaudhary, V., Wang, W., Li, X., Mai, Y., Zhang, Y., and Koreeda, Y. Holistic evaluation of language models. *Annals of the New York Academy of Sciences*, 1525:140–146, 2023a. URL <https://api.semanticscholar.org/CorpusID:253553585>.
- Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., Kumar, A., Newman, B., Yuan, B., Yan, B., Zhang, C., Cosgrove, C. A., Manning, C. D., Re, C., Acosta-Navas, D., Hudson, D. A., Zelikman, E., Durmus, E., Ladhak, F., Rong, F., Ren, H., Yao, H., WANG, J., Santhanam, K., Orr, L., Zheng, L., Yuksekgonul, M., Suzgun, M., Kim, N., Guha, N., Chatterji, N. S., Khattab, O., Henderson, P., Huang, Q., Chi, R. A., Xie, S. M., Santurkar, S., Ganguli, S., Hashimoto, T., Icard, T., Zhang, T., Chaudhary, V., Wang, W., Li, X., Mai, Y., Zhang, Y., and Koreeda, Y. Holistic evaluation of language models. *Transactions on Machine Learning Research*, 2023b. ISSN 2835-8856. URL <https://openreview.net/forum?id=iO4LZibEqW>. Featured Certification, Expert Certification.
- Marcum, J. A statistical theory of target detection by pulsed radar. *IRE Transactions on Information Theory*, 6(2): 59–267, 1960.
- Perlit, Y., Bandel, E., Gera, A., Arviv, O., Ein-Dor, L., Shnarch, E., Slonim, N., Shmueli-Scheuer, M., and Choshen, L. Efficient benchmarking (of language models). In *North American Chapter of the Association for Computational Linguistics*, 2023. URL <https://api.semanticscholar.org/CorpusID:261076362>.
- Perlit, Y., Bandel, E., Gera, A., Arviv, O., Ein-Dor, L., Shnarch, E., Slonim, N., Shmueli-Scheuer, M., and Choshen, L. Efficient benchmarking (of language models). In Duh, K., Gomez, H., and Bethard, S. (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 2519–2536, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.139. URL <https://aclanthology.org/2024.naacl-long.139/>.
- Polo, F. M., Weber, L., Choshen, L., Sun, Y., Xu, G., and Yurochkin, M. tinybenchmarks: evaluating llms with fewer examples. *ArXiv*, abs/2402.14992, 2024. URL <https://api.semanticscholar.org/CorpusID:267897919>.
- Rodriguez, P., Barrow, J., Hoyle, A. M., Lalor, J. P., Jia, R., and Boyd-Graber, J. Evaluation examples are not equally informative: How should that change NLP leaderboards? In Zong, C., Xia, F., Li, W., and Navigli, R. (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 4486–4503, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.346. URL <https://aclanthology.org/2021.acl-long.346/>.
- Rosenblatt, F. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65 6:386–408, 1958. URL <https://api.semanticscholar.org/CorpusID:12781225>.
- Srivastava, A., Rastogi, A., Rao, A., Shueb, A. A. M., Abid, A., Fisch, A., Brown, A. R., Santoro, A., Gupta, A., Garriga-Alonso, A., Kluska, A., Lewkowycz, A., Agarwal, A., Power, A., Ray, A., Warstadt, A., Kocurek, A. W., Safaya, A., Tazarv, A., Xiang, A., Parrish, A., Nie, A., Hussain, A., Askell, A., Dsouza, A., Slone, A., Rahane, A., Iyer, A. S., Andreassen, A., Madotto, A., Santilli, A., Stuhlmuller, A., Dai, A. M., La, A., Lampinen, A. K., Zou, A., Jiang, A., Chen, A., Vuong, A., Gupta, A., Gottardi, A., Norelli, A., Venkatesh, A., Gholami-davoodi, A., Tabassum, A., Menezes, A., Kirubakaran, A., Mullokandov, A., Sabharwal, A., Herrick, A., Efrat, A., Erdem, A., Karakacs, A., Roberts, B. R., Loe, B. S., Zoph, B., Bojanowski, B., Ozyurt, B., Hedayatnia, B., Neyshabur, B., Inden, B., Stein, B., Ekmekci, B., Lin, B. Y., Howald, B. S., Orinon, B., Diao, C., Dour, C., Stinson, C., Argueta, C., Ram’irez, C. F., Singh, C., Rathkopf, C., Meng, C., Baral, C., Wu, C., Callison-Burch, C., Waites, C., Voigt, C., Manning, C. D., Potts, C., Ramirez, C., Rivera, C. E., Siro, C., Raffel, C., Ashcraft, C., Garbacea, C., Sileo, D., Garrette, D. H., Hendrycks, D., Kilman, D., Roth, D., Freeman, D., Khashabi, D., Levy, D., Gonz’alez, D. M., Perszyk, D. R., Hernandez, D., Chen, D., Ippolito, D., Gilboa, D., Dohan, D., Drakard, D., Jurgens, D., Datta, D., Ganguli, D., Emelin, D., Kleyko, D., Yuret, D., Chen, D., Tam, D., Hupkes, D., Misra, D., Buzan, D., Mollo, D. C., Yang, D., Lee, D.-H., Schrader, D., Shutova, E., Cubuk, E. D., Segal, E., Hagerman, E., Barnes, E., Donoway, E., Pavlick, E., Rodolà, E., Lam, E., Chu, E., Tang, E., Erdem, E., Chang, E., Chi, E. A., Dyer, E., Jerzak, E. J., Kim, E., Manyasi, E. E., Zheltonozhskii, E., Xia, F., Siar, F., Mart’inez-Plumed, F., Happ’e, F., Chollet, F., Rong, F., Mishra, G., Winata, G. I., de Melo, G., Kruszewski, G., Parascandolo, G., Mariani, G., Wang, G. X., Jaimovitch-L’opez, G., Betz, G., Gur-Ari, G., Galijasevic, H., Kim, H., Rashkin, H., Hajishirzi, H., Mehta, H., Bogar, H., Shevlin, H., Schutze, H., Yakura, H., Zhang, H., Wong, H. M., Ng, I.,

- Noble, I., Jumelet, J., Geissinger, J., Kernion, J., Hilton, J., Lee, J., Fisac, J. F., Simon, J. B., Koppel, J., Zheng, J., Zou, J., Koco'n, J., Thompson, J., Wingfield, J., Kaplan, J., Radom, J., Sohl-Dickstein, J. N., Phang, J., Wei, J., Yosinski, J., Novikova, J., Bosscher, J., Marsh, J., Kim, J., Taal, J., Engel, J., Alabi, J. O., Xu, J., Song, J., Tang, J., Waweru, J. W., Burden, J., Miller, J., Balis, J. U., Batchelder, J., Berant, J., Frohberg, J., Rozen, J., Hernández-Orallo, J., Boudeman, J., Guerr, J., Jones, J., Tenenbaum, J. B., Rule, J. S., Chua, J., Kanclerz, K., Livescu, K., Krauth, K., Gopalakrishnan, K., Ignatyeva, K., Markert, K., Dhole, K. D., Gimpel, K., Omondi, K., Mathewson, K. W., Chiafullo, K., Shkaruta, K., Shridhar, K., McDonell, K., Richardson, K., Reynolds, L., Gao, L., Zhang, L., Dugan, L., Qin, L., Contreras-Ochando, L., Morency, L.-P., Moschella, L., Lam, L., Noble, L., Schmidt, L., He, L., Col'on, L. O., Metz, L., cSenel, L. K., Bosma, M., Sap, M., ter Hoeve, M., Farooqi, M., Faruqui, M., Mazeika, M., Baturan, M., Marelli, M., Maru, M., Quintana, M. J. R., Tolkiehn, M., Giulianelli, M., Lewis, M., Potthast, M., Leavitt, M. L., Hagen, M., Schubert, M., Baitemirova, M., Arnaud, M., McElrath, M., Yee, M. A., Cohen, M., Gu, M., Ivanitskiy, M., Starritt, M., Strube, M., Swkedrowski, M., Bevilacqua, M., Yasunaga, M., Kale, M., Cain, M., Xu, M., Suzgun, M., Walker, M., Tiwari, M., Bansal, M., Aminnaseri, M., Geva, M., Gheini, M., MukundVarma, T., Peng, N., Chi, N. A., Lee, N., Krakover, N. G.-A., Cameron, N., Roberts, N., Doiron, N., Martinez, N., Nangia, N., Deckers, N., Muennighoff, N., Keskar, N. S., Iyer, N., Constant, N., Fiedel, N., Wen, N., Zhang, O., Agha, O., Elbaghdadi, O., Levy, O., Evans, O., Casares, P. A. M., Doshi, P., Fung, P., Liang, P. P., Vicol, P., Alipoormolabashi, P., Liao, P., Liang, P., Chang, P., Eckersley, P., Htut, P. M., Hwang, P., Milkowski, P., Patil, P. S., Pezeshkpour, P., Oli, P., Mei, Q., Lyu, Q., Chen, Q., Banjade, R., Rudolph, R. E., Gabriel, R., Habacker, R., Risco, R., Milliere, R., Garg, R., Barnes, R., Saurous, R. A., Arakawa, R., Raymaekers, R., Frank, R., Sikand, R., Novak, R., Sitelew, R., Bras, R. L., Liu, R., Jacobs, R., Zhang, R., Salakhutdinov, R., Chi, R., Lee, R., Stovall, R., Teehan, R., Yang, R., Singh, S., Mohammad, S., Anand, S., Dillavou, S., Shleifer, S., Wiseman, S., Gruetter, S., Bowman, S. R., Schoenholz, S. S., Han, S., Kwatra, S., Rous, S. A., Ghazarian, S., Ghosh, S., Casey, S., Bischoff, S., Gehrmann, S., Schuster, S., Sadeghi, S., Hamdan, S. S., Zhou, S., Srivastava, S., Shi, S., Singh, S., Asaadi, S., Gu, S. S., Pachchigar, S., Toshniwal, S., Upadhyay, S., Debnath, S., Shakeri, S., Thormeyer, S., Melzi, S., Reddy, S., Makini, S. P., Lee, S.-H., Torene, S. B., Hatwar, S., Dehaene, S., Divic, S., Ermon, S., Biderman, S., Lin, S., Prasad, S., Piantadosi, S. T., Shieber, S. M., Mishnerghi, S., Kiritchenko, S., Mishra, S., Linzen, T., Schuster, T., Li, T., Yu, T., Ali, T., Hashimoto, T., Wu, T.-L., Desbordes, T., Rothschild, T., Phan, T., Wang, T., Nkinyili, T., Schick, T., Kornev, T., Tunduny, T., Gerstenberg, T., Chang, T., Neeraj, T., Khot, T., Shultz, T., Shaham, U., Misra, V., Demberg, V., Nyamai, V., Raunak, V., Ramasesh, V. V., Prabhu, V. U., Padmakumar, V., Srikumar, V., Fedus, W., Saunders, W., Zhang, W., Vossen, W., Ren, X., Tong, X., Zhao, X., Wu, X., Shen, X., Yaghoobzadeh, Y., Lakretz, Y., Song, Y., Bahri, Y., Choi, Y., Yang, Y., Hao, Y., Chen, Y., Belinkov, Y., Hou, Y., Hou, Y., Bai, Y., Seid, Z., Zhao, Z., Wang, Z., Wang, Z. J., Wang, Z., and Wu, Z. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *ArXiv*, abs/2206.04615, 2022. URL <https://api.semanticscholar.org/CorpusID:263625818>.
- Subramani, N., Eisner, J., Sveigliato, J., Van Durme, B., Su, Y., and Thomson, S. MICE for CATs: Model-internal confidence estimation for calibrating agents with tools. In Chiruzzo, L., Ritter, A., and Wang, L. (eds.), *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 12362–12375, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-189-6. URL <https://aclanthology.org/2025.naacl-long.615/>.
- Vivek, R., Ethayarajh, K., Yang, D., and Kiela, D. Anchor points: Benchmarking models with much fewer examples. *ArXiv*, abs/2309.08638, 2023. URL <https://api.semanticscholar.org/CorpusID:262045288>.
- Xiao, Z., Zhang, S., Lai, V., and Liao, Q. V. Evaluating evaluation metrics: A framework for analyzing NLG evaluation metrics using measurement theory. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 10967–10982, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.676. URL <https://aclanthology.org/2023.emnlp-main.676/>.
- Ye, Q., Fu, H., Ren, X., and Jia, R. How predictable are large language model capabilities? a case study on BIG-bench. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 7493–7517, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.503. URL <https://aclanthology.org/2023.findings-emnlp.503/>.
- Zouhar, V., Cui, P., and Sachan, M. How to select datapoints for efficient human evaluation of nlg models? *ArXiv*, abs/2501.18251, 2025. URL <https://api.semanticscholar.org/CorpusID:275993492>.

A. Dataset Similarity Metrics

Here are more details about the dataset similarity metrics used in our analysis. Consider a pair of datasets (D_i, D_j) ; we use seven similarity measures SIM.

PEARSON CORRELATION: Measures linear relationships between dataset performance vectors but is sensitive to outliers and assumes linearity. We compute the PEARSON CORRELATION:

$$\rho_{D_i, D_j} = \frac{\sum (R_{m, D_i} - \bar{R}_{D_i})(R_{m, D_j} - \bar{R}_{D_j})}{\sqrt{\sum (R_{m, D_i} - \bar{R}_{D_i})^2 \sum (R_{m, D_j} - \bar{R}_{D_j})^2}} \quad (9)$$

where R_{m, D_i} is the performance of model m on dataset D_i and \bar{R}_{D_i} is the average performance across all models for dataset D_i ,

SPEARMAN CORRELATION: A ranked correlation that handles non-linear monotonic relationships but is sensitive to small perturbations that flip ranks.

$$\rho_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} \quad (10)$$

where d_i is the rank difference between model performances on datasets D_i and D_j ,

KENDALL-TAU CORRELATION: Another ranked correlation that measures agreement in the orderings of data but also sensitive to small perturbations.

$$\tau = \frac{C - D}{C + D} \quad (11)$$

where C represents concordant model rankings across two datasets, and D represents discordant rankings.

COSINE SIMILARITY: Measures the cosine of the angle between performance vectors, capturing directional similarity regardless of magnitude differences between datasets.

$$\text{COSINE_SIMILARITY}(D_i, D_j) = \frac{B[:, i] B[:, j]}{\|B[:, i]\| \|B[:, j]\|} \quad (12)$$

LP NORM SIMILARITIES: We define a family of similarity measures based on Lp norms with exponential normalization:

$$\text{LP_SIMILARITY}(D_i, D_j) = \exp(-\|B[:, i] - B[:, j]\|_p) \quad (13)$$

We employed various distance-based similarity measures using exponential normalization to capture different aspects of performance similarity between datasets. We specifically use L1 (MANHATTAN), L2 (EUCLIDEAN), and L3 (MINKOWSKI) norms. The exponential normalization ensures all similarities are bounded in $(0, 1]$, with identical performance patterns yielding similarity 1 and increasingly dissimilar patterns approaching 0.

WASSERSTEIN SIMILARITY: Measures the minimum "cost" of transforming one performance distribution into another, capturing both shape and statistical differences. We also exponentially normalize this such that similarities are bounded in $(0, 1]$.

$$\text{WASSERSTEIN_SIMILARITY}(D_i, D_j) = \exp\left(\frac{-W_1(B[:, i], B[:, j])}{\max_{k, l} W_1(B[:, k], B[:, l])}\right) \quad (14)$$

JENSEN-SHANNON SIMILARITY:

$$\text{JENSEN_SHANNON_SIMILARITY}(D_i, D_j) = 1 - \sqrt{\frac{D_{KL}(P_i || M) + D_{KL}(P_j || M)}{2}} \quad (15)$$

Here $M = \frac{1}{2}(P_i + P_j)$ and P_i, P_j are normalized distributions of $B[:, i], B[:, j]$. Jensen-Shannon similarity provides a symmetric measure based on information theory that quantifies dataset distributional differences.

B. Proportions Better than Random

To evaluate the efficacy of each of the non-random system, we measure the proportion of the 1000 random runs that each system outperforms. We measure this across two metrics on all three benchmarks. The first metric, “AUC”, is measured by SCAUC. The second metric, “Max2”, looks at the SCAUC up until a representative dataset S is discovered that achieves at least 95% coverage (S^*). This is compared to the random baseline for the same number of datasets ($|S^*|$ under a similarity function SIM) across all 1000 runs. Note that SCAUC cannot be computed if $|S^*|=1$. In these cases, we consider the first two dataset instead. These results are below in Table 6.

We find that greedy minimum and greedy maximum both underform random. On MMLU, all similarity measures outperform random, but on HELM and BigBenchLite, only about half the systems outperform random on SCAUC on average. Our representative dataset discovery algorithm seems to generally outperform random early on (*i.e.*, the first few datasets) until S^* is discovered, with most systems outperforming random on “Max2.”

Similarity Methods	HELM		MMLU		BigBenchLite	
	AUC (\uparrow)	Max2 (\uparrow)	AUC (\uparrow)	Max2 (\uparrow)	AUC (\uparrow)	Max2 (\uparrow)
RANDOM (baseline)	–	–	–	–	–	–
GREEDY MINIMUM	0.095	0.095	0.000	0.000	0.000	0.000
GREEDY MAXIMUM	0.030	0.028	0.002	0.223	0.521	0.504
PEARSON	0.626	0.970	0.972	1.000	0.068	0.086
SPEARMAN	0.213	0.415	0.108	0.994	0.057	0.072
KENDALL-TAU	0.552	0.970	0.330	0.994	0.141	0.154
COSINE	0.427	0.409	0.124	0.880	0.115	0.132
MANHATTAN (L1)	0.685	0.501	0.825	0.924	0.727	0.849
EUCLIDEAN	0.595	0.944	0.931	0.868	0.692	0.832
MINKOWSKI (L3)	0.538	0.973	0.939	0.965	0.815	0.928
WASSERSTEIN	0.581	0.507	0.962	0.965	0.692	0.791
JENSEN-SHANNON	0.377	0.409	0.060	0.919	0.155	0.303

Table 6: Proportion of methods that perform better than or equal to random baseline across three evaluation metrics. Values closer to 1.0 indicate better performance relative to random selection. **Bold** indicates the best performing method for each metric within each dataset.