TUNEAHEAD: PREDICTING FINE-TUNING PERFOR-MANCE BEFORE TRAINING BEGINS

Anonymous authors
Paper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026027028

029

031

033

034

037

038

039 040

041

042

043

044

046 047

048

051

052

ABSTRACT

Fine-tuning large language models (LLMs) is compute-intensive and error-prone: model performance depends sensitively on data quality and hyperparameter choices, and naïve runs can even degrade model performance. This raises a fundamental question: Can we predict fine-tuning performance before training begins? We present TUNEAHEAD, a lightweight framework for pre-hoc prediction of finetuning performance. TUNEAHEAD encodes each fine-tuning run as a meta-feature vector that combines static dataset descriptors with dynamic probe features from a short simulated run. A gradient-boosting predictor maps these features to performance predictions, while SHAP-based attributions provide interpretable diagnostics that reveal which specific features are driving performance. Across 1,300+ fine-tuning runs on Qwen2.5-7B-Instruct, TUNEAHEAD consistently outperforms strong baselines such as scaling-law predictors, proxy models, and early-stop extrapolation. On a held-out test set of 370 runs, by defining 'success' as exceeding a performance threshold, it accurately predicted 89.4% of successful runs (110/123) and 91.0% of failure runs (225/247), enabling practitioners to proactively avoid costly unsuccessful runs before training begins. This leads to computational savings of 58.4% in total.

1 Introduction

Fine-tuning large language models (LLMs) has become the standard path to domain adaptation, but it remains costly and unpredictable: performance depends sensitively on data quality and hyperparameter choices, and naïve runs can even *degrade* downstream performance in real-world pipelines (Barnett et al., 2024). For practitioners, the key question is not only *how* to fine-tune, but increasingly *whether* a run is worth doing at all.

Predicting fine-tuning success. Consider a healthcare provider deciding whether to fine-tune a general LLM on a clinical dataset: a failed run may consume hundreds of GPU hours yet underperform the base model. Without prediction, practitioners often discover *only after training* that performance falls short, wasting both time and budget (Figure 1(A)). This raises a crucial question: *can we predict fine-tuning performance before training begins?*

Prior art and their limitations. Scaling-law analyses (Kaplan et al., 2020) capture general trends across models and datasets but offer limited insight for a *specific* dataset. Proxy models such as COSMOS (Wang et al., 2025) and short-horizon extrapolations (Kuramoto & Suzuki, 2025) demonstrate that low-cost prediction is feasible. However, they aggregate all features into a single score, conflating the base model's own characteristics and limitations with dataset properties, leaving practitioners unable to answer the crucial question of 'why' a run might fail and thus unable to make targeted improvements to avoid costly failures.

Predicting with TUNEAHEAD. We introduce **TUNEAHEAD**, a diagnostic prediction framework that predicts fine-tuning performance *before* training begins. The core idea is to capture two complementary categories of low-cost features. The first one is *static dataset descriptors*, which are computed from the dataset itself to provide a foundational, model-agnostic assessment of its intrinsic quality (*e.g.*, lexical diversity, data size). The second is *dynamic probe features*, which are extracted from a short probe run (*e.g.*, early loss decay, gradient stability); their unique advantage is capturing the model-specific learnability of the data, revealing early signs of optimization instability or data-model mismatch that are invisible to static analysis alone. A lightweight gradient-boosting predictor maps

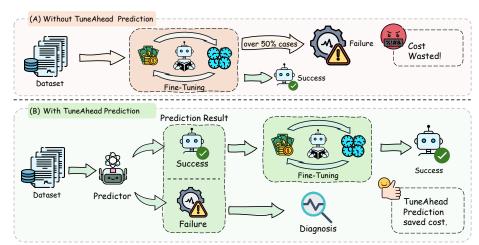


Figure 1: Predicting fine-tuning performance: **(A)** Without **TUNEAHEAD**: failed runs are only identified after training, wasting computational resources and time. **(B)** With **TUNEAHEAD**: low-cost features predict performance in advance, enabling go/no-go decisions and diagnosis for the failure cases.

these features to expected performance, while SHAP-based attribution (Lundberg & Lee, 2017) converts predictions into explanations that reveal which dataset properties matter most. We define a run as 'success' if its predicted score exceeds a predefined practical threshold, and 'failure' otherwise.

With **TUNEAHEAD** (see Figure 1(B)), practitioners can detect unpromising runs *before* training, cutting substantial compute costs while gaining actionable guidance for data refinement. By preventing wasted resources, **TUNEAHEAD** overcomes key limitations of prior work, offering a training-free, interpretable, and dataset-aware approach to reliable fine-tuning.

Figure 2(A) shows that without prediction, all 10 runs (including failures) require ~30 GPU-hours. With TUNEAHEAD (Figure 2(B)), failures are flagged in advance (hatched/blue), so only promising runs are trained, reducing compute to ~18.7 GPU-hours—a 37.4% saving for 4/6 success/failure cases.

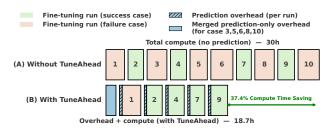


Figure 2: Compute time for 10 runs without (A) vs. with prediction (B).

Contributions. We make the following notable contributions.

- (1) The problem of predicting fine-tuning performance. We cast fine-tuning outcome prediction as a *pre-hoc*, diagnostic meta-learning task. This formulation supports early go/no-go decisions and principled dataset ranking before expensive training is attempted. (Sec. 3)
- (2) The TUNEAHEAD framework. We design a hybrid feature space that combines static dataset descriptors with dynamic probe features, and pairs it with a lightweight predictor and SHAP-based attributions, yielding both accurate predictions and interpretable diagnostics. (Sec. 4)
- (3) Extensive experiments. We conduct over 1,300 fine-tuning runs on Qwen2.5-7B-Instruct (Qwen Team, 2025; Yang et al., 2025), evaluated on Massive Multitask Language Understanding (MMLU) benchmark (Hendrycks et al., 2021). **TUNEAHEAD** consistently outperforms SOTA solutions such as scaling-law predictors, proxy models, and early-stop extrapolation. (Sec. 5)

2 Related Work

Our work sits at the junction of two threads: *LLM performance prediction* and *dataset-level quality assessment*, with connections to fine-tuning dynamics and interpretability. We position **TUNEAHEAD** as a pre-hoc, dataset-level approach that is both *predictive* and *diagnostic*.

LLM performance prediction. Early efforts predict performances by extrapolating short training curves (Domhan et al., 2015), an approach that struggles with modern LLM fine-tuning where dynamics can be non-monotonic and late-emerging. Scaling-law analyses (Kaplan et al., 2020) and recent fine-tuning dynamics tools (*e.g.*, LENSLLM (Zeng et al., 2025)) provide useful macro-level guidance, but they do not explain performance for a *particular* dataset/objective pair. A parallel line uses inexpensive surrogates: proxy heads or smaller models can predict final accuracy with low cost (COSMOS (Wang et al., 2025), PROXYLM (Anugraha et al., 2024), probing-based predictors (Zhu et al., 2022)). However, these methods typically return a single, entangled score that mixes model bias with data properties, offering limited visibility into *why* a run succeeds or fails. Our formulation departs from prior work by combining pre-hoc prediction with explicit, dataset-level diagnosis.

Data quality assessment. A complementary literature scores data quality at the *instance* level—*e.g.*, Data Cartography (Swayamdipta et al., 2020) and Data Shapley (Ghorbani & Zou, 2019)—and more recently curates instruction-tuning data via refinement pipelines (Refine-n-Judge (Cayir et al., 2025)). Others move toward holistic descriptors: Dataset Nutrition Labels (Holland et al., 2018), distributional measures like MAUVE (Pillutla et al., 2021), and generative teaching evaluations (GENTLE (Aoyama et al., 2023)). They improve transparency but generally stop short of *predicting* a dataset's fine-tuning payoff. **TUNEAHEAD** closes this gap by treating each dataset as a meta-instance and tying aggregated *static* descriptors and *early* interaction features to downstream performance.

Fine-tuning dynamics and interpretability. Work on early training signals (*e.g.*, gradient/loss dynamics) informs our choice of low-cost probes (Jastrzebski et al., 2020; Hao et al., 2019). For interpretability, we adopt model-agnostic SHAP attributions (Lundberg & Lee, 2017) to move beyond predicting toward actionable diagnosis at the dataset level.

3 PROBLEM OF PREDICTING FINE-TUNING PERFORMANCE

This section formalizes the task of predicting fine-tuning performance before actual training. Given a dataset–hyperparameter pair, we seek a **low-cost** predictor that approximates the **expensive ground-truth** score that would be obtained after completing full fine-tuning and evaluation. Specifically, let M be a base LLM (e.g., Qwen2.5-7B-Instruct), and A denote a parameter-efficient fine-tuning algorithm (e.g., LoRA). Fine-tuning model M on the dataset–hyperparameter pair (D_i , H_j) produces an adapted model:

$$M'_{i,j} = A(M, D_i, H_j).$$

We then evaluate $M'_{i,j}$ on a downstream benchmark T (e.g., MMLU) to obtain the ground-truth performance score $R_{i,j}$. However, acquiring this score is expensive as it requires a full fine-tuning and evaluation cycle, which motivates the need for prediction.

We thus seek a low-cost prediction function F that consumes a meta-feature vector $V_{i,j}$ describing the dataset and hyperparameter configuration, producing a predicted performance score $P_{i,j}$ that well approximates the ground-truth score:

$$P_{i,j} = F(V_{i,j})$$
 with $P_{i,j} \approx R_{i,j}$.

The prediction function F is trained on a meta-dataset of past fine-tuning experiments drawn from an empirical distribution Dist over pairs (D_i, H_j) . Minimizing an appropriate loss Δ (such as mean squared error), we solve

$$\min_{F} \mathbb{E}_{(D_i, H_j) \sim \text{Dist}} \left[\Delta (F(V_{i,j}), R_{i,j}) \right].$$

This formulation supports several key practical goals: (i) *making go/no-go decisions* before expensive fine-tuning; (ii) *ranking dataset and hyperparameter settings* for resource allocation; and (iii) *diagnosing fine-tuning performance* by linking predictions to dataset and hyperparameter characteristics.

Why pre-hoc failure prediction is feasible. Failed fine-tuning runs often leave clear, low-cost features. Examples include dataset-model mismatch (*e.g.*, high reference perplexity), redundancy or limited diversity (flat or noisy short-horizon progress), and unstable optimization dynamics (volatile gradients and irregular loss decay). Notably, failures are often easier to detect than successes: even a single **strong deficiency can reliably indicate likely failure**, enabling early *rule-out* and data-centric remedies at minimal cost.

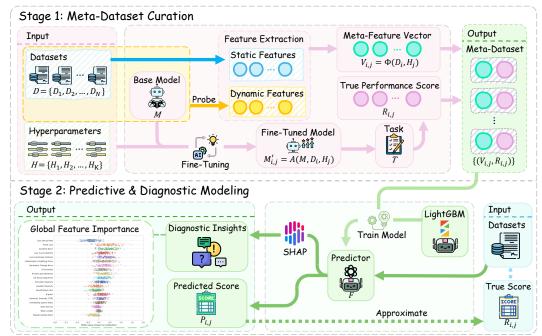


Figure 3: TUNEAHEAD Overview. Stage 1 (Meta-dataset curation) builds meta-feature vectors $V_{i,j}$ by combining static features with dynamic features. Stage 2 (Predictive & Diagnostic Modeling) maps $V_{i,j}$ to performance predictions and uses SHAP for interpretable diagnostics.

4 THE TUNEAHEAD FRAMEWORK

4.1 DESIGN GOALS AND FRAMEWORK OVERVIEW

Design Goals. We identify three key goals that guide our framework design for predicting fine-tuning performance: (G1) Low-cost yet informative features: The meta-features must respect a computational cost budget while still effectively making predictions. (G2) Reliable and generalizable prediction: Predictions must be accurate, well-calibrated, and generalizable across diverse datasets. (G3) Diagnostic interpretability: Predictions must come with human-interpretable attributions that highlight actionable guidance for targeted improvement.

Framework Overview. To satisfy these three goals, TUNEAHEAD is structured into two complementary stages: meta-dataset curation (stage 1) and predictive & diagnostic modeling (stage 2). Figure 3 illustrates the framework overview. Stage 1 constructs a compact meta-feature vector $V_{i,j}$ for each dataset-hyperparameter pair by combining static features that summarize dataset-intrinsic properties and dynamic features (G1) that capture early training behavior via a short, fixed-budget probe run (100 steps). These features are fed into stage 2 for predicting model performance $P_{i,j}$ (G2) along with detailed diagnostic attributions (G3). These explanations identify specific failure modes such as data mismatch, redundancy, or instability, helping enable early rejection of failure runs and guiding focused data-centric improvements.

4.2 STAGE 1: META-DATASET CURATION

The predictive capability of **TuneAhead** heavily depends on how well the meta-feature vector $V_{i,j}$ characterizes each (D_i, H_j) . To balance informativeness with efficiency, $V_{i,j}$ integrates two complementary categories of features: $static\ features$ derived from the dataset itself, providing a model-agnostic prior on dataset quality and $dynamic\ features$ probed from the base model M via a short, fixed-budget run, exposing early signs of instability or mismatch. Together, these features provide a low-cost yet discriminative representation that captures both dataset quality and model learnability, making failure runs easier to detect.

Static features. These are dataset-intrinsic descriptors that require no model training, providing prior insights into dataset properties. We finally selected 14 features that showed robust association with

fine-tuning performance across datasets. The feature selection procedure is detailed in Appendix A.3, while Appendix B provides the full list of features and explains the dataset properties each one captures. The selected static features fall into four categories:

- *Global statistics*: These fundamental quantitative descriptors reflect dataset scale and quality, including dataset size, token length distribution, input–output ratios, and duplication rates.

• Lexical diversity: These metrics measure the richness of the vocabulary. For example, type-token ratio (TTR) quantifies the proportion of unique tokens relative to the total tokens in the dataset:

$$TTR(D_i) = \frac{|\text{unique_tokens}(D_i)|}{|\text{total_tokens}(D_i)|}.$$

• Semantic diversity: These metrics measure the dataset's semantic variability. A high diversity score suggests the dataset covers a wide range of topics and concepts. For example, **pairwise cosine** distance measures the average cosine distance between embeddings e(s):

$$\mathrm{SemDiv}_{\mathrm{PCD}}(D_i) = \frac{2}{|D_i|(|D_i|-1)} \sum_{j < k} \Bigl(1 - \frac{e(s_j) \cdot e(s_k)}{\|e(s_j)\| \|e(s_k)\|} \Bigr).$$

• *Model-Based Complexity*: These metrics quantify dataset complexity w.r.t. the frozen base model \mathcal{M}_0 . For example, **reference perplexity** measures the average perplexity across the dataset:

$$PPL(D_i, \mathcal{M}_0) = \exp\left(-\frac{1}{\sum_{s \in D_i} |s|} \sum_{s \in D_i} \log p_{\mathcal{M}_0}(s)\right).$$

Dynamic features. While static features describe the dataset in isolation, they cannot capture how the base model *actually interacts* with it. We therefore run a standardized **100-step probe** fine-tuning with configuration (D_i, H_j) . This produces low-cost features that approximate early training dynamics. Altogether, we retain 10 dynamic features, which cost less than 5% of a full fine-tuning run. These 10 dynamic features fall into three categories:

• Loss-based indicators: These metrics capture the initial model-data alignment and the subsequent learning progress from the perspective of the optimization objective, like *initial loss* ℓ_0 and the *loss decay rate* which is computed as:

$$LossDecay(D_i, H_i) = slope \left(LinReg(\{(\ell_t, t)\}_{t=0}^{T-1}) \right),$$

where ℓ_t is the training loss at step t.

• *Gradient-based indicators:* These metrics reflect the magnitude, stability, and direction of the learning signals during the early optimization phase like *Gradient Norm*.

• *Model-based indicators*: These metrics assess geometric properties of the loss landscape and changes in the model's internal state, which are often correlated with generalization potential.

4.3 STAGE 2: PREDICTION AND DIAGNOSTIC MODEL

Based on the meta-feature vector extracted in stage 1, we aim to train a lightweight predictor that both well approximates the fine-tuning performance $R_{i,j}$ (G2) while also explaining why a run is likely to succeed or fail (G3). We adopt LightGBM, a gradient-boosted tree model particularly well-suited for heterogeneous, tabular meta-features. As demonstrated in Appendix C, LightGBM achieves accuracy comparable to state-of-the-art alternatives (e.g., SVR) while providing significantly better interpretability and scalability. These properties make it a principled design choice rather than a simple off-the-shelf baseline.

In addition to its powerful prediction capability, LightGBM seamlessly integrates TreeSHAP, a theoretically rigorous method for Shapley value attribution. The SHAP framework decomposes each prediction $P_{i,j}$ into additive contributions from individual meta-features. Unlike black-box or proxy baselines that provide only an opaque overall score, our model delivers transparent attributions. For example, a predicted failure can now be traced back to low lexical diversity (static feature) or unstable gradient norms (dynamic probe feature), which guide targeted improvement.

Table 1: Main results on predicting fine-tuning performances (MMLU). Arrows indicate direction: RMSE \downarrow , $R^2 \uparrow$, $r \uparrow$, Acc@kpp \uparrow . TUNEAHEAD Pareto-dominates all baselines.

Method	RMSE ↓	SE \perp $R^2 \uparrow$		Accuracy $(\pm k pp) \uparrow$		
		10	$r \uparrow$	k=1	k=2	k=3
ProxyLM	2.11	0.98	0.99	40.7	67.9	85.8
Early-Stop Extrapolation	7.43	0.81	0.90	11.2	23.9	32.8
Loss-Rate (Linear)	3.33	0.96	0.98	29.9	50.0	67.5
Reference-PPL (Linear)	6.58	0.85	0.92	8.6	22.0	32.8
TUNEAHEAD-Static-Only	3.50	0.96	0.99	14.9	32.8	49.3
TUNEAHEAD-Dynamic-Only	3.38	0.96	0.99	19.8	35.8	55.6
TUNEAHEAD (Full)	1.47	0.99	0.99	50.0	82.5	95.1

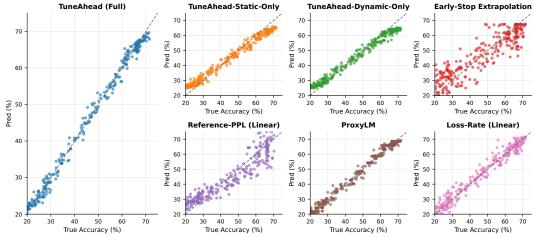


Figure 4: Predicted vs True accuracy across methods. The diagonal line (y=x) indicates a perfect prediction.

5 Experiments

Setup. We validate **TUNEAHEAD** across diverse data and hyperparameter settings on the MMLU task, using Qwen2.5-7B-Instruct, Llama-3-8B-Instruct, and Qwen2-0.5B as base models. All implementation details, data curation protocols, model training procedures, calibration experiments, and baseline implementations are moved to Appendix A. Ground-truth labels are seed-averaged MMLU test accuracy from full LoRA fine-tuning; unless stated, we average over three seeds and evaluate on a held-out test split. For clarity, we label a run as 'success' if its MMLU score exceeds 55%, and 'failure' otherwise. This threshold is chosen only to balance cases for illustration and has no effect on the experiment results, since our predictor outputs a performance score.

Baselines. We compare TUNEAHEAD against (1) literature/practice-inspired baselines and (2) ablation variants: (i) *Early-Stop Extrapolation* — linear extrapolation from the 100-step probe validation loss (Domhan et al., 2015; Adriaensen et al., 2023); (ii) *Loss-Rate (Linear)* — regressors using early loss-decrease rates under standard schedulers (Luo et al., 2025); (iii) *Reference-PPL (Linear)* — regressors on reference perplexity of the fine-tuning data under the frozen base model (Gururangan et al., 2020; Harada et al., 2025); (iv) *ProxyLM* (Anugraha et al., 2024) — regress on proxy models' accuracy (*e.g.*, small LMs) optionally combined with dataset features (*e.g.*, TTR, vocabulary size and average token length); (v) **TUNEAHEAD-***Static-Only* — LightGBM trained on static dataset-intrinsic features only; (vi) **TUNEAHEAD-***Dynamic-Only* — LightGBM trained on dynamic 100-step probe features only. All baselines share the same train/val/test splits and are tuned on the validation set;

Metrics. We report complementary metrics on the test set, all in percentage-point (pp) units for accuracy predictions: RMSE measures absolute error in the model performance prediction; R^2 quantifies explained variance; $Pearson\ r$ measures linear correlation between $P_{i,j}$ and $R_{i,j}$; Acc@kpp is the fraction of predictions with $|P_{i,j}-R_{i,j}| \le k$ percent (we report $k \in \{1,2,3\}$). We additionally provide per-domain breakdowns in Appendix A.7, calibration curves in Appendix A.8, and 95% CIs (bootstrap) with paired permutation tests for significance in Appendix A.9.

Exp-1: Predicting Fine-Tuning Performance and Generalization. In the first set of experiments, we establish the end-to-end predictive strength of **TUNEAHEAD** when both dataset properties and

Table 2: Cross-model generalization results for **TUNEAHEAD** on two additional meta-datasets. Arrows indicate direction: RMSE \downarrow , $R^2 \uparrow$, $r \uparrow$, Acc@kpp \uparrow .

Base Model	RMSE ↓	$R^2 \uparrow$	$r \uparrow$	Accuracy $(\pm k pp) \uparrow$		
	γ	10			k=2	
Llama-3-8B-Instruct Qwen2-0.5B	5.02 3.75		0.93 0.95		55.8 58.4	73.3 74.6

LoRA hyperparameters vary. This set of experiments is designed to validate the accuracy and generalizability of our predictor (G2).

We construct a meta-dataset of over 1,300 complete fine-tuning runs spanning heterogeneous instruction-tuning sources and typical LoRA settings. Unless otherwise specified, all experiments reported in the main text use Qwen2.5-7B-Instruct as the base LLM, with ground truth defined as the *seed-averaged* MMLU test accuracy from full fine-tunes. Curation and training details are provided in Appendix A. To further test generalizability across architectures and scales, we also construct additional meta-datasets on Llama-3-8B and Qwen2-0.5B, experiment results reported in Table 2.

TUNEAHEAD sets a new accuracy bar for this prediction task. As Table 1 shows, on the held-out test set, it cuts RMSE by 30% relative to the strongest non-TUNEAHEAD baseline $(2.11\rightarrow1.47$, ProxyLM) and by 80% relative to Early-Stop Extrapolation $(7.43\rightarrow1.47)$. Tight-tolerance accuracy improves markedly: Acc@1pp +9.3 points (+22.9% rel.), Acc@2pp +14.6 points (+21.5%), and Acc@3pp +9.3 points (+10.8%) over the best baseline, while maintaining near-perfect ranking correlation (Pearson = 0.99). These gains reflect the complementarity of static data descriptors and dynamic 100-step probe signals: either source alone yields $\sim 3.4-3.5$ RMSE, whereas their combination reaches 1.47 (-56% vs. dynamic-only; -58% vs. static-only). Figure 4 visually corroborates these numerical results. The plot for TuneAhead (Full) exhibits the tightest clustering of points along the diagonal, confirming its superior accuracy. In contrast, even the strongest baseline, ProxyLM, shows larger deviations, particularly at the tails of the distribution.

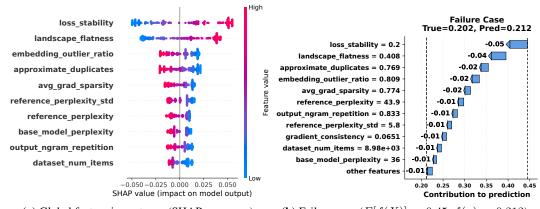
Generalization Across Architectures and Scales. A critical question is whether TUNEAHEAD's predictive power is confined to the Qwen2.5-7B-Instruct model. To proactively assess this, we constructed two additional, smaller-scale meta-datasets using Llama-3-8B (different architecture, 400 runs) and Qwen2-0.5B (smaller scale, 450 runs). As shown in Table 2, TUNEAHEAD continues to capture predictive signal in both cases, achieving $R^2 = 0.86$ on Llama-3-8B and $R^2 = 0.91$ on Qwen2-0.5B, with reasonably high Acc@kpp. Although RMSE is higher, accuracy is lower than in our primary experiments (due to the significantly smaller meta-datasets), the results provide compelling preliminary evidence that the TUNEAHEAD framework is not a model-specific solution.

<u>Summary.</u> (i) Heuristics based solely on early learning curves are *not* sufficient for high-precision prediction; (ii) dataset-intrinsic signals provide complementary, non-redundant information; (iii) their integration in **TUNEAHEAD** yields **low-error**, **tightly calibrated** predictions (95.1% within ± 3 pp), which is the regime practitioners care about for reliable pre-screening without full fine-tuning; (iv) preliminary cross-model experiments suggest that **TUNEAHEAD** generalizes beyond Qwen2.5-7B-Instruct, maintaining predictive accuracy across different model scales and architectures.

Exp-2: Diagnosis. Beyond predictive accuracy, a core contribution of **TUNEAHEAD** is its ability to provide diagnostic insights (G3). In this section, we use TreeSHAP to analyze the trained model and understand the key drivers of fine-tuning success or failure.

Global Feature Importance. The SHAP analysis in Figure 5(a) identifies the most influential features of fine-tuning performance. *loss stability* and *landscape flatness* emerge as the most critical factors, confirming that a stable initial learning phase is paramount. Following these, data quality metrics like *embedding outlier ratio* and *approximate duplicates* are the strongest negative predictors, quantifying the high cost of noisy data. We trained predictors under multiple random seeds and computed TreeSHAP attributions for each. While the exact ranking of features varied slightly across seeds, the same core set of top features consistently emerged, with only minor shifts in their weights.

Case Study: Diagnosing a Predicted Failure Run. While the summary plot reveals global trends, **TUNEAHEAD**'s utility shines in diagnosing individual runs. To demonstrate its practical value, we conduct an in-depth analysis of a representative failure case, which our model correctly predicted.



(a) Global feature importance (SHAP summary). (b) Failure case (E[f(X)] = 0.45, f(x) = 0.212). Figure 5: (a) SHAP summary plot ranking the global importance of meta-features for predicting fine-tuning success; (b) SHAP waterfall plot for a representative failure case (the model correctly predicted low performance).

The SHAP plot Figure 5(b) explains how meta-features contribute to pushing the performance prediction from the average prediction (E[f(X)] = 0.45) down to its final low score (f(x) = 0.212).

The diagnosis reveals a multi-faceted failure, jointly driven by poor learning dynamics and low data quality. On the dynamics side, *loss stability* is extremely low (0.20), whereas successful runs typically range from 0.5 to 0.8; values below 0.3 almost always fail, indicating an unstable trajectory. *Landscape flatness* is also poor (0.408), compared to successful runs clustering around 0.6, with values below 0.5 linked to sharp minima and poor generalization. On the data side, the duplicate rate is 0.769, far above the successful-run median of 0.387. The *embedding outlier ratio* is similarly extreme at 0.809, while successful runs typically stay below 0.5 (median 0.413).

This diagnosis provides a set of targeted, evidence-based prescriptions for the practitioner: (1) **Optimization instability**: Lower the learning rate or adjust optimizer hyperparameters to encourage stable convergence, and (2) **Data quality defects**: Apply semantic de-duplication to reduce redundancy and use embedding-based outlier removal to clean the dataset before fine-tuning.

As a proof of actionability, we applied simple adjustments guided by the diagnosis. Concretely, on the model side, we reduced the learning rate from 3×10^{-5} to 1×10^{-5} and adjusted the optimizer hyperparameters by lowering the AdamW momentum parameters (β_1 from 0.9 to 0.85 and β_2 from 0.999 to 0.98), which encourages more stable convergence. On the data side, we applied SemDeDup (Abbas et al., 2023) for semantic de-duplication, which removes near-duplicates within a cosine threshold of 0.95, and performed embedding-based outlier removal by filtering out samples whose sentence embeddings lie beyond 3 standard deviations from the mean in the representation space. These adjustments improved the run's final MMLU score significantly, from 20.2% to 48.7%.

Exp-3: Ablation Study. To validate our core design choices, we conduct two key ablation studies.

(i) Static vs. Dynamic Feature Ablation. This ablation validates our central hypothesis on the synergy between static and dynamic features. While Exp-1 established the superior performance of the full model, a deeper look at the ablation results in Table 1 and Figure 4 reveals the nature of this synergy. This result demonstrates that a hybrid view can better achieve reliability.

To evaluate this synergy, we partitioned runs into buckets (Table 3) based on whether the subset predictors succeeded (with Acc@2pp metric). The full model demonstrated robust improvements across all scenarios: (i) When only the dynamic predictor succeeded, the full model boosted Acc@2pp from 38.7% to 75.0%; (ii) conversely, when only the static predictor succeeded, it raised Acc@2pp from 37.3% to a perfect 100.0%. (iii) Even when both subset predictors succeeded, the full model still improved Acc@2pp from 40% to 82.1%, (iv) and when both failed, it dramatically mitigated errors, improving Acc@2pp from 12.5% to 70.2%. These rescue effects confirm that static and dynamic features capture different failure modes (data-level flaws vs. model-level instabilities), and their combination is the key to achieving robust prediction.

(ii) Impact of Probe Budget. A key hyperparameter in our framework is the length of the dynamic probe run. To justify our choice of 100 steps, we examine the trade-off between prediction accuracy and the computational cost of the probe. As illustrated in Figure 6, a clear "elbow point" emerges

Table 3: Complementarity analysis on the test set. We report RMSE↓ and Acc@2pp↑; Buckets partition runs by whether static-only and dynamic-only predictions are correct.

(i) Rescued-by-Dynamic			(ii) Rescued-by-Static			
Model	RMSE	Acc@2pp	Model	RMSE	Acc@2pp	
TUNEAHEAD (Full)	1.43	75.0	TUNEAHEAD (Full)	0.69	100.0	
Static-only	4.90	0.0	Static-only	2.92	37.3	
Dynamic-only	3.22	38.7	Dynamic-only	4.60	0.0	
(iii) Both-su	bsets-corr	ect	(iv) Both-subsets-wrong			
Model	RMSE	Acc@2pp	Model	RMSE	Acc@2pp	
TUNEAHEAD (Full)	1.49	82.1	TUNEAHEAD (Full)	1.61	70.2	
Static-only	3.34	36.2	Static-only	4.63	12.5	
Dynamic-only	3.16	40.6	Dynamic-only	5.09	12.5	

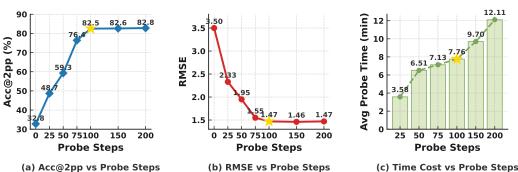


Figure 6: Effect of probe length on prediction accuracy, stability, and time cost. (a) Accuracy at 2pp steadily improves with longer probe runs but exhibits diminishing returns beyond 100 steps. (b) RMSE decreases sharply in the early stage and stabilizes after 100 steps. (c) Average probe time cost grows near-linearly with probe length, with 200 steps requiring about 1.5 times the cost of 100 steps.

around the 100-step mark, indicating a point of diminishing returns. Specifically, increasing the probe budget from 0 to 100 steps yields a significant improvement in accuracy, with RMSE dropping from 3.50 to 1.47 and Acc@2pp rising from 32.8% to 82.5%. However, extending the budget further to 200 steps provides only a negligible gain, with Acc@2pp improving by just 0.3 percentage points to 82.8%. This minimal accuracy improvement comes at a substantial cost, as the average probe time increases linearly, nearly rising from 7.76 minutes to 12.11 minutes. Therefore, we select a 100-step probe as our default configuration, as it offers a robust balance between high predictive power and the low computational overhead mandated by our design goals (G1).

6 CONCLUSION, LIMITATIONS, AND FUTURE WORK

Conclusion. We introduced **TUNEAHEAD**, a framework that predicts fine-tuning performance by combining static and dynamic meta-features. Our experiments show this approach is highly effective, enabling practitioners to avoid costly failed runs and make principled, data-driven decisions.

Limitation. Our primary limitation is generalizability. While **TUNEAHEAD** showed strong performance on several models (see Table 2), these experiments do not yet constitute a comprehensive validation across model families, scales, and task distributions. Specifically, the predictive power of our dynamic features, such as loss stability, is intrinsically tied to the fine-tuning setup. Their importance could shift significantly with different architectures (*e.g.*, dense vs. MoE) or training patterns (*e.g.*, PEFT vs. full-parameter tuning). This implies that while the **TUNEAHEAD** framework is general, a trained predictor instance is likely specific to a model family and requires re-training on a new meta-dataset. Additionally, our compute-saving estimates may vary in other deployments.

Future Work. Our future priorities are twofold. First, we will enhance generalizability not only by expanding our meta-dataset, but more profoundly, by quantifying our feature space's stability across training regimes to develop a meta-learning based predictor that **adapts to new model families with minimal re-training**. The second is to leverage **TUNEAHEAD**'s diagnostics to develop systems that use its feature attributions not just as a passive report, but as an active tool for **automated data quality assessment** and **dataset quality improvements**.

REPRODUCIBILITY STATEMENT

We have made extensive efforts to ensure reproducibility. All dataset curation steps, feature definitions, model training, and evaluation protocols are detailed in Appendix A & B. An anonymized code zip file is provided as supplementary material, containing the full implementation of our framework, including modularized training and evaluation code, SHAP-based diagnostic analysis, and calibration routines. To facilitate reproduction, the meta-dataset we provided illustrates the exact data format and enables full, end-to-end execution of our pipeline. All hyperparameters, experimental settings, and statistical testing procedures are explicitly specified in the main text and the Appendix, ensuring that others can reproduce our reported results.

ETHICS STATEMENT

The TuneAhead framework presented in this work is intended to improve the computational efficiency and resource management of machine learning research and development. We do not foresee any direct negative societal impacts from its application. However, we acknowledge that, like any tool that accelerates model development, its application to sensitive domains or potentially harmful models should be undertaken with caution and ethical oversight. Furthermore, while TuneAhead aims to predict fine-tuning success, the meta-features it relies on (e.g., perplexity, gradient norms) could inadvertently reflect biases present in the base models (e.g., Qwen) or the datasets used for the meta-dataset. Future work could investigate the fairness implications of these predictive features across different demographic and linguistic groups. All research was conducted in compliance with the ICLR Code of Ethics.

LLM USAGE

Our use of Large Language Models (LLMs) in this research was strictly limited to assistance roles, with all core algorithmic ideas, experimental design, and analysis conducted by the human authors.

We used Gemini 2.5 Pro and ChatGPT-5 for English grammar polishing, and asked ChatGPT-5 for suggestions on figure layout and color palettes.

For coding, we consulted ChatGPT-5 and Claude 4.0 for targeted support, such as debugging specific issues and drafting boilerplate scripts (*e.g.*, fine-tuning scaffolds, classic ML baselines). All final code was authored, audited, and verified by the authors.

A portion of the instruction-tuning data in our meta-dataset was synthetically generated using GPT-40 mini and Doubao to increase data diversity. The generation process involved providing seed examples from public datasets and prompting the models to create new, topically related instruction-following pairs. All synthetically generated data underwent a human-led quality control and filtering process before inclusion.

No non-public or sensitive data was used in prompts. LLMs did not contribute any novel algorithmic ideas, hypotheses, or scientific claims presented in this paper. The authors take full responsibility for all content. LLMs do not meet the criteria for authorship.

REFERENCES

Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S. Morcos. Semdedup: Data-efficient learning at web-scale through semantic deduplication, 2023. URL https://arxiv.org/abs/2303.09540.

Steven Adriaensen, Herilalaina Rakotoarison, Samuel Müller, and Frank Hutter. Efficient bayesian learning curve extrapolation using prior-data fitted networks, 2023. URL https://arxiv.org/abs/2310.20447.

Roee Aharoni and Yoav Goldberg. Unsupervised domain clusters in pretrained language models, 2020. URL https://arxiv.org/abs/2004.02105.

- Miltiadis Allamanis, Earl T. Barr, Premkumar Devanbu, and Charles Sutton. A survey of machine learning for big code and naturalness, 2018. URL https://arxiv.org/abs/1709.06182.
 - David Anugraha, Genta Indra Winata, Chenyue Li, Patrick Amadeus Irawan, and En-Shiun Annie Lee. Proxylm: Predicting language model performance on multilingual tasks via proxy models, 2024. URL https://arxiv.org/abs/2406.09334.
 - Tatsuya Aoyama, Shabnam Behzad, Luke Gessler, Lauren Levine, Jessica Lin, Yang Janet Liu, Siyao Peng, Yilun Zhu, and Amir Zeldes. Gentle: A genre-diverse multilayer challenge set for english nlp and linguistic evaluation, 2023. URL https://arxiv.org/abs/2306.01966.
 - Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S. Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, and Simon Lacoste-Julien. A closer look at memorization in deep networks, 2017. URL https://arxiv.org/abs/1706.05394.
 - Amittai Axelrod, Xiaodong He, and Jianfeng Gao. Domain adaptation via pseudo in-domain data selection. In Regina Barzilay and Mark Johnson (eds.), *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pp. 355–362, Edinburgh, Scotland, UK., July 2011. Association for Computational Linguistics. URL https://aclanthology.org/D11-1033/.
 - Scott Barnett, Zac Brannelly, Stefanus Kurniawan, and Sheng Wong. Fine-tuning or fine-failing? debunking performance myths in large language models, 2024. URL https://arxiv.org/abs/2406.11201.
 - Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001. doi: 10.1023/A: 1010933404324.
 - Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwag, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models, 2023. URL https://arxiv.org/abs/2301.13188.
 - Derin Cayir, Renjie Tao, Rashi Rungta, Kai Sun, Sean Chen, Haidar Khan, Minseok Kim, Julia Reinspach, and Yue Liu. Refine-n-judge: Curating high-quality preference chains for llm-fine-tuning, 2025. URL https://arxiv.org/abs/2508.01543.
 - Daniel Cer, Yinfei Yang, Sheng yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. Universal sentence encoder, 2018. URL https://arxiv.org/abs/1803.11175.
 - Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. Electra: Pre-training text encoders as discriminators rather than generators, 2020. URL https://arxiv.org/abs/2003.10555.
 - Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. Documenting large webtext corpora: A case study on the colossal clean crawled corpus, 2021. URL https://arxiv.org/abs/2104.08758.
 - Tobias Domhan, Jost Tobias Springenberg, and Frank Hutter. Speeding up automatic hyperparameter optimization of deep neural networks by extrapolation of learning curves. In *Proceedings of the 24th International Conference on Artificial Intelligence*, IJCAI'15, pp. 3460–3468. AAAI Press, 2015. ISBN 9781577357384.
 - Harris Drucker, Chris J. C. Burges, Linda Kaufman, Alex Smola, and Vladimir Vapnik. Support vector regression machines. In *Advances in Neural Information Processing Systems*, volume 9, pp. 155–161. MIT Press, 1997. URL https://papers.nips.cc/paper/1238-support-vector-regression-machines.
 - Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners, 2021. URL https://arxiv.org/abs/1911.11134.
 - Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks, 2017. URL https://arxiv.org/abs/1703.03400.

- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks, 2019. URL https://arxiv.org/abs/1803.03635.
 - Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning, 2019. URL https://arxiv.org/abs/1904.02868.
 - Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In Geoffrey Gordon, David Dunson, and Miroslav Dudík (eds.), *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pp. 315–323, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR. URL https://proceedings.mlr.press/v15/glorot11a.html.
 - Simon Guiroy, Vikas Verma, and Christopher Pal. Towards understanding generalization in gradient-based meta-learning, 2019. URL https://arxiv.org/abs/1907.07287.
 - Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. Don't stop pretraining: Adapt language models to domains and tasks, 2020. URL https://arxiv.org/abs/2004.10964.
 - Yaru Hao, Li Dong, Furu Wei, and Ke Xu. Visualizing and understanding the effectiveness of bert, 2019. URL https://arxiv.org/abs/1908.05620.
 - Yuto Harada, Yusuke Yamauchi, Yusuke Oda, Yohei Oseki, Yusuke Miyao, and Yu Takagi. Massive supervised fine-tuning experiments reveal how data, layer, and training factors shape llm alignment quality, 2025. URL https://arxiv.org/abs/2506.14681.
 - Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure, 2019. URL https://arxiv.org/abs/1812.04606.
 - Dan Hendrycks, Chris Burns, Santiago Basart, Andy Zou, Mateusz Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations (ICLR)*, 2021.
 - Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md. Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable, empirically, 2017. URL https://arxiv.org/abs/1712.00409.
 - Arthur E. Hoerl and Robert W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970. doi: 10.1080/00401706.1970.10488634.
 - Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022. URL https://arxiv.org/abs/2203.15556.
 - Sarah Holland, Ahmed Hosny, Sarah Newman, Joshua Joseph, and Kasia Chmielinski. The dataset nutrition label: A framework to drive higher data quality standards, 2018. URL https://arxiv.org/abs/1805.03677.
 - Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration, 2020. URL https://arxiv.org/abs/1904.09751.
 - Stanislaw Jastrzebski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun Cho, and Krzysztof Geras. The break-even point on optimization trajectories of deep neural networks, 2020. URL https://arxiv.org/abs/2002.09572.
 - Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, March 2023. ISSN 1557-7341. doi: 10.1145/3571730. URL http://dx.doi.org/10.1145/3571730.

- Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. Exploring the limits of language modeling, 2016. URL https://arxiv.org/abs/1602.02410.
- 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, 2020. URL https://arxiv.org/abs/2001.08361.
- Krishnateja Killamsetty, Durga Sivasubramanian, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. Grad-match: Gradient matching based data subset selection for efficient deep model training, 2021. URL https://arxiv.org/abs/2103.00123.
- Toshiki Kuramoto and Jun Suzuki. Predicting fine-tuned performance on larger datasets before creating them. In Owen Rambow, Leo Wanner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, Steven Schockaert, Kareem Darwish, and Apoorv Agarwal (eds.), *Proceedings of the 31st International Conference on Computational Linguistics: Industry Track*, pp. 204–212, Abu Dhabi, UAE, January 2025. Association for Computational Linguistics. URL https://aclanthology.org/2025.coling-industry.17/.
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better, 2022. URL https://arxiv.org/abs/2107.06499.
- Aodi Li, Liansheng Zhuang, Xiao Long, Minghong Yao, and Shafei Wang. Seeking consistent flat minima for better domain generalization via refining loss landscapes. *arXiv preprint arXiv:2412.13573*, 2024. URL https://arxiv.org/abs/2412.13573.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. A diversity-promoting objective function for neural conversation models, 2016. URL https://arxiv.org/abs/1510.03055.
- Zhuohan Li, Eric Wallace, Sheng Shen, Kevin Lin, Kurt Keutzer, Dan Klein, and Joseph E. Gonzalez. Train large, then compress: Rethinking model size for efficient training and inference of transformers, 2020. URL https://arxiv.org/abs/2002.11794.
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL https://aclanthology.org/W04-1013/.
- Marco Loog and Tom Viering. A survey of learning curves with bad behavior: or how more data need not lead to better performance, 2022. URL https://arxiv.org/abs/2211.14061.
- Scott Lundberg and Su-In Lee. A unified approach to interpreting model predictions, 2017. URL https://arxiv.org/abs/1705.07874.
- Kairong Luo, Haodong Wen, Shengding Hu, Zhenbo Sun, Zhiyuan Liu, Maosong Sun, Kaifeng Lyu, and Wenguang Chen. A multi-power law for loss curve prediction across learning rate schedules, 2025. URL https://arxiv.org/abs/2503.12811.
- Wolfgang Maass. Networks of spiking neurons: The third generation of neural network models. *Neural Networks*, 10(9):1659–1671, 1997. ISSN 0893-6080. doi: https://doi.org/10.1016/S0893-6080(97)00011-7. URL https://www.sciencedirect.com/science/article/pii/S0893608097000117.
- Nikita Moghe, Arnisa Fazla, Chantal Amrhein, Tom Kocmi, Mark Steedman, Alexandra Birch, Rico Sennrich, and Liane Guillou. Machine translation meta evaluation through translation accuracy challenge sets, 2024. URL https://arxiv.org/abs/2401.16313.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization, 2018. URL https://arxiv.org/abs/1808.08745.
- Curtis G. Northcutt, Lu Jiang, and Isaac L. Chuang. Confident learning: Estimating uncertainty in dataset labels, 2022. URL https://arxiv.org/abs/1911.00068.

- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks, 2013. URL https://arxiv.org/abs/1211.5063.
 - Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. Mauve: Measuring the gap between neural text and human text using divergence frontiers, 2021. URL https://arxiv.org/abs/2102.01454.
 - Qwen Team. Qwen2.5-7b-instruct. Hugging Face model card, 2025. URL https://huggingface.co/Qwen/Qwen2.5-7B-Instruct. Accessed: 2025-08-27.
 - Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(1), January 2020. ISSN 1532-4435.
 - Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability, 2017. URL https://arxiv.org/abs/1706.05806.
 - Adriano Rivolli, Luís P. F. Garcia, Carlos Soares, Joaquin Vanschoren, and André C. P. L. F. de Carvalho. Characterizing classification datasets: a study of meta-features for meta-learning, 2019. URL https://arxiv.org/abs/1808.10406.
 - David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, 1986. doi: 10.1038/323533a0.
 - Alex J. Smola and Bernhard Schölkopf. A tutorial on support vector regression. *Statistics and Computing*, 14(3):199–222, 2004. doi: 10.1023/B:STCO.0000035301.49549.88.
 - Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In Sanjoy Dasgupta and David McAllester (eds.), *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pp. 1139–1147, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL https://proceedings.mlr.press/v28/sutskever13.html.
 - Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A. Smith, and Yejin Choi. Dataset cartography: Mapping and diagnosing datasets with training dynamics, 2020. URL https://arxiv.org/abs/2009.10795.
 - Anna Vettoruzzo, Mohamed-Rafik Bouguelia, Joaquin Vanschoren, Thorsteinn Rögnvaldsson, and KC Santosh. Advances and challenges in meta-learning: A technical review, 2023. URL https://arxiv.org/abs/2307.04722.
 - Jiayu Wang, Aws Albarghouthi, and Frederic Sala. Cosmos: Predictable and cost-effective adaptation of llms, 2025. URL https://arxiv.org/abs/2505.01449.
 - Shixian Wen and Laurent Itti. Overcoming catastrophic forgetting problem by weight consolidation and long-term memory, 2018. URL https://arxiv.org/abs/1805.07441.
 - Lei Wu, Zhanxing Zhu, and Weinan E. Towards understanding generalization of deep learning: Perspective of loss landscapes. *arXiv preprint arXiv:1706.10239*, 2017. URL https://arxiv.org/abs/1706.10239.
 - An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Huan Lin, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2025. doi: 10.48550/arXiv.2412.15115. URL https://arxiv.org/abs/2412.15115.
 - Mark Yatskar. A qualitative comparison of CoQA, SQuAD 2.0 and QuAC. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2318–2323, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1241. URL https://aclanthology.org/N19-1241/.

Xinyue Zeng, Haohui Wang, Junhong Lin, Jun Wu, Tyler Cody, and Dawei Zhou. Lensllm: Unveiling fine-tuning dynamics for llm selection, 2025. URL https://arxiv.org/abs/2505.03793.

Zheng Zhao, Shay B. Cohen, and Bonnie Webber. Reducing quantity hallucinations in abstractive

summarization, 2020. URL https://arxiv.org/abs/2009.13312.

Zining Zhu, Soroosh Shahtalebi, and Frank Rudzicz. Predicting fine-tuning performance with probing, 2022. URL https://arxiv.org/abs/2210.07352.

Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences, 2020. URL https://arxiv.org/abs/1909.08593.

A ADDITIONAL EXPERIMENTAL DETAILS

A.1 META-DATASET CURATION

Base Model and Task. All fine-tuning runs use **Qwen-2.5-7B-Instruct**. Downstream performance is measured on **MMLU** for its coverage of knowledge and reasoning across domains.

Dataset Collection $\{D_i\}$. We use public instruction-tuning sources (e.g., Alpaca, Dolly) and both programmatically generated variants and LLM-synthesized examples seeded by these sources, yielding 1,200+ dataset versions. Controlled transformations: (1) **Sub-sampling**: sizes 500–25k; (2) **Domain slicing**: STEM vs. humanities, etc.; (3) **Noise injection**: label noise 0–20%. These produce wide variation in statistical, semantic, and structural properties.

Hyperparameter Space $\{H_j\}$. For each D_i , we sweep LoRA hyperparameters: learning rate $\in \{1e-5, 2e-5, 3e-5\}$, batch size $\in \{8, 16, 32\}$. This covers common practitioner settings and captures realistic data—hyperparameter interactions.

Ground Truth Protocol. For every (D_i, H_j) we run full LoRA fine-tuning to convergence and evaluate on the MMLU test set to obtain $R_{i,j}$. Each $R_{i,j}$ is the mean over 3 seeds to reduce stochasticity.

A.2 META-FEATURE EXTRACTION

We compute a feature vector $V_{i,j}$ per experiment.

Static features. Dataset-intrinsic signals (e.g., semantic diversity, label balance, length stats). Semantic embeddings use all-MiniLM-L6-v2 for efficiency/quality trade-off. Reference Perplexity is computed under the frozen base model (Qwen-2.5-7B-Instruct).

Dynamic probe features. A 100-step probe run per experiment with AdamW and linear scheduler; we log losses/gradients each step and compute early-optimization descriptors.

Preprocessing. We standardize the full feature matrix by z-score before model fitting.

A.3 SHAP-GUIDED FEATURE SELECTION

Setup. We train a LightGBM predictor F_{θ} on a candidate feature set \mathcal{F} and compute SHAP values on the validation set $D_{\text{val}} = \{x_i\}_{i=1}^n$. SHAP ensures an additive decomposition:

$$F_{\theta}(x_i) = \phi_{i,0} + \sum_{f \in \mathcal{F}} \phi_{i,f}.$$

For each feature f, we summarize its global contribution by the mean absolute SHAP and a direction-consistency statistic:

$$s_f := \frac{1}{n} \sum_{i=1}^{n} |\phi_{i,f}|, \qquad \rho_f := \operatorname{SpearmanCorr}(x_{i,f}, \phi_{i,f}),$$

and encode the hypothesized sign by $\eta_f \in \{+1, -1\}$ (whether larger $x_{i,f}$ should increase or decrease the prediction). We use $c_f := \eta_f \, \rho_f$ for signed consistency.

Step 1: Preliminary filtering. Start from 50+ candidates (30 static features and 27 dynamic features). Through small ablations and sanity checks, remove obviously weak or duplicate descriptors to obtain a screened pool \mathcal{F}_1 .

Step 2: SHAP value computation. Fit F_{θ} on \mathcal{F}_1 and compute SHAP values $\{\phi_{i,f}\}$ on D_{val} via TreeExplainer. We inspect beeswarm/bar plots to understand global effects.

Step 3: Global contribution analysis. Retain features that are both strong and consistent with theory. Concretely, keep f if

$$s_f \geq Q_{0.15}(\{s_f\})$$
 and $c_f := \eta_f \, \rho_f \geq 0.20$,

and the Spearman correlation passes significance (p < 0.05). Here $Q_{0.15}$ is the 15th percentile of the empirical $\{s_f\}$.

Step 4: Iterative pruning with CV safeguard. From the retained set, perform backward elimination: at each round remove the weakest candidate (smallest s_f or negative c_f), retrain F_{θ} , and accept the removal only if the cross-validated error does not worsen beyond a fixed tolerance:

$$\Delta RMSE_{cv} = RMSE_{cv}^{new} - RMSE_{cv}^{old} \le \varepsilon, \qquad \varepsilon = 0.01$$

Optionally, use a robust tolerance tied to fold variability:

$$\varepsilon = \min(SE(\Delta), 0.01), \quad SE(\Delta) = \operatorname{sd}(\{\Delta_k\}_{k=1}^K)/\sqrt{K},$$

and additionally require no significant degradation via a paired test (t or Wilcoxon), $p \ge 0.05$. Stop when no feature can be dropped without violating the criterion.

Outcome. This SHAP-guided pipeline yields a compact, non-redundant meta-feature vector (14 static & 10 dynamic in **TUNEAHEAD**), balancing informativeness and stability while aligning with theoretical expectations.

A.4 PREDICTION MODEL TRAINING AND EVALUATION

Training and Evaluation Details. We split the 1,300+ fine-tuning experiments into 28% test, with the remaining 80% further divided into train/validation/calibration subsets (46/14/12%). Unless otherwise noted, the split uses a fixed random seed (=36). The predictor is a Light-GBM gradient-boosted decision tree (GBDT) with the following hyperparameters fixed across all experiments: learning_rate=0.05, num_leaves=4, n_estimators=140, subsample=0.6, colsample_bytree=0.6, min_child_samples=20, and ℓ_2 regularization lambda_12=1.0. We use early stopping with patience of 50 rounds. All models are trained with n_jobs=-1 (multi-threading enabled).

For ablation studies, we define static-only features as dataset descriptors (e.g., length, lexical diversity, perplexity), and dynamic-only features as probe-derived signals (e.g., loss decay, gradient variance). The details of all the features can be looked up in Appendix B.

A.5 BASELINES

Ablation variants. TUNEAHEAD-Static-Only (only static features) and **TUNEAHEAD-**Dynamic-Only (only dynamic probe features).

Practical/literature-inspired baselines. (1) *Early-Stop Extrapolation*: linear extrapolation of the 100-step validation loss (Domhan et al., 2015; Adriaensen et al., 2023); (2) *Loss-Rate Features*: rates of loss decrease under different schedules (Luo et al., 2025); (3) *Reference Perplexity*: reference-PPL as dataset difficulty proxy (Gururangan et al., 2020; Harada et al., 2025); (4) *ProxyLM* (Anugraha et al., 2024): regress on proxy-model scores (e.g., SmolLM-135M/360M, BLOOMZ-560M) optionally with dataset features.

A.6 EVALUATION METRICS AND PROTOCOLS

We evaluate **TUNEAHEAD** with four standard regression and tolerance-based metrics:

• RMSE (percentage points): measures absolute error between predicted and ground-truth performance,

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i,j} (P_{i,j} - R_{i,j})^2}.$$

• R^2 : coefficient of determination, quantifying explained variance in ground-truth performance,

$$R^{2} = 1 - \frac{\sum_{i,j} (R_{i,j} - P_{i,j})^{2}}{\sum_{i,j} (R_{i,j} - R_{i,j})^{2}}.$$

• **Pearson** r: measures linear correlation between predictions and ground-truth,

$$r = \frac{\sum_{i,j} (P_{i,j} - \bar{P})(R_{i,j} - \bar{R})}{\sqrt{\sum_{i,j} (P_{i,j} - \bar{P})^2 \sum_{i,j} (R_{i,j} - \bar{R})^2}}.$$

• Acc@kpp: tolerance-based accuracy, defined as the fraction of predictions within k percentage points of ground-truth,

Acc@kpp =
$$\frac{1}{N} \sum_{i,j} \mathbf{1}(|P_{i,j} - R_{i,j}| \le k), \quad k \in \{1, 2, 3\}.$$

For all metrics, we average over three random seeds when obtaining ground-truth labels.

A.7 PER-DOMAIN BREAKDOWN OF PREDICTION PERFORMANCE

To complement the aggregate results in Table 1, we report per-domain breakdowns of prediction accuracy on MMLU. The 57 subjects of MMLU are grouped into seven finer categories (STEM, Social Sciences, Humanities, Arts & Culture, Health & Medicine, Business & Professional, and Other/General Knowledge). This analysis verifies whether TUNEAHEAD consistently generalizes across heterogeneous domains.

Table A.1: Per-domain breakdown of prediction performance on MMLU. Domains are grouped into finer categories for clarity. Metrics include RMSE (pp), Pearson correlation r, and accuracy within ± 2 pp tolerance.

Domain	$RMSE \downarrow$	$r\uparrow$	Acc@2pp↑
STEM (Math, CS, Physics, Bio)	1.62	0.98	80.5
Social Sciences (Econ, Psych, Soc)	1.45	0.99	84.2
Humanities (History, Philosophy, Law)	1.53	0.99	83.1
Arts & Culture (Literature, Linguistics, Art)	1.59	0.98	81.9
Health & Medicine (Clinical, Nutrition, Public Health)	1.48	0.99	82.7
Business & Professional (Mgmt, Exams)	1.51	0.99	82.3
Other / General Knowledge	1.41	0.99	82.7
Overall	1.47	0.99	82.5

As shown in Table A.1, the predictive performance of **TUNEAHEAD** is consistent across domains. RMSE varies only within ± 0.15 across groups, and Pearson correlations remain above 0.98 throughout. Accuracy within $\pm 2pp$ is also stable, ranging between 80–84%. This robustness indicates that the framework does not disproportionately benefit from or fail on particular subject categories, reinforcing its general applicability across diverse fine-tuning scenarios.

A.8 CALIBRATION ANALYSIS

In Section 5, we emphasized that **TUNEAHEAD**'s predictions are not only accurate in terms of correlation with fine-tuning outcomes, but also well-calibrated in absolute values. This property is crucial for practical use cases, because a well-calibrated predictor allows practitioners to directly interpret a predicted score as an approximate probability of fine-tuning success, enabling threshold-based go/no-go decisions without ad hoc post-processing.

Figure A.1 shows the calibration curves of **TUNEAHEAD** compared with representative baselines. Each point corresponds to a bin of predicted scores, with the x-axis showing the mean predicted value and the y-axis showing the empirical success rate within that bin. The dashed line represents perfect calibration. We observe that several baselines (e.g., *Reference-PPL*, *Early-Stop*) systematically deviate from the diagonal, indicating a tendency to either over-estimate or under-estimate success probability. By contrast, **TUNEAHEAD**'s curve (red circles) remains consistently close to the diagonal across the full range, demonstrating superior calibration. This confirms that **TUNEAHEAD** is not only a strong ranker (as shown by the high correlations in Table 1), but also a reliable probability estimator, making its scores directly actionable for practitioners.

A.9 CONFIDENCE INTERVALS OF MAIN RESULTS

In Section 5, Table 1 reported the aggregate RMSE and accuracy of **TUNEAHEAD** and representative baselines. While those results demonstrated clear performance gains, it is also important to examine the robustness of these findings under statistical resampling. To this end, we conducted bootstrap analysis (N=1000 resamples) and computed 95% confidence intervals for both RMSE and accuracy.

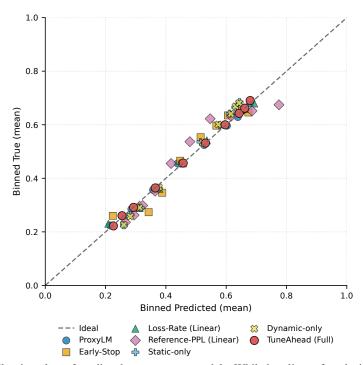


Figure A.1: Calibration plots of predicted scores across models. While baselines often deviate from the ideal diagonal, **TUNEAHEAD** (red) remains closely aligned with perfect calibration. This demonstrates that **TUNEAHEAD**'s predictions are both accurate in ranking and reliable in absolute probability estimation, complementing the aggregate results reported in Table 1.

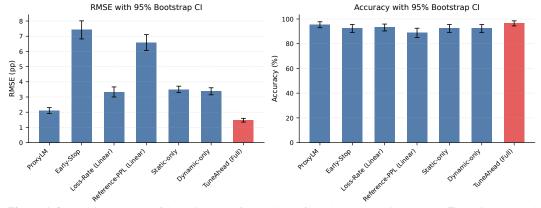


Figure A.2: 95% bootstrap confidence intervals for RMSE (left) and accuracy (right) across **TUNEAHEAD** and baseline predictors. Red bars highlight **TUNEAHEAD** (Full), while blue bars indicate baselines. The results confirm that **TUNEAHEAD** not only achieves the lowest RMSE but also the highest accuracy, with narrow confidence intervals that do not overlap with most baselines. This reinforces the statistical reliability of the improvements reported in Table 1.

As shown in Figure A.2, **TUNEAHEAD** achieves the most stable performance: its RMSE is significantly lower than all baselines, and its accuracy is consistently higher. The narrow error bars demonstrate that these results are not due to random variation, but reflect a statistically robust advantage of combining static and dynamic signals.

A.10 PER-HYPERPARAMETER BREAKDOWN

To verify that the improvements of **TUNEAHEAD** are not tied to a narrow choice of optimization settings, we report test-set performance grouped by key hyperparameters as shown in the tables below.

Table A.2: Per-learning-rate performance of TUNEAHEAD on the test set.

Learning rate	RMSE	$R^2 \uparrow$	$r \uparrow$	Accu	pp) ↑	
	γ				k=2	k=3
2×10^{-4}	1.36	0.99	1.00	51.68	85.23	96.64
3×10^{-4}	1.61	0.99	1.00	47.90	78.99	93.28

Table A.3: Per-batch-size performance of TUNEAHEAD on the test set.

Batch size	RMSE L	$R^2 \uparrow$	$r \uparrow$	Accu	pp) †	
	v	10	. 1	k=1	k=2	k=3
8	1.58	0.99	1.00	51.96	76.47	92.16
16	1.47	0.99	1.00	45.56	84.44	96.67
32	1.32	0.99	1.00	52.63	88.16	97.37

Table A.4: Two-way breakdown by (learning rate \times batch size).

Learning rate	Batch size	RMSE ↓	Batch size RMSE ↓		$r \uparrow$	Accu	racy (±k	pp) ↑
		14.152 γ	10		k=1	k=2	k=3	
2×10^{-4}	8	1.37	0.99	1.00	56.16	82.19	95.89	
2×10^{-4}	16	1.32	0.99	1.00	50.00	87.50	97.50	
2×10^{-4}	32	1.39	0.99	1.00	44.44	88.89	97.22	
3×10^{-4}	8	2.03	0.98	0.99	41.38	62.07	82.76	
3×10^{-4}	16	1.59	0.99	1.00	42.00	82.00	96.00	
3×10^{-4}	32	1.27	0.99	1.00	57.50	85.00	95.00	

Overall, the per-hyperparameter breakdown reinforces the robustness of TUNEAHEAD. First, across different learning rates (Table A.2), the predictor achieves consistently low RMSE (1.3–1.6pp) and nearly perfect correlation ($r\approx 1.0$). The slight drop in Acc@2pp for 3×10^{-4} (79% vs. 85% at 2×10^{-4}) suggests that steeper learning rates may induce more variance in the fine-tuning trajectories, but the deviations remain small relative to the overall gains reported in Table 1.

Second, when grouping by batch size (Table A.3), we find that larger batches generally yield slightly lower RMSE and higher calibration accuracy. For instance, batch size 32 achieves Acc@2pp of 88%, compared to 76% for batch size 8. This trend aligns with the intuition that smaller batches introduce more stochasticity, yet TUNEAHEAD is able to maintain strong predictive reliability across the range.

Finally, the two-way grid (Table A.4) confirms that performance remains strong across all (learning rate × batch size) combinations. Even in the most challenging regime (learning rate 3×10^{-4} , batch size 8), RMSE stays within 2pp and correlation remains high (r=0.99). Conversely, more stable settings such as (2×10^{-4} , 16 or 32) yield near-ideal calibration with Acc@3pp above 97%. These results demonstrate that **TuneAhead**'s improvements are not contingent on narrow hyperparameter choices, but generalize broadly across optimization regimes.

B META-FEATURE COMPENDIUM

This appendix provides a detailed compendium of all candidate static and dynamic meta-features engineered for the TUNEAHEAD framework. For each feature, we include its **definition**, **formula** (when applicable), **acquisition method**, **predictive hypothesis** (**signal**), and relevant literature references.

B.1 STATIC META-FEATURES

Static features are computed from the dataset D_i prior to training, sometimes using the frozen base model \mathcal{M}_0 .

B.1.1 GLOBAL STATISTICS

• Token Lengths (Mean & Std Dev)

$$\mu_{\mathrm{len}} = rac{1}{N} \sum_{j=1}^{N} |x_j|, \qquad \sigma_{\mathrm{len}} = \sqrt{rac{1}{N} \sum_{j=1}^{N} \left(|x_j| - \mu_{\mathrm{len}}
ight)^2}$$

Acquisition: Compute from input/output token counts. **Signal:** Longer average token lengths can increase model perplexity and memory overhead, while shorter inputs may not provide sufficient context. High variance in sequence lengths often signals heterogeneity in data sources or task types, leading to unstable optimization (Vettoruzzo et al., 2023; Moghe et al., 2024).

• Input-Output Length Ratio

$$r = \frac{1}{N} \sum_{j=1}^{N} \frac{|y_j|}{|x_j|}$$

Acquisition: Average of output-to-input length ratios. **Signal:** This ratio correlates with task paradigms: low ratios characterize compressive tasks such as summarization, while high ratios reflect elaborative tasks like question answering or code generation. Such structural differences shape model adaptation speed and generalization (Lin, 2004; Narayan et al., 2018).

• Special Character & Code Ratio

$$r_{\rm sc} = rac{\# ext{ special or code tokens in } D_i}{\# ext{ total tokens in } D_i}$$

Acquisition: Count of special/code tokens vs. total. **Signal:** Elevated ratios typically indicate domain-specific corpora (e.g., programming or formula-heavy data) or unclean web text. Such distributions require tailored tokenization or model adaptation, as shown in large-scale text-to-text transfer studies (Allamanis et al., 2018; Raffel et al., 2020).

• Approximate Duplicates Ratio

$$r_{\text{dup}} = \frac{|\{(s_i, s_j) : \text{Sim}(E(s_i), E(s_j)) > \tau\}|}{N^2}$$

Acquisition: Identify near-duplicates via embedding similarity threshold τ . **Signal:** High duplication reduces effective dataset diversity, amplifies memorization, and weakens generalization. Deduplication in LLM pretraining has been shown to improve downstream performance and reduce overfitting (Carlini et al., 2023; Lee et al., 2022).

• Embedding Outlier Ratio

$$r_{\mathrm{outlier}} = \frac{|\{E(s): \|E(s) - \mu\| > 3\sigma\}|}{N}$$

Acquisition: Detect large deviations in embedding space. **Signal:** Outlier samples often correspond to mislabeled, noisy, or domain-shifted data. Their presence destabilizes optimization and can severely degrade model robustness. Removing outliers is a key step in modern dataset curation pipelines (Hendrycks et al., 2019; Northcutt et al., 2022; Dodge et al., 2021).

• Dataset Size (Num Items)

$$|D_i| = N$$

Acquisition: Dataset example count. **Signal:** Larger datasets typically improve model performance, but the gains follow a power-law with diminishing returns. Scaling law analyses show that optimal performance requires balancing dataset size with model capacity and compute budget(Kaplan et al., 2020; Hoffmann et al., 2022).

B.1.2 LEXICAL DIVERSITY FEATURES

• Type-Token Ratio (TTR)

$$TTR(D_i) = \frac{|Vocab(D_i)|}{|Tokens(D_i)|}$$

Acquisition: Compute vocabulary diversity in D_i . **Signal:** A low TTR indicates repetitive content and limited lexical variety, which can impair a model's ability to generalize to unseen expressions. High TTR reflects lexical richness but may also introduce noise or rare tokens. TTR is a long-established measure of lexical richness and a standard meta-feature in dataset characterization (Rivolli et al., 2019).

N-Gram Repetition Rate

$$r_n = \frac{1}{N} \sum_{j=1}^{N} \frac{\text{\# repeated } n\text{-grams in } y_j}{|y_j|}$$

Acquisition: Fraction of repeated n-grams in outputs. **Signal:** High repetition rates often signal low-quality or synthetic outputs, reducing effective informational content and promoting degenerate training behavior. Low repetition may indicate dispersed data but could also reduce coherence. Repetition metrics are widely used in text degeneration and quality control studies(Rivolli et al., 2019; Holtzman et al., 2020).

• Instruction Complexity

$$C_{\text{inst}} = \frac{1}{N} \sum_{j=1}^{N} \text{depth}(\text{ParseTree}(x_j))$$

Acquisition: Average parse-tree depth of instructions. **Signal:** Shallow trees suggest trivial instructions that under-challenge the model, while overly deep trees reflect high syntactic complexity that may hinder comprehension or stable learning. Balanced complexity encourages both learnability and generalization (Yatskar, 2019)

B.1.3 Information-theoretic properties

• Reference Perplexity (PPL)

$$PPL(D_i, \mathcal{M}_0) = \exp\left(-\frac{1}{\sum_{s \in D_i} |s|} \sum_{s \in D_i} \log p_{\mathcal{M}_0}(s)\right)$$

Acquisition: Measure using frozen base model. **Signal:** High perplexity indicates distributional mismatch between dataset and pretraining corpus, leading to slower convergence and increased adaptation cost. Low perplexity suggests greater alignment with prior knowledge. PPL remains a widely accepted proxy for domain mismatch and learning difficulty(Wu et al., 2017; Jozefowicz et al., 2016).

• Input-Output Semantic Similarity (IO Similarity)

$$Sim(x, y) = \frac{E(x) \cdot E(y)}{\|E(x)\| \|E(y)\|}$$

Acquisition: Cosine similarity of embeddings. **Signal:** Low similarity may reflect irrelevant or hallucinated outputs, while very high similarity often indicates trivial paraphrasing lacking informativeness. Moderate levels of similarity are most effective for meaningful adaptation. Embedding-based similarity has been widely studied in representation learning (Clark et al., 2020; Cer et al., 2018).

Output Semantic Diversity

$$Div(D_i) = \frac{2}{N(N-1)} \sum_{j < k} (1 - Sim(y_j, y_k))$$

Acquisition: Average pairwise output dissimilarity. **Signal:** Low diversity signals redundancy and narrow coverage, while excessive diversity may indicate incoherence or noisy task signals. Balanced semantic diversity provides both robustness and coverage, promoting better generalization (Ziegler et al., 2020; Li et al., 2016).

• LM-Data Vocabulary Alignment (KL Divergence)

$$\mathrm{KL}(P\|Q) = \sum_{w \in V} P(w) \log \frac{P(w)}{Q(w)}$$

Acquisition: Compare dataset and reference corpus word frequencies. **Signal:** Large divergence highlights domain shift, suggesting that the dataset vocabulary departs from the pretraining distribution. This increases the adaptation burden and may reduce efficiency. KL divergence is a standard measure in domain adaptation and data selection (Aharoni & Goldberg, 2020; Axelrod et al., 2011).

• Answer Groundedness

$$g(y,x) = \frac{|n\text{-grams}(y) \cap n\text{-grams}(x)|}{|n\text{-grams}(y)|}$$

Acquisition: Ratio of overlapping n-grams between output and input. **Signal:** Low groundedness suggests hallucination or irrelevant generation, while overly high groundedness may reduce informativeness by copying excessively. Moderate grounding balances fidelity with informativeness, ensuring both reliability and novelty (Ji et al., 2023; Zhao et al., 2020).

B.2 Dynamic Probe Meta-Features

Dynamic features are extracted during a standardized 100-step probe run.

B.2.1 LOSS-BASED INDICATORS

• Initial Loss

$$L_0 = L(\theta_0; D_i)$$

Acquisition: Probe loss at the first optimization step. **Signal:** The initial loss measures how well the pretrained model aligns with the dataset before adaptation. High values suggest a significant domain gap, requiring more updates to adapt, while low values indicate better alignment and easier fine-tuning. It is widely used as a proxy for domain difficulty in transfer learning (Arpit et al., 2017; Hestness et al., 2017).

Loss Decay Rate

$$\alpha = \text{slope}(\text{LinReg}(\{(t, L_t)\}_{t=1}^T))$$

Acquisition: Slope of regression fit on probe loss curve. **Signal:** A steep negative slope indicates strong learnability and rapid adaptation, whereas flat or unstable curves suggest noisy or hard-to-learn data. Early loss decay is highly predictive of final model performance (Wu et al., 2017; Hestness et al., 2017; Loog & Viering, 2022).

• Loss Curve Stability

$$\sigma_L^2 = \frac{1}{T} \sum_{t=1}^{T} \left(L_t - \overline{L} \right)^2$$

Acquisition: Variance of loss during probe. **Signal:** Low variance implies stable optimization and smoother convergence, while high fluctuations often reflect noisy data, poor learning rates, or unstable alignment between model and task. Stable trajectories are associated with better generalization(Li et al., 2024; Wu et al., 2017).

B.2.2 Gradient-based indicators

• Gradient Norm (Mean & Variance)

$$\mu_g = \frac{1}{T} \sum_{t=1}^{T} \|g_t\|_2, \quad \sigma_g^2 = \frac{1}{T} \sum_{t=1}^{T} (\|g_t\|_2 - \mu_g)^2$$

Acquisition: Gradient norms across probe steps. **Signal:** Gradient norms reflect the strength of learning signals. Very large norms can cause instability or gradient explosion, while very small norms may stall learning or trap the model in poor local minima. Both mean and variance provide insight into learning dynamics (Pascanu et al., 2013; Sutskever et al., 2013; Killamsetty et al., 2021).

Gradient Consistency

$$c_t = \frac{g_t \cdot g_{t+1}}{\|g_t\| \|g_{t+1}\|}$$

Acquisition: Cosine similarity between sequential gradients. **Signal:** High consistency indicates coherent optimization paths, suggesting the model is learning a stable objective. Low or negative alignment suggests noisy or conflicting signals, slowing convergence. Gradient alignment is also linked to meta-learning generalization (Finn et al., 2017; Killamsetty et al., 2021; Guiroy et al., 2019)

Gradient Sparsity

$$s = \frac{|\{k : |g_k| < \epsilon\}|}{|q|}$$

Acquisition: Proportion of near-zero gradient components. **Signal:** High sparsity indicates that only a small subset of parameters is being updated, potentially limiting adaptation. Moderate sparsity may improve efficiency and generalization, but excessive sparsity may signal model–data mismatch (Frankle & Carbin, 2019; Evci et al., 2021; Killamsetty et al., 2021).

B.2.3 MODEL-BASED INDICATORS

Parameter Change Norm

$$\Delta \theta = \|\theta_T - \theta_0\|_2$$

Acquisition: L2 norm of parameter changes during probe. **Signal:** Parameter change magnitude reflects adaptation strength. Moderate changes indicate healthy learning, while excessive shifts may reflect instability or overfitting. This feature has been used to analyze learning dynamics and compression strategies (Li et al., 2020; Raghu et al., 2017).

Loss Landscape Flatness Proxy

$$\Delta L = L(\theta^* + \delta) - L(\theta^*)$$

Acquisition: Apply perturbation δ and measure loss change. **Signal:** Flat minima (small ΔL) correspond to more robust solutions with better generalization under distribution shifts, while sharp minima indicate overfitting and fragility. Flatness has been consistently linked to generalization performance(Wu et al., 2017; Li et al., 2024).

Catastrophic Forgetting Proxy

$$\Delta P = P_{\text{baseline}} - P_{\text{probe}}$$

Acquisition: Performance drop on out-of-domain task during probe. **Signal:** Large drops indicate interference between new and old tasks, a hallmark of catastrophic forgetting. This proxy highlights whether fine-tuning data compromises existing knowledge(Wen & Itti, 2018).

Activation Sparsity

$$s_a = \frac{|\{h : |h| < \epsilon\}|}{|h|}$$

Acquisition: Fraction of near-zero activations in hidden layers. **Signal:** Sparse activations suggest selective use of model capacity. Moderate sparsity improves interpretability and efficiency, while excessive sparsity indicates underutilized capacity or inefficient learning. It is widely used as a proxy for model resource utilization (Glorot et al., 2011; Maass, 1997).

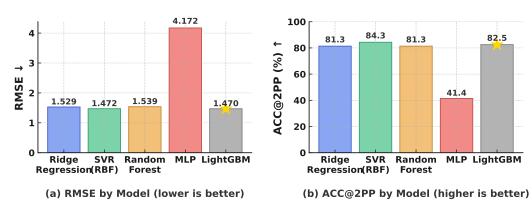


Figure C.1: Performance comparison of predictor models. We evaluate five models on the same meta-feature set. The plots show (a) RMSE (lower is better) and (b) Acc@2pp (higher is better). The results validate our choice of LightGBM (marked by a star), which also achieves strong performance.

C DETAILED ABLATION STUDY ON PREDICTOR CHOICE

Motivation and Setup. A critical design choice in the **TUNEAHEAD** framework is the selection of the prediction model F. The ideal predictor must not only achieve state-of-the-art predictive accuracy (**G2**) but also align with our core principles of providing interpretable diagnostics (**G3**) and ensuring computational efficiency (**G1**). To empirically validate our choice of LightGBM, we conducted a comprehensive comparison against a diverse suite of strong and representative regression models:

- Ridge Regression (Hoerl & Kennard, 1970): a powerful linear model to test the extent of non-linear relationships in the data.
- Support Vector Regressor (SVR, RBF) (Drucker et al., 1997; Smola & Schölkopf, 2004):
 a classic, high-performance kernel-based method adept at capturing complex non-linearities.
- Random Forest (Breiman, 2001): a state-of-the-art ensemble model based on bagging, serving as a direct comparison to LightGBM's boosting approach.
- Multi-Layer Perceptron (MLP) (Rumelhart et al., 1986): a simple but representative neural network baseline for tabular data.

To ensure a fair and rigorous comparison, all models were trained on the identical meta-feature set, and each underwent a systematic hyperparameter search using 5-fold cross-validation on our training set.

Performance Analysis. The results of this comparison are presented in Figure C.1. Our analysis yields two key findings. First, a top tier of models clearly emerges, with LightGBM (RMSE=1.470) and SVR (RMSE=1.472) delivering nearly identical, state-of-the-art predictive accuracy. Random Forest (RMSE=1.589) follows as another strong competitor. This confirms that a GBDT-based approach achieves performance that is on par with the best alternative methods for this task. Second, the simple MLP struggles to generalize effectively (RMSE=4.172), a common outcome on heterogeneous, tabular meta-datasets where GBDT models often excel without extensive architectural tuning.

Justification for Selecting LightGBM. Given the statistically comparable accuracy of the top-performing models (LightGBM and SVR), our final selection was determined by the other two crucial design goals: interpretability and scalability.

Interpretability. LightGBM holds a decisive advantage. Its tree-based architecture integrates seamlessly with SHAP, enabling the precise, feature-level diagnostics that are central to TUNEAHEAD's mission. In contrast, while SVR is a powerful predictor, deriving similarly intuitive, local feature-level attributions from a kernel-based model is significantly more complex and less direct.

Scalability. Furthermore, LightGBM is substantially more scalable. Its training time complexity is more favorable than SVR's, particularly as the number of experiments (samples) in the metadataset grows. This computational efficiency is critical for the future development and application

of TUNEAHEAD to even larger and more diverse problem spaces, as discussed in our Future Work (Sec. 6). Conclusion. While SVR demonstrates highly competitive accuracy on our current dataset, Light-GBM's unique combination of top-tier accuracy, superior interpretability, and better scalability makes it the most principled and strategic choice for the TUNEAHEAD framework.