# LANDMARK-GUIDED POLICY OPTIMIZATION FOR MULTI-OBJECTIVE LANGUAGE MODEL SELECTION

# Anonymous authors

Paper under double-blind review

# **ABSTRACT**

Selecting a pretrained large language model (LLM) to fine-tune for a task-specific dataset can be time-consuming and costly. With several candidate models available to choose from, varying in size, architecture, and pretraining data, finding the best often involves extensive trial and error. In addition, the "best" model may not necessarily be the one with the lowest test loss, as practical considerations such as deployment costs, inference throughput, and limited search budgets might also play crucial roles. To address this, we introduce LAMPS (LAnguage Model Pareto Selection), a novel and open-source multi-objective AutoML framework that quickly identifies near-Pareto-optimal pretrained LLMs for a task-specific dataset. It is based on two key ideas: (1) landmark fine-tuning, which generates early performance indicators of the candidate models, and (2) meta-learning via reinforcement learning, which learns an effective selection policy from historical performance data (a meta-dataset). Our results show that, on held-out datasets, LAMPS reduces search time by an average of 71% compared to exhaustive search, while still covering more than 98% of the optimal target space hypervolume.

#### 1 Introduction

Supervised fine-tuning of a pretrained large language model (LLM) on task-specific datasets is currently the dominant paradigm for achieving state-of-the-art performance in several natural language processing (NLP) tasks (Radford et al., 2019), including question answering (Chowdhery et al., 2023), machine translation (Raffel et al., 2020), summarization (Aghajanyan et al., 2020), and classification (Yang, 2019). However, different pretrained models yield varying downstream performance due to differences in size, architecture, pretraining data, and other intrinsic factors. Therefore, as the set of available pretrained models is already extensive, the important question arises: How can we efficiently find the best-performing model for a task-specific dataset?

A common practice in NLP is to select the largest available model, driven by the belief that larger models invariably provide better performance (e.g., accuracy, F1, perplexity, cross-entropy, depending on the downstream task). Although this is generally true, several studies have shown that smaller models are comparable to or even outperform larger ones (Ouyang et al., 2022; Sanh et al., 2020; Hoffmann et al., 2022; Wahba et al., 2023; DeepSeek-AI et al., 2025; Wang et al., 2025). Moreover, in real-world scenarios, always choosing larger models inevitably leads to higher operational costs and greater environmental impact. This underscores the need to incorporate additional factors into the model selection process beyond a single task-specific performance metric.

A multi-objective perspective is, then, essential to capture the broader spectrum of trade-offs that practitioners face when selecting pretrained LLMs for fine-tuning. In the absence of better alternatives, practitioners may turn to exhaustive search. Although theoretically sound, this method quickly becomes prohibitively expensive for a large number of candidate models, especially for target datasets with several million examples. As language models continue to expand in scale and diversity, there is an increasing need for a principled, holistic, and efficient selection strategy, especially with the growing interest in specialized LLM-based AI agents (Gutowska, 2024; Ma et al., 2024).

In this paper, we introduce LAMPS (LAnguage Model Pareto Selection), a novel and open-source multi-objective AutoML framework for selecting LLMs to fine-tune on task-specific datasets. It

integrates two complementary strategies: (1) *landmark fine-tuning*, which generates early performance indicators for candidate models by evaluating them on incrementally larger subsets of the training data; and (2) *meta-learning via reinforcement learning*, which leverages historical model performance data on multiple datasets to learn how to efficiently allocate training resources for new datasets. In other words, this process generates a policy that manages the selection and early stopping of candidate models, adjusting its strategy based on both observed and historical performance to efficiently discard low-potential models and prioritize promising ones.

Our main contributions are as follows: (i) Formulating the language model selection for fine-tuning explicitly as a multi-objective optimization problem; (ii) Introducing LAMPS, a novel and open-source AutoML framework combining landmark fine-tuning, meta-learning, and reinforcement learning to rapidly identify near-Pareto-optimal language models for a new task-specific dataset.

The remainder of the paper is organized as follows. Section 2 reviews relevant related work. Section 3 states the multi-objective optimization problem. The method is proposed in Section 5 and Section 6 presents the experimental setup and main findings. We conclude in Section 7.

#### 2 RELATED WORK

Selecting an appropriate base learner (model, algorithm, pipeline, etc.) for a given task has been a long-standing research topic and is usually called *model selection* (Bozdogan, 1987; Maron & Moore, 1993; McQuarrie & Tsai, 1998; Chapelle et al., 2002; Biem, 2003; Brazdil et al., 2003; Zhao & Yu, 2006; Adankon & Cheriet, 2009). Among the different approaches available, metalearning has been a popular choice (Kalousis & Hilario, 2000; Fürnkranz et al., 2002; Brazdil & Giraud-Carrier, 2018; Jain et al., 2024; de Amorim et al., 2025; Farhadi et al., 2025), mainly due to its ability to transfer knowledge from prior learning experiences, reducing the cost of exploration and improving sample efficiency.

In this section, we provide a brief overview of the related areas that form the foundation of our LAMPS framework.

**Pretrained Model Selection in Deep Learning** Fine-tuning pretrained deep learning models for specific downstream tasks has become the standard approach in both computer vision and natural language processing. Compared to training from scratch, fine-tuning is far more efficient and requires much less data than pretraining (Hepburn, 2018). For this reason, being able to select the right pretrained model efficiently is becoming increasingly relevant due to the considerable computational costs and the rapid introduction of new models with varying sizes, architectures, training data, and capabilities. To the best of our knowledge, the only work that explicitly addresses the selection of LLMs for fine-tuning is by Monteiro et al. (2024), but it neither considers the multi-objective aspects of the model selection nor adjusts its recommendations based on actual fine-tuning learning curves.

**Subsampling Landmarks** A sampling landmark is a performance-based meta-feature, representing the performance of a particular model on samples of available data, providing a quick estimate of its performance (Brazdil et al., 2022; Pfahringer et al., 2000) and, consequently, allowing indirect characterization of the target dataset. One variant is called <u>subsampling landmarks</u>, which considers a sequence of sample sizes in increasing order, effectively representing the early stages of the learning curve (Soares et al., 2001; Fürnkranz & Petrak, 2001). This is conceptually related to the scaling laws observed in deep neural networks (Kaplan et al., 2020) and large language models (Zhang et al., 2024), which describe the predictable relationship between model performance and, among other factors, dataset size. Subsampling landmarks can thus be viewed as a localized and practical proxy for these scaling behaviors, enabling performance forecasting without requiring full-scale training. Similar ideas have been applied for hyperparameter optimization (Domhan et al., 2015; Jamieson & Talwalkar, 2016; Klein et al., 2017; Li et al., 2018), which use partial learning curves to stop training poor configurations early. Such methods, however, remain inherently single-objective and cannot directly address the multi-objective settings considered in this work.

Multi-Task and Meta-Reinforcement Learning Reinforcement learning is a powerful tool for sequential decision-making problems, but it often struggles with generalization to new (unknown)

tasks, requiring large amounts of data to readapt effectively. Two areas address these limitations: multi-task reinforcement learning (MTRL) (Teh et al., 2017; Sodhani et al., 2021) and metareinforcement learning (Meta-RL) (Finn et al., 2017; Nichol et al., 2018; Wang et al., 2024). MTRL trains a single policy across a distribution of tasks, leveraging shared structure to improve generalization and learning efficiency. In contrast, Meta-RL focuses on learning a policy that can rapidly adapt to new tasks using limited data, typically by encoding task-specific information into its internal state or parameters. In this work, we focus on MTRL, as our goal is to evaluate policies on previously unseen datasets without further adaptation at test time.

Multi-Objective Reinforcement Learning Multi-objective reinforcement learning (MORL) extends standard RL by optimizing policies with respect to multiple, often conflicting objectives rather than a single reward. Prior research on MORL, often combined with meta-learning, has largely relied on scalarization or objective preferences, requiring weight sweeps across many preferences to approximate the Pareto front (Lu et al., 2024; Wang et al., 2024; Liu & Qian, 2021; Chen et al., 2019). Because each weight vector defines a different scalar objective, changing preferences generally requires another sweep (i.e., additional fine-tuning runs), so computation grows with each revision. By contrast, we target Pareto coverage in a single, efficient run.

# 3 PROBLEM STATEMENT

Consider a target dataset  $\mathcal{D}$  and a set  $\mathcal{X}$  of candidate pretrained language models to be fine-tuned. Then, given n metrics of interest (objectives), the problem can be formulated as the following multi-objective optimization problem:

$$\min_{x \in \mathcal{X}} \quad (f_1(x, \mathcal{D}), \dots, f_n(x, \mathcal{D}))$$
s.t.  $f_i(x, \mathcal{D}) \le f_i^{\max} \text{ for all } i = 1, \dots, n,$  (1)

where  $f_i(x, \mathcal{D})$  represents the value of the *i*-th objective function after fine-tuning the pretrained model  $x \in \mathcal{X}$  on the task-specific dataset  $\mathcal{D}$ , and  $f_i^{\max}$  denotes an arbitrary upper bound for that objective.

Common objectives may include final test loss, training time (cost), inference throughput, number of model parameters, and resource usage (i.e., number of GPUs). It is very common that some objectives conflict with each other. For example, achieving a lower test loss may require longer training time or more GPUs. For this reason, there is typically no single solution that is optimal across all objectives. Hence, the notion of optimality is based on Pareto-dominance, or simply dominance, as defined below.

**Definition 1** (Weak dominance). A solution  $x_1 \in \mathcal{X}$  weakly dominates another solution  $x_2 \in \mathcal{X}$ , denoted  $x_1 \succeq x_2$ , if  $f_i(x_1, \mathcal{D}) \leq f_i(x_2, \mathcal{D})$  for all  $i \in \{1, \dots, n\}$ . That is,  $x_1$  is not worse than  $x_2$  in all objectives.

**Definition 2** (Pareto-dominance). A solution  $x_1 \in \mathcal{X}$  dominates another solution  $x_2 \in \mathcal{X}$ , denoted  $x_1 \succ x_2$ , if  $f_i(x_1, \mathcal{D}) \leq f_i(x_2, \mathcal{D})$  for all  $i \in \{1, \dots, n\}$ , with at least one of these inequalities holding strictly. That is, there is  $j \in \{1, \dots, n\}$  such that  $f_j(x_1, \mathcal{D}) < f_j(x_2, \mathcal{D})$ . In other words,  $x_1$  dominates  $x_2$  if  $x_1$  is not worse than  $x_2$  in all objectives, but it is better in at least one of them.

**Definition 3** (Pareto-optimal). A model  $x^* \in \mathcal{X}$  is Pareto-optimal if there is no other  $x \in \mathcal{X}$  that dominates  $x^*$ .

One way to evaluate and compare sets of candidate solutions is to use the *hypervolume indicator* (Guerreiro et al., 2021; Emmerich et al., 2005), which quantifies the volume of the objective space weakly dominated by a set of solutions and bounded above by a given reference point  $r = [f_1^{\max}, \dots, f_n^{\max}]^{\top}$ . For any subset  $X \subset \mathcal{X}$ , the hypervolume indicator is denoted as  $H_{\mathcal{D}}(X, r)$ . Intuitively, each solution in X defines a box in the objective space, with one corner at the objective values of the solution and the opposite corner at the reference point r. It is defined formally as follows:

**Definition 4** (Hypervolume indicator). Given a set of points  $S \subset \mathbb{R}^n$  and a reference point  $r \in \mathbb{R}^n$ , the hypervolume indicator of S is the measure of the region weakly dominated by S and bounded

above by r, i.e.,

162

163

164

166 167

168

169 170

171

172

173 174

175 176

177

178

179

181

182

183

185

186

187

188

189

190

191

192

193

194

195

196 197

199

200

201 202

203 204

205

206

207

208

209

210

211

212

213

214

215

$$H(S,r) = \Lambda \left( \bigcup_{\substack{p \in S \\ p \le r}} [p,r] \right),$$

where  $\Lambda(\cdot)$  denotes the Lebesgue measure and  $[p,r] = \{q \in \mathbb{R}^n \mid \forall i=1,\ldots,n: p_i \leq q_i \leq r_i\}$ denotes the box delimited below by  $p \in S$  and above by r.

It has been shown that maximizing the hypervolume indicator is equivalent to finding the Pareto optimal set (Guerreiro et al., 2021; Liu et al., 2019). Figure 1 illustrates this with a practical comparison, showing that the Pareto-optimal set has the highest hypervolume. Thus, the problem in (1) can be reformulated as a single-objective problem as follows:

8.0

0.7

0.6

0.5

0.4

Test Loss

$$\max_{X \subset \mathcal{X}} \quad H_{\mathcal{D}}(X, r) \tag{2}$$

A trivial solution would involve fine-tuning all models on the target dataset (i.e.,  $X = \mathcal{X}$ ), but this is computationally intractable. To encourage computational efficiency, we introduce a regularization term penalizing the number of selected pretrained models:

$$\max_{X \subset \mathcal{X}} \quad H_{\mathcal{D}}(X, r) - \lambda |X| \tag{3}$$

where  $\lambda > 0$  is a user-defined penalty factor. To ensure that the optimal solution for the problem in Equation 3 contains exactly all Paretooptimal solutions,  $\lambda$  must satisfy the following theorem, proved in Appendix G:

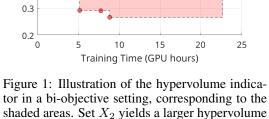
**Theorem 1** (Condition on  $\lambda$ ). *The optimal so*lution  $X^* \subset \mathcal{X}$  of problem 3 contains only and exactly all Pareto-optimal solutions if and only if:

$$0 < \lambda \le \min_{x \in X^*, X \subset X^*} \Delta H_{\mathcal{D}}(x \mid X), \quad (4)$$

where  $\Delta H_{\mathcal{D}}(x \mid X)$  denotes the incremental

In other words, the penalty  $\lambda$  must be smaller than or equal to the smallest incremental hypervolume gained by including a new Pareto-optimal solution into the subset of selected solution candidates. If this condition holds, the optimal solution set will include only all Pareto-optimal solutions. The

next sections present empirical strategies for quickly providing near-Pareto optimal solutions.



**Objective Space** 

Pareto-optimal

Solution set X,

than  $X_1$ , which is closer to the true Pareto front. hypervolume obtained by adding the Pareto-optimal solution x to the subset  $X \subseteq X^*$ .

# LANDMARK FINE-TUNING

Fine-tuning a pretrained model on a task-specific dataset is inevitable if one desires to evaluate its true performance and determine its suitability for a given application. However, as discussed earlier, evaluating every candidate model is computationally expensive. Prior work on hyperparameter optimization suggests that evaluating models for only a single epoch can already be a good proxy for its final performance (Egele et al., 2023). However, training for just one epoch may still consume significant resources, particularly for large models and datasets.

To mitigate this inefficiency and allow for even earlier identification of unpromising candidates, we propose landmark fine-tuning, a lightweight fine-tuning strategy based on subsampling landmarks to obtain early estimates of objective values  $f_i(x, \mathcal{D})$  for  $i \in \{1 \dots n\}$ .

Given that the target dataset  $\mathcal{D}$  has a training and a test split, namely  $\mathcal{D}^{train}$  and  $\mathcal{D}^{test}$ , the core idea is to split  $\mathcal{D}^{\text{train}}$  into K exponentially larger subsets  $\mathcal{D}_1 \dots \mathcal{D}_K$ . Each subset  $\mathcal{D}_k$  contains  $\left|\frac{1}{2^{(K-k)}}|\mathcal{D}^{\text{train}}|\right|$  samples, where  $\mathcal{D}_k \subset \mathcal{D}_{k+1}$  for  $k=1\ldots K-1$ .

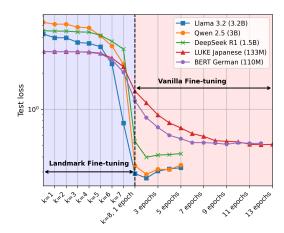


Figure 2: Landmark fine-tuning on the 20 Newsgroups dataset using K=8. Larger models start off worse but eventually outperform smaller ones. Notably, models that improve quickly early on tend to achieve lower final loss, suggesting that the initial segments of the learning curves can help predict the overall performance.

The process starts by fine-tuning a pretrained model on  $\mathcal{D}_1$  for a *single epoch* and evaluating it on the entire  $\mathcal{D}^{\text{test}}$ . Next, it continues the fine-tuning process on the subsequent (larger) subset  $\mathcal{D}_2$ , repeating this process up to  $\mathcal{D}_K$  (100% of the training dataset). After that, we continue fine-tuning the model for more epochs until convergence or other stop criterion.

Figure 2 shows a practical example of landmark fine-tuning with K=8, depicting the learning curves (test cross-entropy loss) of five different pretrained models fine-tuned on the 20 Newsgroups dataset. Two non-English LLMs are included to illustrate the performance of less suitable models on an English dataset. Notice that larger models start with higher losses than smaller ones, but eventually overtake them, achieving lower final losses. In addition, among the larger models, those that improve more quickly in the initial steps tend to achieve better final test loss. These observations support the idea that early segments of the training curve can indeed be predictive of final loss, with predictions becoming more accurate as additional curve segments are provided.

# 5 META-LEARNED RESOURCE ALLOCATION VIA REINFORCEMENT LEARNING

Although landmark fine-tuning provides early performance estimates, it is still necessary to determine when to continue training a candidate model or not, based on partial information collected so far. To address this, we train a reinforcement learning agent on a meta-dataset of historical fine-tuning trajectories, covering a diverse set of pretrained LLMs and downstream tasks. The agent learns to allocate training resources by tracking how performance evolves across landmark steps, enabling fast and generalizable identification of near-Pareto-optimal models.

**Observation space** The observation space defines the information available to the RL agent at each decision step. At each time step t, the RL agent observes, for every candidate model, the objectives of interest (e.g., the elapsed training time and test loss), together with the number of fine-tuning steps that each candidate has completed.

Action space The action space specifies the set of decisions available to the RL agent at each step. For each time step t, the agent selects an action  $a_t \in \{1, \ldots, m\}$ , representing the index of a candidate pretrained model, where m is the total number of candidates. Each action corresponds to allocating one additional fine-tuning step to the selected model. To improve exploration efficiency, we apply invalid action masking for terminated models (Huang & Ontañón, 2022). A binary mask specifies which models remain available for selection. The policy then samples only from this valid subset by setting the probability of invalid actions to zero. This prevents wasted trials on completed

models and makes the exploration phase more efficient, as the agent can focus its decisions on candidates that may still yield improvements.

**Termination condition** An episode corresponds to the full search process and terminates when all Pareto-optimal models have been fully fine-tuned, thereby achieving the maximum hypervolume. Thanks to invalid action masking, the episode is guaranteed to terminate within a finite number of steps, preventing the agent from getting stuck in infinite allocations to unproductive models.

**Training algorithm** For training the policy, we adopted Distral (distill and transfer learning), a framework for multi-task RL where the knowledge gained in one task is distilled into a shared policy, then transferred to other tasks via regularization using a Kullback-Leibler (KL) divergence (Teh et al., 2017). As the underlying optimizer, we adopted Proximal Policy Optimization (PPO) (Schulman et al., 2017), combining its stability with cross-task transfer from Distral.

**Rewards** The reward function links our multi-objective search problem to the policy's learning process. Let  $X_t \subseteq \mathcal{X}$  be the set of fully fine-tuned models by time step t, and let T be the length of the episode. Inspired by equation 3, we could initially define a sparse reward function

$$r_t = \begin{cases} H_{\mathcal{D}}(X_t) - \lambda |X_t| & \text{if } t = T\\ 0 & \text{otherwise} \end{cases}$$
 (5)

so that PPO would maximize

$$\max_{\theta} \quad \mathbb{E}_{\rho \sim \pi_{\theta}} \left[ \sum_{t=0}^{T} \gamma^{t} r_{t} \right] = \mathbb{E}_{\rho \sim \pi_{\theta}} \left[ H_{\mathcal{D}}(X_{T}) - \lambda |X_{T}| \right], \tag{6}$$

where  $\rho$  is a trajectory sampled using policy  $\pi_{\theta}$ , and  $\gamma$  is the discount factor.

Because an episode terminates only after all Pareto-optimal models have been fully fine-tuned,  $H_{\mathcal{D}}(X_T)$  is identical for every trajectory and, therefore, constant. The objective thus collapses to minimizing the expected number of models evaluated, i.e.,  $\mathbb{E}[-|X_T|]$ . Notice that  $\lambda$  also vanishes in this sparse reward setting, so we do not need to estimate it. Finally, to make the reward positive and incentivize faster convergence to the optimal, we adopted the following sparse reward function:

$$r_t = \begin{cases} \frac{|\mathcal{X} \setminus X_t|}{\Delta_t} & \text{if } t = T\\ 0 & \text{otherwise} \end{cases}$$
(7)

where  $\Delta_t$  is the cumulative wall-clock time spent up to time step t. In other words, we seek to maximize the number of pretrained models not fully fine-tuned, divided by the time spent to find all Pareto-optimal models. This produces a positive and well-scaled learning signal and preserves the optimal solution of equation 3, as the highest reward is obtained when  $X_T = X^*$ . Appendix C presents additional evidence showing that the proposed reward signal in equation 7 leads PPO to converge to the optimal solution during the training phase.

**Meta-dataset** To meta-train a policy capable of efficiently identifying (or approximating) the Pareto-optimal set for new task-specific datasets, we conducted a fine-tuning campaign and constructed a meta-dataset containing fully recorded learning curves of 70 pretrained LLMs, each landmark fine-tuned on multiple datasets (see Appendix D). This setup enables the agent to query arbitrary trajectories during its training, allowing the use of on-policy algorithms such as PPO.

**Deployment (search procedure)** Given a trained policy  $\pi_{\theta}$  and a target dataset  $\mathcal{D}$ , the search procedure of LAMPS is outlined in Algorithm 1. The process begins by constructing the initial state  $s_0$  through zero-shot evaluation of all candidate models on the test split. It also serves as a sanity check to ensure that each model is available, downloaded properly, and compatible with the available hardware (and drivers) where the search will be performed. The policy then proceeds by selecting and executing new actions until the search budget is exhausted. In the end, dominated solutions are filtered out, so that only the best trade-offs are presented to the user.

<sup>&</sup>lt;sup>1</sup>During policy training, the agent has access to privileged information to determine when the episode is finished.

```
324
           Algorithm 1: LAMPS Search Procedure
325
           Input: Training dataset \mathcal{D}^{train}, test dataset \mathcal{D}^{test}
326
           Input: Policy \pi_{\theta}
327
           Output: Set of non-dominated fine-tuned models \hat{X}^{\star}
328
        1 Initialize time step t \leftarrow 0:
330
        2 Initialize the set of selected models X \leftarrow \emptyset;
331
        3 Evaluate candidate models on \mathcal{D}^{\text{test}} and construct the initial state s_0;
332
        4 while search budget not exhausted do
333
                Select action a_t \leftarrow \arg \max_a \ \pi_{\theta}(a \mid s_t);
        5
334
                Fine-tune model x_{a_t} for one additional fine-tuning step on \mathcal{D}_{k+1}^{\text{train}};
335
                Evaluate updated performance on \mathcal{D}^{\text{test}};
336
                if stopping criterion met for model x_{a_{+}} then
337
                 X \leftarrow X \cup \{x_{a_t}\}
338
                Update environment state s_{t+1};
        10
339
                t \leftarrow t + 1;
340
       12 \hat{X}^* = \{ x \in X \mid \nexists y \in X : y \succ x \};
341
342
       13 return X^*;
343
```

# 6 EXPERIMENTS AND RESULTS

This section presents our experimental setup and main findings, demonstrating how well the trained policy generalizes to held-out datasets. All experiments in this paper were conducted on eight NVIDIA A100 (40 GB) GPUs.

#### 6.1 EXPERIMENTAL SETUP

**Pretrained LLMs** In our experiments, we tested 70 different pretrained language models, spanning models from a few million parameters (ALBERT) to eight billion parameters (DeepSeek-R1). These models cover languages such as English, Japanese, Chinese, German, Dutch, Spanish, and many of which are multilingual. The complete list of pretrained models can be found in Appendix E. We did not considered any Mixture-of-Experts (MoE) models, as they are usually more challenging to fine-tune and more prone to overfitting (Fedus et al., 2022; Shen et al., 2024).

**Fine-tuning Setup** We adopted full-model fine-tuning, which updates all parameters of the pretrained models. Although parameter-efficient methods such as LoRA (Hu et al., 2022) or layerfreezing strategies can significantly reduce computational overhead, full fine-tuning often leads to better downstream performance (Zhang et al., 2024; Shuttleworth et al., 2024). All models were fine-tuned under identical hyperparameter settings. See Appendix B for details.

**Reinforcement Learning Setup** We used the following libraries: Stable-Baselines (SB3) for PPO implementation and invalid action masking (Raffin et al., 2021), and the Gymnasium library for standardized environment definition (Towers et al., 2024).

**Objectives** For the optimization criteria, we focus on two objectives: test loss (measured via cross-entropy) and training time required for completing the fine-tuning. Test loss is a widely accepted proxy for task-specific performance, and training time serves as a practical and measurable approximation for other metrics, such as model size, inference throughput, deployment cost, etc. These choices are not fixed for LAMPS, as the framework is objective-agnostic. Hence, any measurable objectives can be used, as long as the corresponding metrics are recorded in the meta-dataset.

**Baselines** To our knowledge, no prior work has explored the same multi-objective optimization problem. Hence, a direct comparison with other existing methods was not possible. For this reason, we compared LAMPS with three basic baselines:

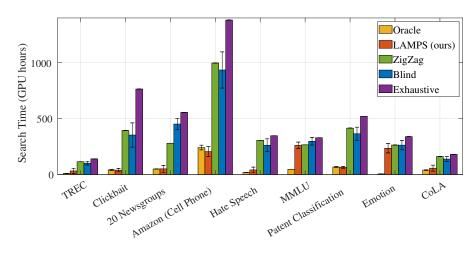


Figure 3: Mean time cost (in GPU hours) to reach 98% of the optimal hypervolume indicator on held-out datasets. For reference, we also show the time to complete an exhaustive search. On average, LAMPS reduces the search time by 71% compared to the exhaustive search, outperforming other feasible methods and being comparable to the ORACLE in 7 out of 9 datasets.

- **Blind**: chooses actions at random. Its performance serves as a lower bound on performance and represents the worst-case scenario.
- **Oracle**: assumes prior knowledge of the Pareto-optimal models for a given task. The performance of this approach represents the best-case scenario. In practice, this information is not available and serves only as a theoretical upper bound.
- **ZigZag**: a simple heuristic that sorts all candidate models by their number of parameters, then selects them in an alternating order (from largest to smallest and vice versa) in an attempt to quickly increase the covered hypervolume.

**Evaluation Method** To evaluate LAMPS's generalization, we employed leave-one-out cross-validation (Hastie et al., 2009), where one dataset is held exclusively for testing. For each fold, the policy is trained on the remaining datasets for a fixed number of steps and then evaluated on the held-out dataset. This allows us to assess how well the learned policy transfers to previously unseen tasks. To ensure robustness, this procedure was repeated five times, and we report the average performance across these runs.

#### 6.2 RESULTS

To evaluate the generalization of LAMPS to unseen datasets, Figure 3 reports the time required to reach 98% of the optimal hypervolume in each held-out dataset. Recall that, in our problem formulation, achieving optimal hypervolume corresponds to identifying all Pareto-optimal models. For reference, we also include the time needed for an exhaustive search to complete. Across all held-out datasets, LAMPS consistently outperforms all other feasible baselines and matches the performance of ORACLE in 7 out of 9 datasets. On average, this translates into a 71% reduction in search time. To illustrate the practical implications, consider the Amazon dataset: running an exhaustive search on a single A100 40GB NVIDIA GPU (\$3.67 hourly) would cost \$5,069.37, whereas LAMPS reduces it to just \$754.55 (an 85% reduction) with only a 2% degradation in the hypervolume.

Figure 4 provides further insight by tracking the progression of the hypervolume over search time. For comparability, hypervolume values are normalized by the maximum hypervolume, and we report the *hypervolume loss* (1 – normalized hypervolume) in logarithmic scale to highlight when the policy reaches optimality. Although LAMPS does not always reach optimality in a timely manner (compared to the other baselines), it clearly achieves near-optimal solutions quickly, eventually faster than ORACLE. This ability to deliver high-quality solutions at a fraction of the cost makes

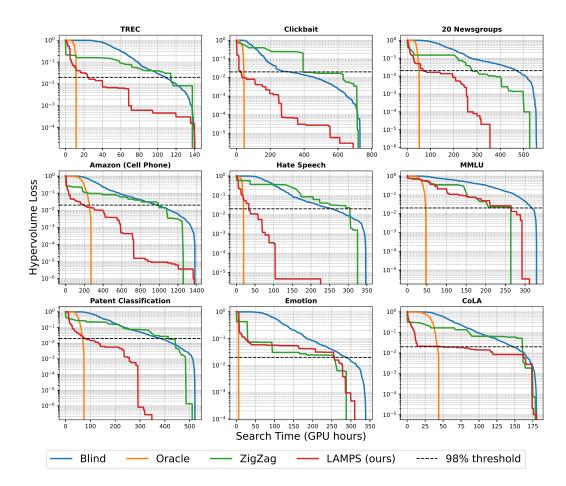


Figure 4: Evolution of the hypervolume indicator on held-out datasets as a function of search budget. LAMPS rapidly identifies near-optimal solutions (dashed line) in seven out of nine cases, demonstrating strong generalization capabilities, even when trained on a small meta-dataset.

LAMPS the best trade-off between efficiency and solution quality, positioning it as a pragmatic and strong tool for practitioners.

Moreover, in multi-objective applications, the end user must ultimately select a preferred solution from the Pareto front, often revisiting trade-offs as requirements, constraints, or business priorities evolve. By quickly providing a diverse set of strong candidates, LAMPS not only accelerates the search, but also enables practitioners to reconsider or change their choice later without having to undergo another expensive search, offering both flexibility and long-term practical value.

## 7 CONCLUSION

We presented LAMPS, a novel and open-source AutoML framework for efficiently selecting pretrained language models for fine-tuning, framing it as a multi-objective optimization problem. By combining landmark fine-tuning and meta-learning via reinforcement learning, LAMPS significantly reduces search costs while maintaining near-optimal performance. Experiments show that LAMPS reduces search time by 71% on average with minimal hypervolume degradation. To our knowledge, this is the first framework to deliver Pareto-efficient selection and fine-tuning for LLMs, establishing a new baseline for cost-aware AutoML and paving the way toward sustainable, high-performance deployment of foundation models.

# REFERENCES

- Mathias M. Adankon and Mohamed Cheriet. Model selection for the ls-svm. application to hand-writing recognition. *Pattern Recognition*, 42(12):3264–3270, 2009. ISSN 0031-3203. doi: https://doi.org/10.1016/j.patcog.2008.10.023. URL https://www.sciencedirect.com/science/article/pii/S0031320308004494. New Frontiers in Handwriting Recognition.
- Armen Aghajanyan, Akshat Shrivastava, Anchit Gupta, Naman Goyal, Luke Zettlemoyer, and Sonal Gupta. Better fine-tuning by reducing representational collapse. *arXiv preprint arXiv:2008.03156*, 2020.
- Alain Biem. A model selection criterion for classification: application to hmm topology optimization. In *Proceedings of the Seventh International Conference on Document Analysis and Recognition Volume 1*, ICDAR '03, pp. 104, USA, 2003. IEEE Computer Society. ISBN 0769519601.
- Hamparsum Bozdogan. Model selection and akaike's information criterion (aic): The general theory and its analytical extensions. *Psychometrika*, 52(3):345–370, 1987. doi: 10.1007/BF02294361. URL https://doi.org/10.1007/BF02294361.
- Pavel Brazdil and Christophe Giraud-Carrier. Metalearning and algorithm selection: progress, state of the art and introduction to the 2018 special issue. *Machine learning*, 107:1–14, 2018.
- Pavel Brazdil, Jan N. van Rijn, Carlos Soares, and Joaquin Vanschoren. *Dataset Characteristics (Metafeatures)*, pp. 53–75. Springer International Publishing, Cham, 2022. ISBN 978-3-030-67024-5. doi: 10.1007/978-3-030-67024-5\_4. URL https://doi.org/10.1007/978-3-030-67024-5\_4.
- Pavel B. Brazdil, Carlos Soares, and Joaquim Pinto Da Costa. Ranking learning algorithms: using ibl and meta-learning on accuracy and time results. *Mach. Learn.*, 50(3):251–277, March 2003. ISSN 0885-6125. doi: 10.1023/A:1021713901879. URL https://doi.org/10.1023/A:1021713901879.
- Abhijnan Chakraborty, Bhargavi Paranjape, Sourya Kakarla, and Niloy Ganguly. Stop clickbait: Detecting and preventing clickbaits in online news media. In *Advances in Social Networks Analysis and Mining (ASONAM)*, 2016 IEEE/ACM International Conference on, pp. 9–16. IEEE, 2016.
- Olivier Chapelle, Vladimir Vapnik, and Yoshua Bengio. Model selection for small sample regression. *Machine Learning*, 48(1):9–23, 2002. doi: 10.1023/A:1013943418833. URL https://doi.org/10.1023/A:1013943418833.
- Xi Chen, Ali Ghadirzadeh, Mårten Björkman, and Patric Jensfelt. Meta-learning for multi-objective reinforcement learning. In 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 977–983. IEEE Press, 2019. doi: 10.1109/IROS40897.2019.8968092. URL https://doi.org/10.1109/IROS40897.2019.8968092.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113, 2023.
- Thomas Davidson, Dana Warmsley, Michael Macy, and Ingmar Weber. Automated hate speech detection and the problem of offensive language. In *Proceedings of the 11th International AAAI Conference on Web and Social Media*, ICWSM '17, pp. 512–515, 2017.
- Lucas B. V. de Amorim, George D. C. Cavalcanti, and Rafael M. O. Cruz. Meta-scaler: A meta-learning framework for the selection of scaling techniques. *IEEE Transactions on Neural Networks and Learning Systems*, 36(3):4805–4819, 2025. doi: 10.1109/TNNLS.2024.3366615.
- DeepSeek-AI et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.

- Gianna M. Del Corso, Antonio Gullí, and Francesco Romani. Ranking a stream of news. In *Proceedings of the 14th International Conference on World Wide Web*, WWW '05, pp. 97–106, New York, NY, USA, 2005. Association for Computing Machinery. ISBN 1595930469. doi: 10.1145/1060745.1060764. URL https://doi.org/10.1145/1060745.1060764.
  - 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.
  - Romain Egele, Isabelle Guyon, Yixuan Sun, and Prasanna Balaprakash. Is one epoch all you need for multi-fidelity hyperparameter optimization? *arXiv preprint arXiv:2307.15422*, 2023.
  - Michael Emmerich, Nicola Beume, and Boris Naujoks. An emo algorithm using the hypervolume measure as selection criterion. In *Proceedings of the Third International Conference on Evolutionary Multi-Criterion Optimization*, EMO'05, pp. 62–76, Berlin, Heidelberg, 2005. Springer-Verlag. ISBN 3540249834. doi: 10.1007/978-3-540-31880-4\_5. URL https://doi.org/10.1007/978-3-540-31880-4\_5.
  - Armin Farhadi, Roya Hatami, Mohammad Robat Mili, Christos Masouros, and Mehdi Bennis. A meta-learning approach for energy-efficient resource allocation and antenna selection in star-bdris aided wireless networks. *IEEE Wireless Communications Letters*, pp. 1–1, 2025. doi: 10. 1109/LWC.2025.3543780.
  - William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: scaling to trillion parameter models with simple and efficient sparsity. *J. Mach. Learn. Res.*, 23(1), January 2022. ISSN 1532-4435.
  - Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning Volume 70*, ICML'17, pp. 1126–1135. JMLR.org, 2017.
  - Johannes Fürnkranz and Johann Petrak. An evaluation of landmarking variants. In *Working Notes of the ECML/PKDD 2000 Workshop on Integrating Aspects of Data Mining, Decision Support and Meta-Learning*, pp. 57–68. Citeseer, 2001.
  - Johannes Fürnkranz, Josef Petrak, Pavel Brazdil, and Carlos Soares. On the use of fast subsampling estimates for algorithm recommendation. Technical report, Austrian Research Institute for Artificial Intelligence, 2002.
  - Alec Go, Richa Bhayani, and Lei Huang. Twitter sentiment classification using distant supervision. *CS224N project report, Stanford*, 1(12):2009, 2009.
  - Andreia P. Guerreiro, Carlos M. Fonseca, and Luís Paquete. The hypervolume indicator: Computational problems and algorithms. *ACM Comput. Surv.*, 54(6), July 2021. ISSN 0360-0300. doi: 10.1145/3453474. URL https://doi.org/10.1145/3453474.
  - Anna Gutowska. What are AI agents?, July 2024. URL https://www.ibm.com/think/topics/ai-agents. Accessed: 2025-02-07.
  - Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining, inference, and prediction.* Springer, New York, 2 edition, 2009. ISBN 978-0-387-84857-0. URL https://hastie.su.domains/ElemStatLearn/.
  - Jason Hepburn. Universal language model fine-tuning for patent classification. In Sunghwan Mac Kim and Xiuzhen (Jenny) Zhang (eds.), *Proceedings of the Australasian Language Technology Association Workshop 2018*, pp. 93–96, Dunedin, New Zealand, December 2018. URL https://aclanthology.org/U18-1013/.
  - 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, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack William Rae, and

Laurent Sifre. An empirical analysis of compute-optimal large language model training. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=iBBcRU1OAPR.

- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9.
- Shengyi Huang and Santiago Ontañón. A closer look at invalid action masking in policy gradient algorithms. *The International FLAIRS Conference Proceedings*, 35, May 2022. ISSN 2334-0762. doi: 10.32473/flairs.v35i.130584. URL http://dx.doi.org/10.32473/flairs.v35i.130584.
- Nishant Jain, Arun S. Suggala, and Pradeep Shenoy. Improving generalization via meta-learning on hard samples. In 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 27590–27599, 2024. doi: 10.1109/CVPR52733.2024.02606.
- Kevin Jamieson and Ameet Talwalkar. Non-stochastic best arm identification and hyperparameter optimization. In Arthur Gretton and Christian C. Robert (eds.), *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pp. 240–248, Cadiz, Spain, 09–11 May 2016. PMLR. URL https://proceedings.mlr.press/v51/jamieson16.html.
- A. Kalousis and M. Hilario. Model selection via meta-learning: a comparative study. In *Proceedings* 12th IEEE Internationals Conference on Tools with Artificial Intelligence. ICTAI 2000, pp. 406–413, 2000. doi: 10.1109/TAI.2000.889901.
- 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. *OpenAI blog*, 2020. URL https://arxiv.org/abs/2001.08361.
- Aaron Klein, Stefan Falkner, Simon Bartels, Philipp Hennig, and Frank Hutter. Fast Bayesian Optimization of Machine Learning Hyperparameters on Large Datasets. In Aarti Singh and Jerry Zhu (eds.), *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pp. 528–536. PMLR, 20–22 Apr 2017. URL https://proceedings.mlr.press/v54/klein17a.html.
- Ken Lang. 20 newsgroups dataset, 1995. URL http://qwone.com/~jason/ 20Newsgroups/. Accessed: 2025-02-26.
- Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *Journal of Machine Learning Research*, 18(185):1–52, 2018. URL http://jmlr.org/papers/v18/16-558.html.
- Fei-Yu Liu and Chao Qian. Prediction guided meta-learning for multi-objective reinforcement learning. In 2021 IEEE Congress on Evolutionary Computation (CEC), pp. 2171–2178. IEEE Press, 2021. doi: 10.1109/CEC45853.2021.9504972. URL https://doi.org/10.1109/CEC45853.2021.9504972.
- Yang Liu, Jingxuan Wei, Xin Li, and Minghan Li. Generational distance indicator-based evolutionary algorithm with an improved niching method for many-objective optimization problems. *IEEE Access*, 7:63881–63891, 2019. doi: 10.1109/ACCESS.2019.2916634.
- Junlin Lu, Patrick Mannion, and Karl Mason. A meta-learning approach for multi-objective reinforcement learning in sustainable home environments. In *European Conference on Artificial Intelligence (ECAI)* 2024, January 2024.
- Wei Ma, Daoyuan Wu, Yuqiang Sun, Tianwen Wang, Shangqing Liu, Jian Zhang, Yue Xue, and Yang Liu. Combining fine-tuning and llm-based agents for intuitive smart contract auditing with justifications. *arXiv preprint arXiv:2403.16073*, 2024.

- Oded Maron and Andrew W. Moore. Hoeffding races: accelerating model selection search for classification and function approximation. In *Proceedings of the 7th International Conference on Neural Information Processing Systems*, NIPS'93, pp. 59–66, San Francisco, CA, USA, 1993. Morgan Kaufmann Publishers Inc.
  - Allan D R McQuarrie and Chih-Ling Tsai. Regression and time series model selection. WORLD SCIENTIFIC, 1998. doi: 10.1142/3573. URL https://www.worldscientific.com/doi/abs/10.1142/3573.
  - Marcio Monteiro, Charu Karakkaparambil James, Marius Kloft, and Sophie Fellenz. Characterizing text datasets with psycholinguistic features. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 14977–14990, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
  - Alex Nichol, Joshua Achiam, and John Schulman. On first-order metalearning algorithms, 2018. URL https://openai.com/index/ on-first-order-meta-learning-algorithms/.
  - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 27730–27744. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper\_files/paper/2022/file/blefde53be364a73914f58805a001731-Paper-Conference.pdf.
  - Bernhard Pfahringer, Hilan Bensusan, and Christophe G. Giraud-Carrier. Meta-learning by land-marking various learning algorithms. In *Proceedings of the Seventeenth International Conference on Machine Learning*, ICML '00, pp. 743–750, San Francisco, CA, USA, 2000. Morgan Kaufmann Publishers Inc. ISBN 1558607072.
  - Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
  - 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.
  - Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021. URL http://jmlr.org/papers/v22/20-1364.html.
  - Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, 2020. URL https://arxiv.org/abs/1910.01108.
  - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347.
  - Sheng Shen, Le Hou, Yanqi Zhou, Nan Du, Shayne Longpre, Jason Wei, Hyung Won Chung, Barret Zoph, William Fedus, Xinyun Chen, Tu Vu, Yuexin Wu, Wuyang Chen, Albert Webson, Yunxuan Li, Vincent Y Zhao, Hongkun Yu, Kurt Keutzer, Trevor Darrell, and Denny Zhou. Mixture-of-experts meets instruction tuning: A winning combination for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=6mLjDwYte5.
  - Reece Shuttleworth, Jacob Andreas, Antonio Torralba, and Pratyusha Sharma. Lora vs full fine-tuning: An illusion of equivalence. *arXiv preprint arXiv:2410.21228*, 2024.

- Carlos Soares, Johann Petrak, and Pavel Brazdil. Sampling-based relative landmarks: systematically test-driving algorithms before choosing. In Pavel Brazdil and Alípio Jorge (eds.), *Progress in Artificial Intelligence*, pp. 88–95, Berlin, Heidelberg, 2001. Springer Berlin Heidelberg. ISBN 978-3-540-45329-1.
- Shagun Sodhani, Amy Zhang, and Joelle Pineau. Multi-task reinforcement learning with context-based representations. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9767–9779. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/sodhani21a.html.
- Yee Whye Teh, Victor Bapst, Wojciech Marian Czarnecki, John Quan, James Kirkpatrick, Raia Hadsell, Nicolas Heess, and Razvan Pascanu. Distral: robust multitask reinforcement learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, pp. 4499–4509, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu, Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
- Yasmen Wahba, Nazim Madhavji, and John Steinbacher. A comparison of svm against pre-trained language models (plms) for text classification tasks. In Giuseppe Nicosia, Varun Ojha, Emanuele La Malfa, Gabriele La Malfa, Panos Pardalos, Giuseppe Di Fatta, Giovanni Giuffrida, and Renato Umeton (eds.), *Machine Learning, Optimization, and Data Science*, pp. 304–313, Cham, 2023. Springer Nature Switzerland. ISBN 978-3-031-25891-6.
- Qi Wang, Chengwei Zhang, and Bin Hu. Dynamic programming with meta-reinforcement learning: a novel approach for multi-objective optimization. *Complex & Intelligent Systems*, 10(4): 5743–5758, 2024. doi: 10.1007/s40747-024-01469-1. URL https://doi.org/10.1007/s40747-024-01469-1.
- Shangshang Wang, Julian Asilis, Omer Faruk Akgül, Enes Burak Bilgin, Ollie Liu, and Willie Neiswanger. Tina: Tiny reasoning models via lora, 2025. URL https://arxiv.org/abs/2504.15777.
- Zhilin Yang. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv* preprint arXiv:1906.08237, 2019.
- Biao Zhang, Zhongtao Liu, Colin Cherry, and Orhan Firat. When scaling meets LLM finetuning: The effect of data, model and finetuning method. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=5HCnKDeTws.
- Peng Zhao and Bin Yu. On model selection consistency of lasso. *Journal of Machine Learning Research*, 7(90):2541–2563, 2006. URL http://jmlr.org/papers/v7/zhao06a.html.

# A LAMPS: GETTING STARTED

This section demonstrates how to use LAMPS for a new, unseen dataset. The provided policy was meta-trained across all datasets described in the main paper for a total of 45M steps, minimizing the following objectives: test loss and training time. Before running it, make sure to have sufficient disk space (at least 2TB) for intermediate storage of models and checkpoints. In addition, some models hosted on Hugging Face may require license agreements or explicit acceptance terms. Ensure that the necessary access is granted to your user account prior to execution.

Listing 1: Running LAMPS for a new dataset.

```
# Create the Python environment
conda create -n lamps python=3.12
conda activate lamps

# Install dependencies
pip install -r requirements.txt

# Initiate the search using the trained policy
python eval.py --policy "policies/ALL-DISTRAL-45M_steps.zip" \
    --dataset "stanfordnlp/imdb" \
    --input-col "text" \
    --target-col "label"
```

# B HYPERPARAMETERS

#### B.1 FINE-TUNING

We used the Trainer module from Hugging Face's **transformers** library for fine-tuning. The key hyperparameters and settings were as follows:

• Optimizer: AdamW • Learning rate:  $7 \times 10^{-6}$ 

• Batch size: Automatically determined based on available hardware

Early stopping patience: 3 epochsMixed precision: Enabled (BF16)

All unspecified settings followed the default values defined in Trainer module.

#### B.2 PPO

For the PPO algorithm, we used the implementation from Stable Baselines3 library. The key hyper-parameters and settings were as follows:

Learning rate: 1 × 10<sup>-4</sup>
 Minibatch size: 256
 Num. epochs: 15
 Discount (γ): 0.99
 GAE parameter (λ): 0.97
 Clip range: 0.20
 VF coeff. c<sub>1</sub>: 0.5

• Entropy coeff.  $c_2$ : 0.23

All policies were trained using the Gymnasium environment API with invalid action masking.

# C CONVERGENCE ANALYSIS OF THE REWARDS

In order to provide additional evidence that the reward function defined in equation 7 effectively guides the agent toward the optimal solution set, according to the original multi-objective problem in equation 1, Figure 5 presents a typical reward evolution observed during training, for both single task and multi-task RL (MTRL) using PPO.

For better interpretability and comparison, reward values are normalized such that a value of 3000 corresponds to the optimal reward, when the agent exclusively evaluates Pareto-optimal solutions, achieving maximal hypervolume in minimal time. The learned policy exhibits a consistent upward trend in reward, eventually converging to the optimal value.

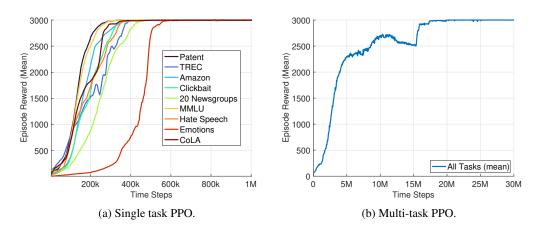


Figure 5: Normalized reward progression during policy training using PPO algorithm. As expected, multi-task RL takes longer, but it also converges to the optimal reward.

# D DATASETS

This section describes the datasets used in our experiments. When a dataset did not provide predefined training and test splits, we reserved 20% of the data for testing. Not that all datasets correspond to text classification tasks. Although this focus may appear restrictive, they span a broad spectrum of linguistic competencies and thus form a diverse, challenging testbed for rigorously assessing the generalization and transfer learning capabilities of LAMPS across various tasks and linguistic demands.

**TREC** A classic question classification benchmark with 6 coarse-grained classes (e.g., abbreviation, entity, description and abstract concept, human being, location, and numeric value). Task: Question classification. License: N/A (widely used academic benchmark; originally from UIUC).

**Clickbait** Contains news headlines labeled as either "clickbait" or "non-clickbait". Derived from social media posts (Chakraborty et al., 2016). Task: Binary classification. License: N/A.

**20 Newsgroups** A collection of 20,000 newsgroup emails across 20 different topics (Lang, 1995). Task: Topic classification. License: CC BY 4.0.

**Amazon Reviews (cell-phone)** Subset of the Amazon Product Review 2013 dataset, filtered for the "Cell Phone reviews" category. Includes star ratings from 1 to 5 and contains 78,930 reviews. Task: Sentiment classification (5 classes). License: N/A (Amazon public data, widely used in academia).

**Hate Speech and Offensive Language** A corpus of over 24,000 tweets manually annotated as hate speech, offensive but not hateful, or neither (Davidson et al., 2017). Task: Offensive language classification (3 classes). License: MIT License.

(academic benchmark from the GLUE suite).

864

865

866

867 868

870 871

872

873 874

875

876

877

MIT License.

USPTO data).

878 Ε PRETRAINED LANGUAGE MODELS 879 880 Below is the list of pretrained models used during all experiments of this paper: 883 BERT: 1. google-bert/bert-large-cased-whole-word-masking 885 2. google-bert/bert-large-uncased-whole-word-masking-fine-tuned-squad 3. google-bert/bert-large-uncased-whole-word-masking 4. google-bert/bert-large-uncased 888 5. google-bert/bert-large-cased-whole-word-masking-fine-tuned-squad 889 6. google-bert/bert-large-cased 890 7. google-bert/bert-base-uncased 891 8. google-bert/bert-base-multilingual-uncased 892 9. google-bert/bert-base-multilingual-cased 10. google-bert/bert-base-german-dbmdz-uncased 893 11. google-bert/bert-base-german-dbmdz-cased 894 12. google-bert/bert-base-german-cased 895 13. google-bert/bert-base-chinese 14. google-bert/bert-base-cased 897 GPT: 899 900 1. openai-community/gpt2 901 2. openai-community/gpt2-medium 902 3. openai-community/gpt2-large 903 4. openai-community/gpt2-x1 904 RoBERTa: 905 906 1. FacebookAI/roberta-base 907 2. FacebookAI/roberta-large 908 3. FacebookAI/xlm-roberta-base 909 4. FacebookAI/xlm-roberta-large 910 5. FacebookAI/xlm-roberta-large-fine-tuned-conll02-dutch 911 6. FacebookAI/xlm-roberta-large-fine-tuned-conll02-spanish 912 7. FacebookAI/xlm-roberta-large-fine-tuned-conll03-english 913 8. FacebookAI/xlm-roberta-large-fine-tuned-conll03-german 914 915 OPT: 916 1. facebook/opt-125m 917 2. facebook/opt-350m

MMLU Massive Multitask Language Understanding, a benchmark covering 57 diverse subject

areas from elementary math to law and philosophy. Task: Multi-choice question answering. License:

**Patent Classification** Consisting of 35, 000 Patent abstracts labeled with Cooperative Patent Classification (CPC) codes (9 classes). Task: Topic classification. License: Public domain (based on

**Emotion** A dataset of 20,000 Twitter messages in English annotated with one of six basic emo-

tions (anger, fear, joy, love, sadness, surprise). Task: Emotion classification. License: MIT License.

**CoLA** Corpus of Linguistic Acceptability, a dataset of English sentences labeled as grammati-

cally acceptable or unacceptable. Task: Acceptability classification (binary). License: Unknown

918	3. facebook/opt-1.3b
919	4. facebook/opt-2.7b
920	5. facebook/opt-6.7b
921	1
922	Llama:
923	
924	1. meta-llama/Llama-3.2-1B
925	2. meta-llama/Llama-3.2-1B-Instruct
926	3. meta-llama/Llama-3.2-3B 4. meta-llama/Llama-3.1-8B
927	4. Ilicta-nama/Liama-3.1-6D
928	DistilBERT:
929	Distribution
930	<ol> <li>distilbert/distilbert-base-multilingual-cased</li> </ol>
931	2. distilbert/distilbert-base-german-cased
932	3. distilbert/distilbert-base-uncased-distilled-squad
933	4. distilbert/distilbert-base-cased-distilled-squad
934	<ol> <li>distilbert/distilbert-base-cased</li> <li>distilbert/distilbert-base-uncased</li> </ol>
935	7. distilbert/distilroberta-base
936	8. distilbert/distilgpt2
937	o. distributed distribute
938	ALBERT:
939	
940	1. albert/albert-xlarge-v2
941	2. albert/albert-xxlarge-v2
942	3. albert/albert-xxlarge-v1
943	4. albert/albert-xlarge-v1
944	<ul><li>5. albert/albert-large-v2</li><li>6. albert/albert-large-v1</li></ul>
945	7. albert/albert-base-v2
946	8. albert/albert-base-v1
947	
948	LUKE:
949	
950	1. studio-ousia/mluke-large
951	<ol> <li>studio-ousia/mluke-large-lite</li> <li>studio-ousia/mluke-base-lite</li> </ol>
952	4. studio-ousia/mluke-base
953	5. studio-ousia/luke-japanese-base
954	6. studio-ousia/luke-japanese-base-lite
955	7. studio-ousia/luke-japanese-large-lite
956	8. studio-ousia/luke-japanese-large
957	9. studio-ousia/luke-large-lite
958	10. studio-ousia/luke-base-lite
959	
000	11. studio-ousia/luke-large
960	11. studio-ousia/luke-large 12. studio-ousia/luke-base
961	12. studio-ousia/luke-base
961 962	
961 962 963	12. studio-ousia/luke-base
961 962 963 964	12. studio-ousia/luke-base  DeepSeek:  1. deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B 2. deepseek-ai/DeepSeek-R1-Distill-Qwen-7B
961 962 963 964 965	<ul><li>12. studio-ousia/luke-base</li><li>DeepSeek:</li><li>1. deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B</li></ul>
961 962 963 964 965 966	12. studio-ousia/luke-base  DeepSeek:  1. deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B 2. deepseek-ai/DeepSeek-R1-Distill-Qwen-7B 3. deepseek-ai/DeepSeek-R1-Distill-Llama-8B
961 962 963 964 965 966 967	12. studio-ousia/luke-base  DeepSeek:  1. deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B 2. deepseek-ai/DeepSeek-R1-Distill-Qwen-7B
961 962 963 964 965 966 967 968	12. studio-ousia/luke-base  DeepSeek:  1. deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B 2. deepseek-ai/DeepSeek-R1-Distill-Qwen-7B 3. deepseek-ai/DeepSeek-R1-Distill-Llama-8B  Qwen:
961 962 963 964 965 966 967 968 969	12. studio-ousia/luke-base  DeepSeek:  1. deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B 2. deepseek-ai/DeepSeek-R1-Distill-Qwen-7B 3. deepseek-ai/DeepSeek-R1-Distill-Llama-8B  Qwen: 1. Qwen/Qwen2.5-0.5B
961 962 963 964 965 966 967 968	12. studio-ousia/luke-base  DeepSeek:  1. deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B 2. deepseek-ai/DeepSeek-R1-Distill-Qwen-7B 3. deepseek-ai/DeepSeek-R1-Distill-Llama-8B  Qwen:

# F ADDING NEW MODELS TO THE META-DATASET

To incorporate a new model into the recommendation pool of LAMPS, it must first be integrated into the meta-dataset. We refer to this process as *model fingerprinting*. Because LAMPS relies on meta-learning, it is necessary to observe the actual performance of the new model on known datasets before the system can generalize its behavior to unseen datasets. This integration requires two steps:

- The new LLM must be fine-tuned on all datasets currently included in the meta-dataset, with all relevant metrics recorded.
- 2. The reinforcement learning policy must be retrained on the expanded meta-dataset.

Currently, complete retraining is the recommended procedure for reliable integration of new models. Although incremental training strategies could further reduce the computational overhead, the cost of full retraining is already negligible compared to the fine-tuning runs required to expand the metadataset.

The ideal number and diversity of datasets in the meta-dataset remains an open research question. A smaller set of datasets facilitates the addition of new models, since each integration requires fewer fine-tuning runs. Conversely, a larger and more diverse collection typically improves the generalization ability of the learned policy to unseen tasks. How to balance these competing goals remains an open challenge for future work.

# G Proof of Theorem 1

*Proof.* We first prove that the maximizer  $X_{\lambda} = \arg\max_{X \subset \mathcal{X}} H_{\mathcal{D}}(X,r) - \lambda |X|$  is a subset of Pareto solutions  $X^*$ , that is, for any  $\lambda > 0$ ,  $X_{\lambda} \subset X^*$ . This is proved by contradiction. Suppose that there exists a  $x \in X_{\lambda}$  that is not Pareto-optimal. Then, there exists a  $x^* \in \mathcal{X}$  dominating x such that  $\Lambda([x,r]) < \Lambda([x^*,r])$  holds. Denote  $X_{\lambda}^*$  the set obtained from  $X_{\lambda}$  by replacing x with  $x^*$ . By the definition of  $H_{\mathcal{D}}$ , we know  $H_{\mathcal{D}}(X_{\lambda},r) < H_{\mathcal{D}}(X_{\lambda}^*,r)$ . Then, it holds

$$H_{\mathcal{D}}(X_{\lambda}, r) - \lambda |X_{\lambda}| = H_{\mathcal{D}}(X_{\lambda}, r) - \lambda |X_{\lambda}^{*}| < H_{\mathcal{D}}(X_{\lambda}^{*}, r) - \lambda |X_{\lambda}^{*}|.$$

This contradicts the assumption that  $X_{\lambda}$  is the maximizer of problem (3). Hence, for any  $\lambda>0$ , we know  $X_{\lambda}\subset X^*$ . Below we prove the if part and the only if part respectively. **The** if part: In this part, we prove that if equation 4 holds, then the optimal solution  $X^*\subset \mathcal{X}$  of problem equation 3 contains only and exactly all Pareto-optimal solutions. Let  $X_{\lambda}=\arg\max_{X\subset\mathcal{X}}H_{\mathcal{D}}(X,r)-\lambda|X|$ . From the above discussion we know  $X_{\lambda}\subset X^*$ . Suppose  $|X^*|-|X_{\lambda}|=s$ . We denote  $\{x_{i_1},\ldots,x_{i_s}\}\subset X^*$  the subset of  $X^*$  such that  $\{x_{i_1},\ldots,x_{i_s}\}\cap X_{\lambda}=\varnothing$ . We define  $X_k=X_{\lambda}\cup\{x_{i_1},\ldots,x_{i_k}\}$  for all  $k\in\{0,1,\ldots,s\}$ . Then, we know  $X_s=X^*$ ,  $X_0=X_{\lambda}$  and  $|X_{k+1}|-|X_k|=1$  for all  $k\in\{0,1,\ldots,s-1\}$ . Note that

$$H_{\mathcal{D}}(X^*, r) - \lambda |X^*| - \left(H_{\mathcal{D}}(X_{\lambda}, r) - \lambda |X_{\lambda}|\right)$$

$$= H_{\mathcal{D}}(X^*, r) - H_{\mathcal{D}}(X_{\lambda}, r) - \lambda \left(|X^*| - |X_{\lambda}|\right)$$

$$= H_{\mathcal{D}}(X_s, r) - H_{\mathcal{D}}(X_0, r) - s\lambda$$

$$= \sum_{k=1}^{s} \left(H_{\mathcal{D}}(X_k, r) - H_{\mathcal{D}}(X_{k-1}, r) - \lambda\right)$$

$$\geq \sum_{k=1}^{s} \left(H_{\mathcal{D}}(X_k, r) - H_{\mathcal{D}}(X_{k-1}, r) - \min_{x \in X, X \subset X^*} \Delta H_{\mathcal{D}}(x|X)\right)$$

$$\geq 0, \tag{8}$$

where the last second inequality used equation 4 and the last inequality used the definition of  $\min_{x \in X, X \subset X^*} \Delta H_{\mathcal{D}}(x|X)$ . **The only if part:** To prove this part of the result, we only need to show that there exists an optimization problem whose Pareto solution set  $X^*$  with  $|X^*| = s$  satisfies that for any sequence of subsets  $\{X_i\}_{i=1}^s$  satisfying  $X_i \subset X^*$  and  $|X_i| = i$ , it holds

$$\max_{i \in \{2, \dots, s\}} H_{\mathcal{D}}(X_i, r) - H_{\mathcal{D}}(X_{i-1}, r) = \min_{x \in X, X \subset X^*} \Delta H_{\mathcal{D}}(x|X). \tag{9}$$

On the other hand, from equation 8 and the definition of  $X^*$  we know

$$H_{\mathcal{D}}(X^*,r) - \lambda |X^*| - \left(H_{\mathcal{D}}(X_{\lambda},r) - \lambda |X_{\lambda}|\right) = \sum_{k=1}^{s} \left(H_{\mathcal{D}}(X_k,r) - H_{\mathcal{D}}(X_{k-1},r)\right) - s\lambda \ge 0.$$

Hence, we have  $\sum_{k=1}^{s} (H_{\mathcal{D}}(X_k, r) - H_{\mathcal{D}}(X_{k-1}, r)) \ge s\lambda$ . Combining this observation with equation 9 together, we get

$$\min_{x \in X, X \subset X^*} \Delta H_{\mathcal{D}}(x|X) \ge \lambda.$$

The proof is completed by noting that equation 9 always holds for arbitrary  $X^*$  with  $|X^*| = 2$ .  $\square$ 

# H THE USE OF LARGE LANGUAGE MODELS (LLMS)

We used large language models solely for surface-level editing: spelling and grammar correction, and minor wording improvements. LLMs were *not* used for idea generation, experiment design, data analysis, coding, mathematical derivations, or substantive content creation.