

# TRAINING MICE TO COMPETE WITH ELEPHANTS: A GUIDE FOR CUSTOMIZING SMALL-SIZED LLMs ON KNOWLEDGE AND SKILLS DATA

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Customizing large language models (LLMs) is increasingly in demand by enterprises and individual developers. It allows LLMs to be tailored for domain expertise, aligned with organizational guidelines, and enhanced for user experience. Effective customization hinges on three core elements: a small-size model, large-scale domain-specific datasets, and an effective training strategy to help the model acquire relevant knowledge and skills from the data. In this paper, we focus on the third element by conducting an in-depth study on fine-tuning LLMs (3B to 7B parameters) using large-scale instruction tuning datasets across multiple knowledge domains and skills. We examine various training configurations and strategies on three pretrained LLMs. Our results question several common training practices, including hyperparameter recommendations from TULU and phased training recommended by Orca.

Key insights from our work include: (i) larger batch sizes paired with lower learning rates lead to improved model performance on benchmarks such as MMLU, MTBench, and Open LLM Leaderboard; (ii) early-stage training dynamics, such as lower gradient norms and higher loss values, are strong indicators of better final model performance, allowing for early termination of sub-optimal runs and significant computational savings; (iii) skipping warmup and using a constant learning rate do not compromise performance; and (iv) stacked training outperforms phased training. With these findings holding robustly across model families and sizes, we hope this study serves as a comprehensive guide for practitioners fine-tuning small LLMs.

## 1 INTRODUCTION

Large language models (LLMs) are growing in size, but bigger is not always better. Small-sized LLMs (3B to 7B parameters) are becoming the backbone of enterprise AI systems due to their adaptability and efficiency (RedHat, 2024; Zhang et al., 2023; Lee, 2024). Customizing these models allows them to be tailored for specific tasks, domains, or organizational needs. Compared to larger LLMs, fine-tuning and deploying these models is faster, more cost-effective, and does not require specialized infrastructure or extensive hardware like GPUs and TPUs. Moreover, they can be hosted on consumer-grade machines, making them accessible to smaller organizations and individual researchers while delivering performance comparable to much larger models on many specialized tasks. Importantly, these models offer users full control over their data and fine-tuned versions, reducing the risk of data breaches or non-compliance with regulations such as GDPR (2016).

Instruction tuning small LLMs using large-scale domain-specific data is an effective method for enhancing domain knowledge and skills. Our work focuses on supervised fine-tuning (SFT), as it is the most widely used approach for instruction tuning and aligns with our goal of customizing LLMs using knowledge and skills data. Traditionally, instruction tuning has focused on enabling LLMs to follow user instructions and improve zero-shot capabilities (Ouyang et al., 2022). More recently, it has been used to fine-tune LLMs for customization by leveraging large-scale knowledge and skills instruction datasets. For example, a telecommunications company might customize LLMs for customer service tasks. In this case, skills data, such as annotated conversations and summarized customer support logs, can improve the model’s ability to classify interactions, detect sentiment, and

054 adjust its language. Meanwhile, knowledge data, such as device troubleshooting guides and service  
055 plans, can teach the model domain-specific information. These customized knowledge and skills  
056 datasets are often accessible and of high quality (Gunasekar et al., 2023). Additionally, a growing  
057 body of work is focused on generating large-scale knowledge and skills instruction data (Sudalairaj  
058 et al., 2024; Wang et al., 2023b; Taori et al., 2023; Xu et al., 2023; Li et al., 2024). However, there  
059 is limited research on how to effectively fine tune small-sized LLMs on large-scale knowledge and  
060 skills data.

061 Practitioners have limited resources to reference when searching for optimal training strategies and  
062 hyper-parameters for instruction-tuning small LLMs on knowledge and skills data. Many LLMs are  
063 closed-source, and even those that are open-source often lack detailed technical reports describing  
064 how to set up hyper-parameters or which configurations were attempted but unsuccessful. As a  
065 result, critical factors like batch size and learning rate, as well as their impact on final model per-  
066 formance, remain unclear. Additionally, phase training is increasingly used for instruction tuning,  
067 where LLMs are fine-tuned progressively, starting with simple instruction-following data (e.g., gen-  
068 eral knowledge from elementary or middle school), then moving to foundational knowledge (e.g.,  
069 graduate-level content), and finally to skills-based data. However, it is unclear how well phase train-  
070 ing outperforms traditional stacked training where all data is combined into a single phase. Ident-  
071 ifying an effective set of hyper-parameters is especially difficult for users with limited computational  
072 resources. This motivates the main question we aim to study:

073 *How can we effectively fine-tune a small-size LLM (3B–7B parameters) on large-scale instruction*  
074 *tuning datasets that cover diverse knowledge and skills?*  
075

076 In this paper, we present a comprehensive empirical study on fine-tuning small-size LLMs and com-  
077 pare our findings with existing research on this topic. We experiment with 3 open-source models—  
078 Granite 3B, Granite 7B, LLaMA 3.2 3B, and Mistral 7B—fine-tuning them on five  
079 datasets: an instruction-following dataset with 308,343 samples, a foundational knowledge dataset  
080 with 231,178 samples, a complex skills dataset with 285,966 samples, the TULU mixture v2 dataset,  
081 and a domain-specific math, reasoning, and coding dataset. Through a series of experiments, we sys-  
082 tematically vary hyper-parameters and training strategies and collect experimental results. Our find-  
083 ings challenge several widely accepted practices, including those recommended by TULU (Wang  
084 et al., 2023a; Ivison et al., 2023), which is often considered a gold standard for LLM fine-tuning. For  
085 example, they use a (small) batch size of 128 samples, which we find to underperform in our experi-  
086 ments. We conjecture that this choice was driven by their computational constraints, as larger batch  
087 sizes can produce models with higher downstream performance but require much longer training  
088 time under limited computing resources. Additionally, while learning rate schedulers with warm-up  
089 and cosine decay are widely used in neural network training, including in TULU, our results show  
that these techniques have minimal impact on model downstream performance.

090 Our key observations are: (i) larger batch sizes combined with lower learning rates improve gener-  
091 alization and performance on benchmarks like MMLU (Hendrycks et al., 2020), MTBench (Zheng  
092 et al., 2023), and Open LLM Leaderboard v2; (ii) early-stage training dynamics, such as lower  
093 gradient norms and higher loss values, are strong indicators of final model performance, enabling  
094 early termination of sub-optimal runs and significant computational savings; (iii) omitting warmup  
095 steps and using constant learning rates does not compromise performance; and (iv) stacked train-  
096 ing frequently outperforms phased training. We also address adaptations for new architectures and  
097 emphasize the importance of efficient data handling techniques, such as bucketing and balanced  
098 compute distribution across GPUs. Our findings aim to provide practitioners with actionable in-  
099 sights to fine-tune LLMs more effectively, optimizing performance while simplifying the training  
100 process. This can benefit the open-source community focused on instruction tuning and serve as a  
101 reference for practitioners with limited computational resources.

## 102 RELATED WORK

103 **Instruction Tuning Data.** Instruction tuning with diverse, large-scale datasets can effectively im-  
104 prove LLM performance across downstream tasks (Wang et al., 2023b; Honovich et al., 2023; Chung  
105 et al., 2024; Isik et al., 2024; Cheng et al., 2024). Recent studies have found that large-scale instruc-  
106 tion tuning data focusing on knowledge and skills is particularly beneficial for adapting LLMs to  
107 customized domains or applications, improving factual recall and reducing hallucinations (Cheng

108 et al., 2023; Allen-Zhu & Li, 2023; Yang et al., 2024). This observation has led to a growing body  
109 of research introducing novel instruction tuning datasets. For instance, several works leveraged  
110 larger, more powerful LLMs (e.g., ChatGPT (OpenAI, 2023; 2022) and Mistral models (Jiang et al.,  
111 2023; 2024)) to distill instruction tuning data from them using seed examples provided by users  
112 (Mitra et al., 2024; Xu et al., 2023; Ding et al., 2023; Peng et al., 2023; Mukherjee et al., 2023).  
113 GLAN (Li et al., 2024) and LAB (Sudalairaj et al., 2024) further advanced this area by proposing  
114 taxonomy-driven frameworks to enhance the diversity of synthetic instruction tuning data. Building  
115 on these datasets, many studies explored strategies to optimize dataset composition, select representa-  
116 tive data subsets, and evaluate data quality before incorporating them into model training (Iverson  
117 et al., 2023; Liu et al., 2023; Li et al., 2023; Xie et al., 2023). While these advancements have driven  
118 rapid progress in instruction-tuned LLMs, limited work has focused on how to effectively use such  
119 data during training to achieve optimal performance, or how training outcomes vary with different  
120 compute budgets (e.g., GPUs and TPUs). In this paper, we fill this gap by conducting an extensive  
121 set of experiments to investigate various training strategies and hyperparameters for customizing  
122 small LLMs on these datasets, analyzing how different configurations interact with available com-  
123 pute resources to impact the downstream performance of fine-tuned models.

124  
125  
**Training Dynamics.** Training configurations and hyper-parameter setups play a pivotal role in  
126 training LLMs, as they directly influence model performance, convergence stability, and resource  
127 efficiency. Most research has focused on the pre-training phase, as it is the most resource-intensive  
128 part of LLM development (Yang et al., 2022; Hägele et al., 2024; Bi et al., 2024; Kaplan et al., 2020;  
129 Rosenfeld et al., 2019; Gunter et al., 2024; Dubey et al., 2024). For example, Sardana et al. (2024);  
130 Hoffmann et al. (2022) introduced scaling laws to determine optimal model sizes for given datasets  
131 and Hägele et al. (2024) proposed novel learning rate schedulers as alternatives to conventional cosine  
132 decay. Additionally, recent research proposed to incorporate instruction tuning data alongside  
133 pre-training data as part of a decay phase in pre-training, linking to the body of research on con-  
134 tinual pre-training (Hu et al., 2024; Ibrahim et al., 2024; Lesort et al., 2021; Scialom et al., 2022).  
135 In contrast, our work shifts the focus to customizing pre-trained LLMs through instruction tuning,  
136 highlighting under-explored challenges in training strategies and hyper-parameter configurations for  
137 this stage. Many instruction tuning studies either omit the reporting of hyperparameters altogether  
138 (Mukherjee et al., 2023) or only provide a selective set of hyperparameters used in successful runs  
139 (Wang et al., 2023a; Iverson et al., 2023; Xu et al., 2023), often without disclosing failed experiments  
140 or alternative configurations explored during their research. In contrast, we conduct extensive ex-  
141 periments, exploring a range of hyper-parameters and training strategies. Our findings challenge  
142 several widely used practices, including TULU, and we hope that our work can serve as a valuable  
143 reference for practitioners and spark discussions on a deeper understanding of training dynamics for  
144 fine-tuning LLMs.

145  
146  
**Traditional Wisdom in Neural Networks Training.** Identifying effective training configurations  
147 to improve model generalization has been an active area of research long before the rise of LLMs  
148 (Zhang et al., 2017; Srivastava et al., 2014; Ioffe & Szegedy, 2015). For example, Jiang et al.  
149 (2019) conducted extensive experiments to analyze how different generalization measures predict  
150 final model performance, offering insights for hyper-parameter tuning. However, many established  
151 findings do not always extend to LLMs. For example, Keskar et al. (2016) found that large batch  
152 sizes led to poor generalization due to sharp minima. In contrast, our work shows that using large  
153 batches results in higher scores on MT-Bench, indicating improved generalization on downstream  
154 performance. This discrepancy arises due to the difference in experimental settings, where both  
155 architecture (CNNs and MLPs vs. Transformers), and task domains (CV vs. NLP) vary significantly;  
156 and importantly LLMs are pre-trained on massive datasets which drastically changes the starting  
157 point for supervised fine tuning (Peng et al., 2023). Additionally, fine-tuning LLMs poses unique  
158 challenges, as it often requires state-of-the-art clusters spanning multiple machines, each equipped  
159 with multiple GPUs, and advanced networking to optimize speed, memory efficiency, and scalability  
160 using frameworks such as DeepSpeed (Rasley et al., 2020), PyTorch’s FSDP (Zhao et al., 2023)  
161 or Megatron-LM (Narayanan et al., 2021). These requirements are not typically encountered in  
conventional deep learning workflows.

## 2 EXPERIMENTAL SETUP

This section outlines the pre-trained LLMs, the datasets curated for fine-tuning these models, the training strategies used, and the hyper-parameters tested in our experiments. Details on the training infrastructure and optimization techniques used in our experiments can be found in Appendix A.3. We directly present the experiments and results in the following sections. For readers interested in the detailed experimental design and hypotheses, please refer to Appendix A.4.

### 2.1 BASE MODELS AND DATASETS

We conduct experiments using three open-source, small-sized LLMs: Granite 3B<sup>1</sup>, Granite 7B, LLaMA 3.2 3B, and Mistral 7B. The Granite models (Mishra et al., 2024), developed by IBM Research, are based on the transformer architecture (Vaswani et al., 2017) and differ primarily in their model sizes. The LLaMA model (Touvron et al., 2023), also based on the transformer architecture, is widely recognized for its efficient scaling laws and serves as a strong baseline for general-purpose LLMs. The Mistral model (Jiang et al., 2023), developed by Mistral AI, is also a transformer-based architecture with 7 billion parameters, incorporating architectural optimizations for improved performance. We include the Granite, LLaMA, and Mistral models to ensure that our findings on fine-tuning strategies generalize across varying architectures and model sizes within the small-sized LLM category.

We curate a comprehensive set of data designed to progressively enhance the base models’ capabilities in instruction following (phase 00), foundational knowledge (phase 05), and complex skills (phase 10). The datasets are organized into three phases, each focused on specific aspects of language understanding and generation (see Appendix A.2 for details). We also conducted experiments with the TULU dataset (Wang et al., 2023a; Ivison et al., 2023), a diverse mix of complex, instruction-tuning data from human and GPT-4 sources. Finally, we test our findings on a synthetically generated Math, Reasoning, and Code dataset, similar to our other datasets, with a focus on tasks in these domains to ensure they hold for domain-specific datasets.

### 2.2 TRAINING STRATEGIES

We explore two training strategies—*sequential phased training* and *stacked training*. Phased training follows the approach adopted by recent instruction tuning research (Sudalairaj et al., 2024; Mitra et al., 2023; Pang et al., 2024), where the base model is fine-tuned on different data types in a pre-determined sequence. This strategy aims to mitigate catastrophic forgetting and allows the model to build progressively on knowledge and skills acquired in earlier stages. In our experiments, models are fine-tuned in multiple phases, each focusing on a specific type of data (see Appendix A.2 for details on the datasets used in each phase). At the end of each phase, the best-performing checkpoint is selected based on evaluation metrics before proceeding to the next phase. Stacked training combines all data from different phases into a single training phase, exposing the model to diverse data simultaneously. This approach simplifies the training pipeline by eliminating the need for phase-wise data curation.

### 2.3 HYPERPARAMETERS

Our experiments explore various hyperparameter configurations to analyze their impact on training dynamics and model performance.

- **Batch Size.** We investigate effective batch sizes of 128 (small), 3,840 (medium), and 7,680 (large) samples. The effective batch size is achieved through a combination of micro-batch sizes and gradient accumulation steps. For instance, on 64 GPUs, we can process a batch of 3,840 samples in a single micro-batch, whereas on 1 GPU or 8 GPUs, we use gradient accumulation to approximate the same batch size. We confirm that gradient accumulation on a single node produces equivalent results to multi-node distributed training, with details in Appendix A.5.10.
- **Learning Rate and Warmup Steps.** We experiment with various goal learning rates:  $2 \times 10^{-5}$ ,  $3 \times 10^{-5}$ ,  $4 \times 10^{-5}$ ,  $6 \times 10^{-5}$ ,  $8 \times 10^{-5}$ , and  $1 \times 10^{-4}$ . Warmup steps are varied among 0,

<sup>1</sup>We got early access to a preview version of the Granite 3B model.

Table 1: Summary of hyperparameter configurations.

Hyperparameter	TULU	TULU++	LAB
Effective Batch Size	128 samples	Same as TULU	3,840 or 7,680 samples
Learning Rate Scheduler	Warmup ratio: 0.03 Linear decay until the end of training	Same as TULU	Warmup ratio: 0.01 (25 steps linear warmup) No decay (constant rate after warmup)
Number of Epochs	3	4	10
Goal Learning Rate	$2 \times 10^{-5}$	$3 \times 10^{-5}$	$2 \times 10^{-5}$ (also tested with higher rates)

25, and 100, corresponding to different numbers of samples processed before reaching the goal learning rate. The learning rate schedule typically involve a linear warmup to the goal learning rate, followed by either a constant learning rate or the cosine decay schedule.

- **Training Configurations.** We consider three main hyperparameters configurations: LAB (Sudalairaj et al., 2024), TULU (Wang et al., 2023a; Ivison et al., 2023), and a new configuration introduced in this paper, TULU++. Details of these configurations are provided in Table 1.

We used the LAB hyperparameter configuration for all experiments where we varied a single factor (e.g., batch size, learning rate, learning rate schedule, training strategy) while keeping all other settings constant to isolate its effect. For comparisons between TULU and LAB, we directly used the respective configurations. LAB and TULU were chosen as primary configurations due to their prominence: TULU is widely regarded as a gold standard for fine-tuning LLMs with high-quality instruction datasets, while LAB introduces a multi-phase tuning framework leveraging knowledge and skills data to reduce reliance on human annotations.

## 2.4 EVALUATION METRICS

**Benchmarks.** To assess the models’ performance and ability to generalize, we use two primary benchmarks: MMLU (Hendrycks et al., 2020) and MTBench (Zheng et al., 2023). MMLU assesses the models’ knowledge and reasoning across a wide range of subjects. It includes questions from 57 subjects spanning STEM, humanities, social sciences, and more, testing the model’s ability to recall factual knowledge and apply reasoning skills to answer multiple-choice questions. MTBench evaluates multi-turn conversational abilities and generalization to unseen tasks. It measures the quality of responses in dialogue settings, focusing on coherence, relevance, informativeness, and adherence to instructions. The benchmark covers diverse tasks such as reasoning, coding, mathematics, and other skill-based domains. Additionally, we evaluated our models on MMLU-Pro, GPQA, MuSR, MATH, IFEval, and BBH from the Open LLM Leaderboard v2<sup>2</sup>. For the comparison with the TULU dataset, we used the same benchmarks as in the TULU paper (Wang et al., 2023a; Ivison et al., 2023): MMLU, GSM8K, BBH, ToxiGen, and TruthfulQA. Finally, we also include evaluations on ARC (Clark et al., 2018) and GSM8K (Cobbe et al., 2021).

**Efficiency Metrics.** In our experiments, *sample efficiency* and *compute efficiency* effectively represent the same metric. This is because, whether we use multiple GPUs for faster fine-tuning or a single GPU, the total GPU compute and aggregate training hours remain the same. Thus, methods that are more sample-efficient also exhibit better compute efficiency.

## 3 MAIN RESULTS

In this section, we present the empirical findings of our experiments, focusing on the impact of different training strategies, batch sizes, and hyperparameter configurations on the fine-tuning performance of LLMs. We present results using the Granite 7B model and provide additional experiments to other model sizes and architectures (Granite 3B, LLaMA 3B, and Mistral 7B models) in Appendix A.5.8 to validate the robustness and generalizability of our findings. We include baseline scores for the Granite and LLaMA base pretrained models in applicable tables to facilitate easier interpretation of fine-tuned performance. MTBench scores are not provided for baseline models, as these benchmarks evaluate instruction-following and conversational capabilities not present in base models.

<sup>2</sup>[https://huggingface.co/docs/leaderboards/open\\_llm\\_leaderboard/about](https://huggingface.co/docs/leaderboards/open_llm_leaderboard/about)

### 3.1 STACKED TRAINING VS. SEQUENTIAL PHASED TRAINING

We conducted a comprehensive comparison between *stacked training* and *sequential phased training* to evaluate their effectiveness in fine-tuning small sized LLMs. The analysis was performed using the Granite 7B model and evaluated on the MMLU, MTBench, ARC, GSM8K, and Leaderboard (BBH, MATH, MuSR) benchmarks. We observed that stacked training outperformed sequential phased training in both performance and sample efficiency across all batch sizes – 128 and 3,840. The detailed comparison of performance across batch sizes is presented in Appendix A.5.1, along with corresponding figures.

Table 2: Comparison of Stacked vs. Phased Training Strategies. Samples indicate the number of data points required to reach peak performance for each benchmark. Cells highlighted in green indicate better scores, and blue indicates higher sample efficiency (fewer samples used).

Benchmark	Score			Samples	
	Granite Base	Stacked	Phased	Stacked	Phased
MMLU	0.48	0.53	0.52	3,694,080	7,859,902
MTBench	-	6.77	6.76	4,392,960	8,057,918
Leaderboard (BBH)	0.09	0.10	0.10	3,694,080	8,057,918
Leaderboard (MATH Lvl 5)	0.01	0.01	0.00	3,694,080	8,057,918
Leaderboard (MuSR)	0.01	0.08	0.07	3,694,080	8,057,918
ARC	0.78	0.76	0.74	3,694,080	8,057,918
GSM8K	0.11	0.39	0.37	3,694,080	8,057,918

As shown in Table 2, we compare the performance of stacked and phased training strategies using the LAB hyperparameter configuration, which provided the best overall results for both approaches. Stacked training achieves slightly better performance on most benchmarks and comparable performance on the rest, while also being more sample-efficient, requiring significantly fewer samples to reach peak performance. Detailed plots and scores over all checkpoints during training are provided in Appendix A.5.1.

These findings suggest that the stacked training approach improves performance by enabling the model to learn from diverse data simultaneously. Additionally, phased training demands extra time and samples to identify the optimal checkpoint for transitioning between phases. This requires running the model longer to determine peak performance. The increased overhead further diminishes the sample efficiency of phased training compared to the stacked approach.

### 3.2 IMPACT OF BATCH SIZE

We investigated the effect of batch size on model performance by experimenting with effective batch sizes of 128, 3,840, and 7,680 samples. The experiments were conducted using the Granite 7B model and evaluated on the MMLU and MTBench benchmarks. To ensure a fair comparison, we ran each experiment for approximately the same number of gradient steps.

**Observations.** Larger batch sizes lead to better final performance but may require more computational resources and training samples. For stacked training, larger batch sizes uniformly resulted in improved performance on both MMLU and MTBench. The consistent gains from larger batch sizes may stem from the increased data diversity within each batch, which proves especially beneficial when working with datasets that combine multiple task types and knowledge domains. For phased training, the batch size of 3,840 samples outperformed the smaller batch size of 128 samples. While larger batch sizes still yield better overall performance, the impact is less pronounced compared to stacked training. This could be due to each phase concentrating on a specific data type, thereby limiting diversity within batches. Table 3 illustrates the performance of different batch sizes in both stacked and phased training on the MMLU and MTBench benchmarks.

**Trade-off.** We observed that models trained with smaller batch sizes achieved higher performance faster in terms of the number of processed samples but plateaued earlier compared to those trained with larger batch sizes. Conversely, models with larger batch sizes required longer training time to reach similar performance levels due to fewer gradient updates per epoch. This trend is illustrated in

Table 3: Comparison of Batch Sizes Across Stacked and Phased Training Strategies on MMLU and MTBench Benchmarks. Green cells indicate better scores, while blue cells highlight higher sample efficiency (fewer samples required).

Benchmark	Strategy	Score				Samples		
		Granite Base	128	4K	8K	128	4K	8K
MMLU	Stacked	0.48	0.516	0.526	0.529	2,099,328	3,694,080	8,885,760
	Phased	0.48	0.513	0.524	-	2,915,233	7,859,902	-
MTBench	Stacked	-	6.406	6.768	6.831	1,799,424	4,392,960	8,586,240
	Phased	-	6.325	6.756	-	2,815,265	8,057,918	-

Appendix A.5.2, where the performance curves for larger batch sizes span a greater number of training samples. However, with extended training, larger batch sizes led to higher final performance.

### 3.3 EFFECT OF LEARNING RATE SCHEDULES ON LARGE BATCH SIZES

We explored whether using a cosine decay learning rate schedule improves model performance when training with large batch sizes. Cosine decay is often thought to facilitate convergence by allowing higher initial learning rates and gradually reducing them. It can be particularly beneficial when training with large batches that may require larger steps to make meaningful progress. We conducted experiments using the Granite 7B model with an effective batch size of 3,840 samples. We compared two learning rate schedules: a constant learning rate and a cosine decay schedule. The learning rate tested was  $2 \times 10^{-5}$ .

Table 4: Effect of Cosine Decay on MMLU and MTBench Scores at Learning Rate  $2 \times 10^{-5}$ . Cells highlighted in green indicate better scores, and blue indicates higher sample efficiency (fewer samples used).

Benchmark	Score			Samples	
	Granite Base	No Decay	Cosine Decay	No Decay	Cosine Decay
MMLU	0.48	0.5242	0.5251	2,475,200	1,188,096
MTBench	-	6.7562	6.6813	2,673,216	1,188,096

**Observations.** As shown in Table 4, the models trained with a constant learning rate (no decay) performed on par with those trained with a cosine decay schedule on both MMLU and MTBench, and in some cases even outperformed them, particularly on MTBench. Detailed plots are provided in Appendix A.5.5.

**Analysis.** Our findings suggest that, contrary to common practice, cosine decay may not improve model performance when fine-tuning small-size LLMs with large batch sizes. Instead, a constant learning rate ensures consistent progress throughout training, under the assumption that the initial rate is suitable for stable training. For practitioners, this implies that using a constant learning rate could simplify the training process without compromising performance, and may even offer slight improvements.

### 3.4 TULU vs. LAB

We compared the TULU and LAB hyperparameter configurations to assess their effectiveness in enhancing the model’s memorization and generalization capabilities. Memorization was evaluated using a subset of the MMLU benchmark focused on factual knowledge domains, while generalization was assessed using the MTBench benchmark, which tests the model’s ability to perform diverse and complex tasks requiring various skills. Detailed plots and scores over all checkpoints during training are provided in Appendix A.5.4, along with performance results for the Leaderboard (BBH, MATH Lvl 5, MuSR), ARC, and GSM8K benchmarks. Tables 8 and 9 show that LAB outperforms TULU across all benchmarks.

**Cross-Dataset Evaluation with the TULU Dataset.** To further investigate the impact of batch size on fine-tuning performance, we conducted an experiment using the TULU dataset (Wang et al., 2023a; Ivison et al., 2023). This dataset is a refined mixture of instruction-tuning data, integrating both human and GPT-4-generated instructions that are complex and cover various domains. We fine-tuned the Granite 7B model on the TULU dataset using two configurations: batch size of 128 as recommended by TULU, and our configuration with a larger batch size of 3,840 (denoted as 4k). The rationale for testing across datasets was to determine whether the advantages of larger batch sizes observed on our datasets would generalize to different fine-tuning datasets. We evaluated the models using the same benchmarks as in the TULU paper: MMLU (Hendrycks et al., 2020), GSM8K (Cobbe et al., 2021), BBH (Suzgun et al., 2022), ToxiGen (Hartvigsen et al., 2022), and TruthfulQA (Lin et al., 2021). The results in Table 5 demonstrate the broad applicability of larger batch sizes in fine-tuning, with the 4k batch size outperforming the 128 batch size across all but one metric (TruthfulQA).

Table 5: Evaluation results for TULU Dataset across batch sizes. Cells highlighted in green indicate better scores. “Ours” refers to the configuration with a batch size of 4k, while “Theirs” uses the original batch size of 128.

Benchmark	Theirs (128 Batch Size)	Ours (4k Batch Size)
MMLU	0.48	0.50
BBH	0.40	0.44
GSM8K	0.25	0.28
ToxiGen	0.54	0.55
TruthfulQA	0.45	0.44

### 3.5 EFFECT OF LEARNING RATE

We examined how different learning rates impact the model’s downstream performance, using Granite 7B as a base model. We used the LAB hyperparameter configuration since it outperformed TULU. We conducted a learning rate sweep from  $2 \times 10^{-5}$  to  $1 \times 10^{-4}$ . All other hyperparameters were kept constant to isolate the effect of the learning rate. We evaluated on MMLU, MTBench, Leaderboard (BBH, MuSR), ARC, and GSM8K benchmarks after the final phase of phased training. As shown in Table 6, the lowest learning rate of  $2 \times 10^{-5}$  yielded the best performance on most benchmarks and comparable performance on the rest. As the learning rate increased, there was a consistent decline in benchmark performance. This trend suggests that lower learning rates enhance the model’s ability to generalize to unseen tasks requiring knowledge, complex reasoning, and instruction following.

Table 6: Effect of Learning Rate Sweep on Benchmark Scores. Cells highlighted in green indicate better scores.

Benchmark	Granite Base Pretrained	Learning Rates			
		2e-5	4e-5	8e-5	1e-4
MMLU	0.48	0.52	0.52	0.52	0.52
MTBench	-	6.76	6.64	6.53	6.47
Leaderboard (BBH)	0.09	0.10	0.09	0.09	0.08
Leaderboard (MuSR)	0.01	0.08	0.07	0.08	0.06
ARC	0.78	0.74	0.75	0.75	0.73
GSM8K	0.11	0.38	0.36	0.37	0.30

Lower learning rates may aid in retaining knowledge from previous training phases (e.g., instruction following and memorization) by preventing abrupt changes to the model’s parameters. This is particularly important when fine-tuning on complex skills in Phase 10, as it requires the model to build upon its existing capabilities without forgetting prior knowledge. Additionally, our experiments revealed that larger batch sizes did not require higher learning rates. Further details are provided in Appendix A.5.3.



### 3.6 EFFECT OF WARMUP STEPS

We investigated the impact of the number of warmup steps on the training process and final model performance. The warmup phase is traditionally considered crucial for stabilizing training, especially when using higher learning rates, by gradually increasing the learning rate from a small value to its target value over a specified number of steps (Goyal et al., 2017). We ran experiments with the Granite 7B model in the stacked setting using LAB hyperparameters—our best configuration—across three warmup setups: 0, 25, and 100 warmup steps, corresponding to approximately 0, 96,000, and 384,000 samples processed before reaching the target learning rate, respectively. As shown in Appendix A.5.6, the model trained without warmup steps achieved better performance on the MMLU benchmark and similar performance on MTBench compared to models trained with 25 or 100 warmup steps. The training curves for all configurations followed a similar trajectory, converging to comparable performance levels within approximately the same number of training steps, indicating that omitting warmup steps does not negatively affect the final model performance. Although omitting the warmup simplifies the training process, it offers no advantage in terms of faster convergence. Furthermore, we monitored training dynamics such as gradient norms and loss values across the different warmup configurations. As shown in Appendix A.5.6, the gradient norms and loss values exhibited similar patterns across all configurations, indicating stable training even without a warmup phase.

### 3.7 EARLY TRAINING DYNAMICS AS PREDICTORS OF FINAL PERFORMANCE

We consistently observed that models exhibiting lower gradient norms and higher training loss values during training achieved better final performance on MMLU and MTBench. Figures 1 and 2 illustrate the correlation between early training dynamics—gradient norms and loss values—and final benchmark performances.

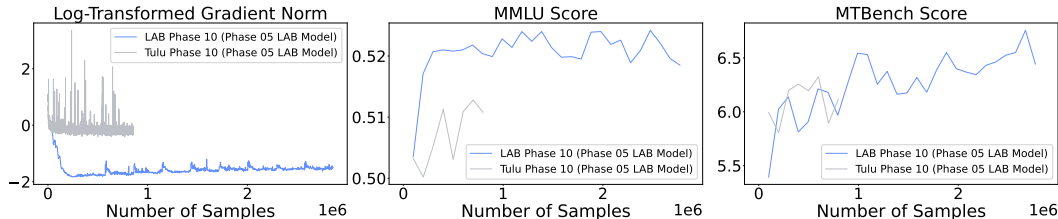


Figure 1: Correlation between early training dynamics and final performance on MMLU and MTBench benchmarks for TULU vs. LAB Phase 10 training.

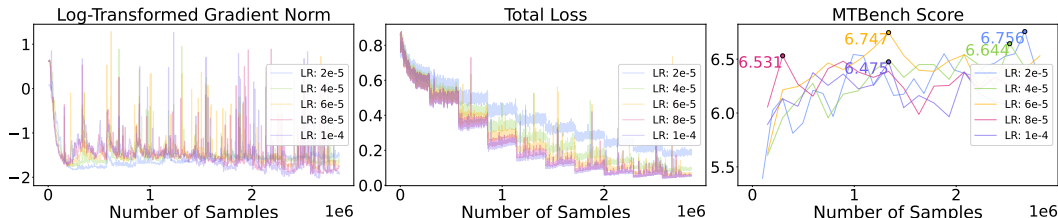


Figure 2: LAB Learning Rate (LR) Sweep: Training Dynamics and MTBench Performance. MMLU results are provided in Appendix A.5.9.

**TULU vs. LAB Phase 10 Training (Figure 1).** The LAB configuration achieved better final performance with lower gradient norms compared to the TULU configuration.

**LAB Learning Rate Sweep Experiments.** Models trained with a learning rate of  $2 \times 10^{-5}$  demonstrated lower gradient norms initially, which increased toward the end of training, and higher loss throughout, ultimately resulting in superior final performance compared to models trained with higher learning rates. For the most effective learning rates, the gradient norm started at its lowest value and increased towards the end of training (Figure 2). Despite the higher gradient norms in the later stages, the associated loss remained higher throughout the entire training for these rates.

This is consistent with the use of lower learning rates, which typically result in higher training loss but better generalization. Figure 2 shows the gradient norms and loss values for different learning rates, along with the final performance on MTBench. The lowest learning rates delivered superior results. Smaller learning rates may enable the model to stabilize the learning process initially and then gradually explore more challenging regions of the loss landscape as training progresses, leading to better generalization and final performance. We hypothesize that lower gradient norm values at the start of training contribute to a smoother and more stable optimization process, preventing the model from overfitting too quickly. This allows for gradual learning, which we hypothesize facilitates better exploration of the loss landscape as training progresses. The subsequent increase in gradient norm during later training stages may indicate that the model is delving into more complex regions of the parameter space, enhancing its ability to generalize. MMLU results are provided in Appendix A.5.9.

The correlation between early training dynamics and final performance holds across different batch sizes, warmup steps, and learning rate schedules (see Appendix A.5.9 for additional results).

## 4 DISCUSSION, GUIDELINES FOR PRACTITIONERS, AND LIMITATIONS

**Balancing Performance and Efficiency.** Our results show a trade-off between performance and computational cost. Configurations such as higher batch sizes or lower learning rates achieve better final performance but take longer to converge. In contrast, hyperparameters yielding lower final performance often dominate early on before plateauing. For those with limited resources, smaller batch sizes or higher learning rates may be more efficient. For example, in stacked training, a 4k batch size outperforms 8k initially, and higher learning rates offer faster learning in early stages. Moreover, we encourage practitioners to monitor early training dynamics, such as gradient norms and loss values, as they correlate strongly with final model performance. Observing lower gradient norms and higher loss values during the initial phases of training can serve as reliable indicators of better generalization capabilities. This allows for early termination of suboptimal runs, conserving computational resources.

**Training Strategy Recommendations.** Based on our empirical evidence, we advocate for stacked over sequential phased training. This recommendation is supported by consistent performance gains and improved sample efficiency observed in both the Granite 7B and Granite 3B models. Stacked training simplifies the fine-tuning process and eliminates the need for phase-wise data management.

**Hyperparameter Selection.** We offer guidance on selecting batch sizes, learning rates, warmup steps, and learning rate schedules. Larger batch sizes (e.g., 4k and 8k) are recommended, as they have demonstrated superior performance across model sizes compared to smaller batch sizes like 128. Low learning rates are crucial for optimal performance. We found that  $2 \times 10^{-5}$  works well for Granite models, while  $1 \times 10^{-6}$  is optimal for Mistral. Lower learning rates allow for more precise adjustments to the model weights, preventing overshooting in the optimization landscape. Practitioners should start with these values and, if necessary, perform a localized search by testing slightly higher or lower learning rates to find the optimal setting for their specific model. This approach significantly reduces the search space. Our experiments indicate that omitting warmup steps and using a constant learning rate instead of cosine decay does not negatively impact performance, simplifying the training process without sacrificing model quality.

**Limitations.** Our experiments focused on small (3B to 7B parameters) LLMs and were conducted on two model architectures: Granite (based on the Llama architecture) and Mistral. While our findings are promising, they may not directly generalize to larger models or other architectures. Future work should explore larger models and a broader range of architectures, such as Gemma and others. Also, we did not explore parameter-efficient fine-tuning techniques like LoRA, alternative pre-training objectives, or different tokenizer configurations, which might affect the applicability of our fine-tuning strategies. Furthermore, our evaluation was centered on synthetic datasets generated based on a comprehensive taxonomy encompassing various knowledge and skills, using benchmarks like MMLU and MTBench that align with this dataset. It would be valuable to assess whether our findings apply to other benchmarks like GSM8K, BBH, IFEval, and ARC. Confirming this across diverse evaluation metrics would strengthen the generalizability of our conclusions. Finally, we acknowledge that our experiments were conducted using a single seed due to computational constraints, which may introduce some noise into the observations.

## REFERENCES

- 540  
541  
542 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and  
543 extraction. *arXiv preprint arXiv:2309.14316*, 2023.
- 544 Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding,  
545 Kai Dong, Qiusi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with  
546 longtermism. *arXiv preprint arXiv:2401.02954*, 2024.
- 547 Daixuan Cheng, Shaohan Huang, and Furu Wei. Adapting large language models via reading com-  
548 prehension. In *The Twelfth International Conference on Learning Representations*, 2023.
- 549 Daixuan Cheng, Yuxian Gu, Shaohan Huang, Junyu Bi, Minlie Huang, and Furu Wei. In-  
550 struction pre-training: Language models are supervised multitask learners. *arXiv preprint*  
551 *arXiv:2406.14491*, 2024.
- 552 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,  
553 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned lan-  
554 guage models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- 555 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
556 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
557 *arXiv preprint arXiv:1803.05457*, 2018.
- 558 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
559 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
560 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 561 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-  
562 efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*,  
563 35:16344–16359, 2022.
- 564 Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong  
565 Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional  
566 conversations. In *Conference on Empirical Methods in Natural Language Processing*, 2023.  
567 URL <https://api.semanticscholar.org/CorpusID:258840897>.
- 568 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
569 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
570 *arXiv preprint arXiv:2407.21783*, 2024.
- 571 GDPR. General data protection regulation (GDPR). <https://gdpr-info.eu>, 2016.
- 572 Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, An-  
573 drew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet  
574 in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- 575 R. L. Graham. Bounds for certain multiprocessing anomalies. *Bell System Technical Journal*, 45(9):  
576 1563–1581, 1966. doi: 10.1002/j.1538-7305.1966.tb01709.x.
- 577 Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth  
578 Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. Textbooks are  
579 all you need. *arXiv preprint arXiv:2306.11644*, 2023.
- 580 Tom Gunter, Zirui Wang, Chong Wang, Ruoming Pang, Andy Narayanan, Aonan Zhang, Bowen  
581 Zhang, Chen Chen, Chung-Cheng Chiu, David Qiu, et al. Apple intelligence foundation language  
582 models. *arXiv preprint arXiv:2407.21075*, 2024.
- 583 Alexander Hägele, Elie Bakouch, Atli Kosson, Loubna Ben Allal, Leandro Von Werra, and Martin  
584 Jaggi. Scaling laws and compute-optimal training beyond fixed training durations. *arXiv preprint*  
585 *arXiv:2405.18392*, 2024.
- 586 Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar.  
587 Toxigen: A large-scale machine-generated dataset for adversarial and implicit hate speech detec-  
588 tion. *arXiv preprint arXiv:2203.09509*, 2022.
- 589  
590  
591  
592  
593

- 594 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
595 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*  
596 *arXiv:2009.03300*, 2020.
- 597  
598 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
599 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training  
600 compute-optimal large language models. In *Proceedings of the 36th International Conference*  
601 *on Neural Information Processing Systems*, pp. 30016–30030, 2022.
- 602 Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. Unnatural instructions: Tuning  
603 language models with (almost) no human labor. In *The 61st Annual Meeting Of The Association*  
604 *For Computational Linguistics*, 2023.
- 605  
606 Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang,  
607 Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zhen Leng Thai, Kaihuo Zhang, Chongyi Wang,  
608 Yuan Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chaochao Jia, Guoyang  
609 Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. Minicpm: Unveiling the potential of small  
610 language models with scalable training strategies. *ArXiv*, abs/2404.06395, 2024. URL <https://api.semanticscholar.org/CorpusID:269009975>.
- 611  
612 Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L Richter, Quentin Anthony, Timothée  
613 Lesort, Eugene Belilovsky, and Irina Rish. Simple and scalable strategies to continually pre-train  
614 large language models. *arXiv preprint arXiv:2403.08763*, 2024.
- 615  
616 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by  
617 reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456.  
618 PMLR, 2015.
- 619  
620 Berivan Isik, Natalia Ponomareva, Hussein Hazimeh, Dimitris Pappas, Sergei Vassilvitskii, and  
621 Sanmi Koyejo. Scaling laws for downstream task performance of large language models. In *ICLR*  
622 *2024 Workshop on Navigating and Addressing Data Problems for Foundation Models*, 2024. URL  
623 <https://openreview.net/forum?id=PXwdsgZjnX>.
- 624  
625 Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep  
626 Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. Camels in a changing cli-  
627 mate: Enhancing lm adaptation with tulu 2. *arXiv preprint arXiv:2311.10702*, 2023.
- 628  
629 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
630 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
631 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 632  
633 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bam-  
634 ford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.  
635 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- 636  
637 Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic  
638 generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019.
- 639  
640 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,  
641 Scott Gray, Alec Radford, Jeff Wu, and Dario Amodei. Scaling laws for neural language  
642 models. *ArXiv*, abs/2001.08361, 2020. URL <https://api.semanticscholar.org/CorpusID:210861095>.
- 643  
644 Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Pe-  
645 ter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv*  
646 *preprint arXiv:1609.04836*, 2016.
- 647  
648 Lisa Lee. Tiny titans: How small language models outperform llms for less, 2024. URL <https://www.salesforce.com/blog/small-language-models/>.
- 649  
650 Timothée Lesort, Massimo Caccia, and Irina Rish. Understanding continual learning settings  
651 with data distribution drift analysis. *ArXiv*, abs/2104.01678, 2021. URL <https://api.semanticscholar.org/CorpusID:233024916>.

- 648 Haoran Li, Qingxiu Dong, Zhengyang Tang, Chaojun Wang, Xingxing Zhang, Haoyang Huang,  
649 Shaohan Huang, Xiaolong Huang, Zeqiang Huang, Dongdong Zhang, et al. Synthetic data  
650 (almost) from scratch: Generalized instruction tuning for language models. *arXiv preprint*  
651 *arXiv:2402.13064*, 2024.
- 652 Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang,  
653 Tianyi Zhou, and Jing Xiao. From quantity to quality: Boosting llm performance with self-  
654 guided data selection for instruction tuning. *ArXiv*, abs/2308.12032, 2023. URL <https://api.semanticscholar.org/CorpusID:261076515>.
- 655  
656 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human  
657 falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- 659 Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data  
660 for alignment? a comprehensive study of automatic data selection in instruction tun-  
661 ing. *ArXiv*, abs/2312.15685, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:266551413)  
662 [CorpusID:266551413](https://api.semanticscholar.org/CorpusID:266551413).
- 663  
664 Mayank Mishra, Matt Stallone, Gaoyuan Zhang, Yikang Shen, Aditya Prasad, Adriana Meza So-  
665 ria, Michele Merler, Parameswaran Selvam, Saptha Surendran, Shivdeep Singh, et al. Gran-  
666 ite code models: A family of open foundation models for code intelligence. *arXiv preprint*  
667 *arXiv:2405.04324*, 2024.
- 668 Arindam Mitra, Luciano Del Corro, Shweti Mahajan, Andres Codas, Clarisse Simoes, Sahaj Agar-  
669 wal, Xuxi Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Aggarwal, Hamid Palangi, Guoqing  
670 Zheng, Corby Rosset, Hamed Khanpour, and Ahmed Awadallah. Orca 2: Teaching small lan-  
671 guage models how to reason, 2023. URL <https://arxiv.org/abs/2311.11045>.
- 672  
673 Arindam Mitra, Luciano Del Corro, Guoqing Zheng, Shweti Mahajan, Dany Rouhana, Andres Co-  
674 das, Yadong Lu, Wei-ge Chen, Olga Vrousos, Corby Rosset, et al. Agentinstruct: Toward gen-  
675 erative teaching with agentic flows. *arXiv preprint arXiv:2407.03502*, 2024.
- 676  
677 Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and  
678 Ahmed Hassan Awadallah. Orca: Progressive learning from complex explanation traces of  
679 gpt-4. *ArXiv*, abs/2306.02707, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:259075316)  
[CorpusID:259075316](https://api.semanticscholar.org/CorpusID:259075316).
- 680  
681 MultipackSampler. Multipack sampler: Multipack distributed sampler for fast padding-free training  
682 of llms, 2024. URL [https://github.com/imoneoi/multipack\\_sampler](https://github.com/imoneoi/multipack_sampler).
- 683  
684 Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vi-  
685 jay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, Amar  
686 Phanishayee, and Matei Zaharia. Efficient large-scale language model training on gpu clus-  
687 ters using megatron-lm. In *Proceedings of the International Conference for High Performance*  
688 *Computing, Networking, Storage and Analysis, SC '21*, New York, NY, USA, 2021. Associa-  
689 tion for Computing Machinery. ISBN 9781450384421. doi: 10.1145/3458817.3476209. URL  
690 <https://doi.org/10.1145/3458817.3476209>.
- 691  
692 OpenAI. Chatgpt: Optimizing language models for dialogue. *OpenAI Blog*, 2022. URL <https://openai.com/blog/chatgpt/>.
- 693  
694 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 695  
696 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong  
697 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-  
698 ton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano,  
699 Jan Leike, and Ryan J. Lowe. Training language models to follow instructions with human  
700 feedback. *ArXiv*, abs/2203.02155, 2022. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:246426909)  
701 [CorpusID:246426909](https://api.semanticscholar.org/CorpusID:246426909).
- 702  
703 Wei Pang, Chuan Zhou, Xiao-Hua Zhou, and Xiaojie Wang. Phased instruction fine-tuning for  
704 large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of*  
705 *the Association for Computational Linguistics ACL 2024*, pp. 5735–5748, Bangkok, Thailand and

- 702 virtual meeting, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.  
703 findings-acl.341. URL <https://aclanthology.org/2024.findings-acl.341>.  
704
- 705 Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning  
706 with gpt-4. *ArXiv*, abs/2304.03277, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:257985497)  
707 [CorpusID:257985497](https://api.semanticscholar.org/CorpusID:257985497).
- 708 Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System opti-  
709 mizations enable training deep learning models with over 100 billion parameters. *Proceedings*  
710 *of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*,  
711 2020. URL <https://api.semanticscholar.org/CorpusID:221191193>.  
712
- 713 RedHat. Llms vs slms, 2024. URL [https://www.redhat.com/en/topics/ai/](https://www.redhat.com/en/topics/ai/llm-vs-slm)  
714 [llm-vs-slm](https://www.redhat.com/en/topics/ai/llm-vs-slm).
- 715 Jonathan S. Rosenfeld, Amir Rosenfeld, Yonatan Belinkov, and Nir Shavit. A constructive prediction  
716 of the generalization error across scales. *ArXiv*, abs/1909.12673, 2019. URL [https://api.](https://api.semanticscholar.org/CorpusID:203592013)  
717 [semanticscholar.org/CorpusID:203592013](https://api.semanticscholar.org/CorpusID:203592013).  
718
- 719 Nikhil Sardana, Jacob Portes, Sasha Doubov, and Jonathan Frankle. Beyond chinchilla-optimal: Ac-  
720 counting for inference in language model scaling laws. In *Forty-first International Conference on*  
721 *Machine Learning*, 2024. URL <https://openreview.net/forum?id=0bmXrtDUu>.  
722
- 723 Thomas Scialom, Tuhin Chakrabarty, and Smaranda Muresan. Fine-tuned language models are  
724 continual learners. In *Conference on Empirical Methods in Natural Language Processing, 2022*.  
725 URL <https://api.semanticscholar.org/CorpusID:252815378>.
- 726 Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.  
727 Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine*  
728 *learning research*, 15(1):1929–1958, 2014.
- 729 Shivchander Sudalairaj, Abhishek Bhandwadar, Aldo Pareja, Kai Xu, David D Cox, and Akash  
730 Srivastava. Lab: Large-scale alignment for chatbots. *arXiv preprint arXiv:2403.01081*, 2024.  
731
- 732 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,  
733 Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks  
734 and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.
- 735 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
736 Liang, and Tatsunori B Hashimoto. Stanford alpaca: An instruction-following llama model, 2023.  
737
- 738 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
739 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
740 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 741 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
742 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Infor-*  
743 *mation Processing Systems*, 2017.  
744
- 745 Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David  
746 Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. How far can camels go? exploring  
747 the state of instruction tuning on open resources. *Advances in Neural Information Processing*  
748 *Systems*, 36:74764–74786, 2023a.
- 749 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and  
750 Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. In  
751 *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume*  
752 *1: Long Papers)*, pp. 13484–13508, 2023b.  
753
- 754 Yong Xie, Karan Aggarwal, and Aitzaz Ahmad. Efficient continual pre-training for building domain  
755 specific large language models. In *Annual Meeting of the Association for Computational Linguis-*  
*tics*, 2023. URL <https://api.semanticscholar.org/CorpusID:265213147>.

756 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and  
757 Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions.  
758 *arXiv preprint arXiv:2304.12244*, 2023.

759  
760 Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ry-  
761 der, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural  
762 networks via zero-shot hyperparameter transfer. *arXiv preprint arXiv:2203.03466*, 2022.

763 Zitong Yang, Neil Band, Shuangping Li, Emmanuel Candès, and Tatsunori Hashimoto. Synthetic  
764 continued pretraining. *arXiv preprint arXiv:2409.07431*, 2024.

765  
766 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding  
767 deep learning requires rethinking generalization. In *International Conference on Learning Rep-  
768 resentations*, 2017. URL <https://openreview.net/forum?id=Sy8gdB9xx>.

769 Vivienne Zhang, Shashank Verma, Neal Vaidya, Abhishek Sawarkar, and Amanda Saunders. Nvidia  
770 ai foundation models: Build custom enterprise chatbots and co-pilots with production-ready llms,  
771 2023.

772 Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright,  
773 Hamid Shojanazeri, Myle Ott, Sam Shleifer, et al. Pytorch fsdp: Experiences on scaling fully  
774 sharded data parallel. *Proceedings of the VLDB Endowment*, 16(12):3848–3860, 2023.

775  
776 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
777 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
778 chatbot arena. *Advances in Neural Information Processing Systems*, 2023.

## 780 A APPENDIX

### 781 A.1 MODEL DETAILS

782  
783  
784 **Granite 3B** is composed of transformer layers (encoder blocks) that include multi-head self-  
785 attention mechanisms and feed-forward networks. It has a smaller hidden size and fewer attention  
786 heads, making it less computationally intensive and faster for both training and inference.

787  
788 **Granite 7B** has more transformer layers and increased hidden dimensions, offers greater rep-  
789 resentational capacity. It also includes more attention heads, enabling it to capture more complex  
790 language patterns and long-range dependencies.

791  
792 **Llama 3.2 3B** employs a scaled-down transformer architecture with fewer layers and a reduced  
793 hidden size compared to larger models in the Llama family. It maintains the core design principles of  
794 its larger counterparts, including rotary positional embeddings and optimized attention mechanisms,  
while balancing performance and efficiency for resource-constrained environments.

795  
796 **Mistral 7B** incorporates advanced attention mechanisms, including multi-query attention, to re-  
797 duce memory usage and increase computational efficiency during inference. Additionally, it lever-  
798 ages training optimizations like mixed-precision training and gradient checkpointing to accelerate  
training and lower resource demands.

### 800 A.2 DATASETS DETAILS

801  
802 The datasets were curated using a taxonomy-driven approach to ensure comprehensive coverage of  
803 instruction-following, foundational knowledge, and compositional skills. The taxonomy hierarchi-  
804 cally organizes tasks into three main branches—knowledge, foundational skills, and compositional  
805 skills—each further divided into granular subcategories. For each subcategory, manually written  
806 instruction-response pairs served as seed examples. These examples guided synthetic data genera-  
807 tion using teacher models (e.g., Mixtral-7x8B) to expand the dataset while maintaining high quality  
808 and diversity. For knowledge data, reliable sources such as textbooks and technical manuals pro-  
809 vided a grounding for synthetic questions and responses. Foundational skills data were drawn from  
public datasets covering essential areas like mathematics, coding, and reasoning. Compositional  
skills were synthesized using a taxonomy-guided approach to combine knowledge and foundational

skills for complex tasks, such as writing detailed emails or generating logical arguments. We provide details about the datasets we used in Table 7.

Table 7: Summary of datasets used in different phases.

Phase	Description	# Samples
<b>Phase 00</b>	Instruction following warmup: simple, template-based instruction-response pairs to transition the base models to instruction-following behavior.	308343
<b>Phase 05</b>	Foundational knowledge acquisition: synthetically generated question-answer pairs from textbooks covering a wide range of disciplines up to graduate-level courses.	231178
<b>Phase 10</b>	Complex skills development: synthetic data generated using a taxonomy of skills, including tasks like poetry, email writing, logical reasoning, coding, and more.	285966
<b>All-Phases</b>	Combination of phases 00, 05, and 10, exposing models to all data types simultaneously.	825487

### A.3 TRAINING INFRASTRUCTURE AND OPTIMIZATION

To handle large batch sizes and optimize computational efficiency, we use an optimized training infrastructure.

**Optimizer.** Across all experiments, we use the Adam optimizer with  $\beta_1 = 0.9$  and  $\beta_2 = 0.95$ . By adjusting  $\beta_2 = 0.95$ , we reduce the emphasis on the variance of past gradients, which is beneficial when training with large batch sizes that provide more stable gradient estimates.

**Batching and Gradient Accumulation.** To achieve the large effective batch sizes required for our experiments, we employed gradient accumulation techniques. Gradient Accumulation involves accumulating gradients over multiple forward and backward passes before performing an optimizer step (i.e., updating model weights). This effectively increases the batch size without necessitating additional memory to store larger batches in a single pass. For instance, in a single-node setup with 8 GPUs, we set a micro-batch size per GPU and used gradient accumulation steps to reach an effective batch size of 3,840 samples. Specifically, if each GPU processes a micro-batch of  $b$  samples and we accumulate gradients over  $k$  steps, the effective batch size  $B$  is  $B = b \times k \times N$ , where  $N$  is the number of GPUs. In multi-node setups with 64 GPUs, we could process the entire batch in a single step without accumulation due to the distributed computational resources. This approach allowed us to simulate very large batch sizes, to investigate their impact on model performance.

**Efficient Distributed Sampling.** We implement a variant of Multipack distributed sampler (MultipackSampler, 2024), which offers significant advantages over naive sampling approaches in distributed training of LLMs. Drawing on concepts from the identical machine scheduling problem (Graham, 1966), our implementation uses an approximate solution at the sample level, achieving near-optimal GPU utilization. Our variant extends the original design by accounting for padding, crucial for non-linear attention mechanisms like scaled dot-product attention (Vaswani et al., 2017), and clustering together samples of similar length. It ensures that even with padding, no GPU exceeds a pre-determined token capacity, which we calculate to maintain an expected micro-batch size that satisfies:

$$E[\text{Effective Batch Size}] = E[\text{Micro Batch Size} \times \text{Gradient Accumulation Steps}]$$

This approach balances computational load across GPUs, resulting in improved training throughput and stability. Additionally, our sampler supports both linear attention mechanisms, such as FlashAttention (Dao et al., 2022), and traditional non-linear attention, making it versatile for various model architectures.



#### 864 A.4 EXPERIMENTAL DESIGN

865 We investigate how training strategies, batch sizes, learning rates, and warmup steps influence LLM  
866 fine-tuning. We systematically vary these factors while holding other parameters constant to isolate  
867 their individual effects.

868  
869 **Impact of Batch Size and Training Strategies.** We examine how different batch sizes influence  
870 model performance and training dynamics in both stacked and phased training settings. Our hy-  
871 pothesis is that larger batch sizes will improve model performance in stacked training by ensuring  
872 sufficient data diversity within each batch, allowing for more robust gradient updates. In contrast,  
873 their impact on phased training may be less pronounced. On the other hand, this improvement may  
874 come at the cost of reduced sample efficiency. While larger batch sizes may achieve better final  
875 performance, they typically require more training samples and computational resources due to the  
876 higher number of samples used per gradient step. Conversely, smaller batch sizes could achieve  
877 comparable performance with fewer samples, especially in phased training, where each batch is  
878 inherently constrained to a specific data type, limiting intra-batch diversity. We formalize these hy-  
879 potheses below and explore the trade-offs between performance gains and computational efficiency  
880 in our experiments.

- 881 • **Hypothesis 1.** Stacked training underperforms at smaller batch sizes due to insufficient diversity  
882 within each batch. A smaller batch may not capture the wide range of data types present in the  
883 combined dataset, leading to less effective learning. In contrast, larger batch sizes in stacked  
884 training could match or surpass phased training by capturing a wider range of signals in each  
885 gradient update.
- 886 • **Hypothesis 2.** While the stacked approach simplifies the training pipeline by eliminating the need  
887 for phase selection of data, it can be less sample efficient. Learning all types of data simultane-  
888 ously could require more steps for the model to adequately learn the complex and diverse patterns  
889 in the combined dataset. This translates to worse sample efficiency, as the model may need more  
890 gradient updates to converge.

891  
892 **Learning Rate Exploration.** We conduct a learning rate sweep to examine its influence on train-  
893 ing dynamics and final model performance. We explore whether larger batch sizes require higher  
894 learning rates, hypothesizing that increased gradient stability at higher batch sizes may allow for  
895 more aggressive learning rate schedules without causing instability. Additionally, we hypothesize  
896 that lower learning rates (e.g.,  $2e-5$ ) might lead to better generalization, as smaller updates allow  
897 the model to converge more gradually. Additionally, we evaluate the effects of different learning rate  
898 schedules (constant vs. cosine decay) on model performance, as cosine decay is widely recognized  
899 for its smooth convergence properties. Our goal is to empirically determine whether this schedule  
900 offers tangible benefits in our specific setting.

901  
902 **Warmup Steps Analysis.** Warmup steps are commonly used in fine-tuning LLMs to stabilize train-  
903 ing. We investigate whether reducing or removing warmup steps can accelerate convergence without  
904 sacrificing final model performance.

905 **Training Dynamics and Early Performance Indicators.** We monitor key training dynamics, such  
906 as gradient norms and loss values, to explore their correlation with the model’s final performance  
907 metrics on benchmarks (MMLU and MTBench). Monitoring gradient norm and loss during training  
908 provides insights into the smoothness of the optimization trajectory, with lower gradient norms sug-  
909 gesting traversal through flatter regions of the loss landscape, which, as discussed later, can influence  
910 final model performance. The goal was to investigate whether early-stage indicators, such as a lower  
911 gradient norm and higher loss values during the initial phase of training or consequently throughout  
912 the entire training, can serve as reliable predictors of better performance on benchmarks. This ap-  
913 proach could allow for the early termination of suboptimal runs, optimizing computational resources  
914 by focusing only on models that demonstrate promising training dynamics. By closely examining  
915 these metrics across multiple learning rate configurations, batch sizes, and training strategies, we  
916 aimed to understand how these dynamics reflect the underlying optimization process and its rela-  
917 tionship to final task performance, without the need for full training to completion. Identifying  
these early indicators is critical for advancing sample-efficient training methodologies, especially in  
large-scale experiments.

**Adaptations for Different Model Architectures.** To evaluate the generalizability of our findings across different model architectures, we extended our experiments to include the Mistral 7B model. Mistral models have architectural optimizations and differ from the Granite models in certain aspects, such as tokenization and layer configurations. Given these differences, we made specific adaptations to the training hyperparameters (learning rate, warmup steps) to verify the robustness of our methodology and check if our experiments are applicable to a range of small sized LLM architectures.

## A.5 ADDITIONAL RESULTS

### A.5.1 STACKED TRAINING VS. SEQUENTIAL PHASED TRAINING

Contrary to one of our initial hypotheses that stacked training might underperform at smaller batch sizes due to insufficient data diversity within each batch, our results indicate that stacked training consistently achieves better or comparable performance to phased training at batch sizes of 128, and 4,000 samples.

Figure 3 illustrates the performance of stacked and phased training strategies on MTBench and MMLU benchmarks across different batch sizes. Notably, stacked training slightly outperforms phased training at each batch size, suggesting that batch size does not significantly impact the difference between the two training strategies. A possible explanation is that, even with smaller batch sizes, the stacked approach provides enough data diversity for effective learning, allowing the model to generalize well across different types of data.

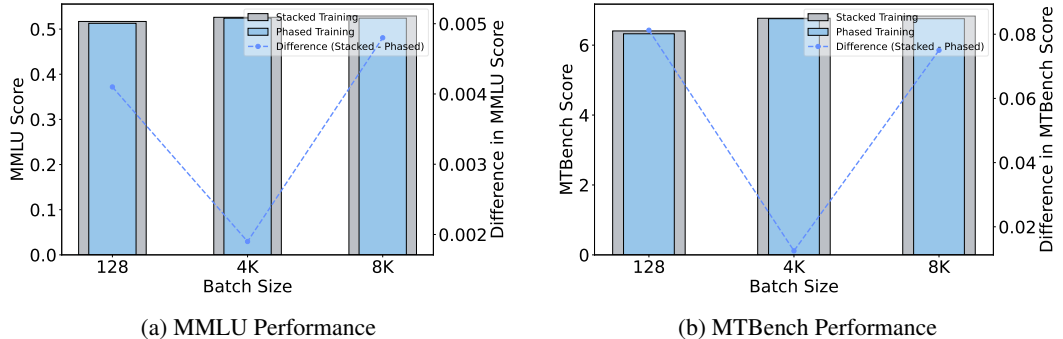


Figure 3: Comparison of stacked and phased training across different batch sizes on MMLU and MTBench benchmarks.

In Figure 4, we compare the performance of both training strategies using the LAB hyperparameter configuration. Figure 4a shows the final MTBench performance, where stacked training outperformed phased training by 0.01 points. Figure 4b illustrates that stacked training is also more sample-efficient, with the best performance points annotated by the number of samples required to reach them. Note that the line for phased training begins partway through, as samples from Phases 00 and 05 were already included. This applies consistently to all similar figures presented in this paper.

In addition to the MTBench results, we include the MMLU performance comparisons here. Figure 5 shows the final MMLU performance using LAB hyperparameters for both stacked and phased training strategies. Stacked training outperformed phased training on the MMLU benchmark by 0.01 points, consistent with the observations from MTBench.

Figure 6 illustrates the sample efficiency comparison for MMLU. Similar to the MTBench results, stacked training achieves higher MMLU performance more quickly than phased training. These results reinforce our findings that stacked training not only improves performance but also enhances sample efficiency on both MMLU and MTBench benchmarks.

To investigate whether phased training might be effective when phases are split based on difficulty, we conducted an additional experiment. In this setup, we partitioned the dataset into two phases based on difficulty, using the length of free-form answers as a proxy for difficulty.

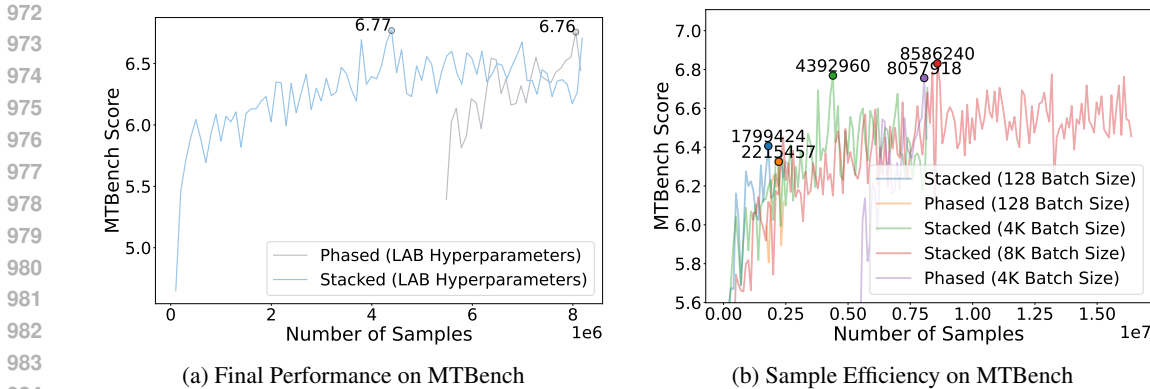


Figure 4: Comparison of stacked and phased training strategies on MTBench using LAB hyperparameters.

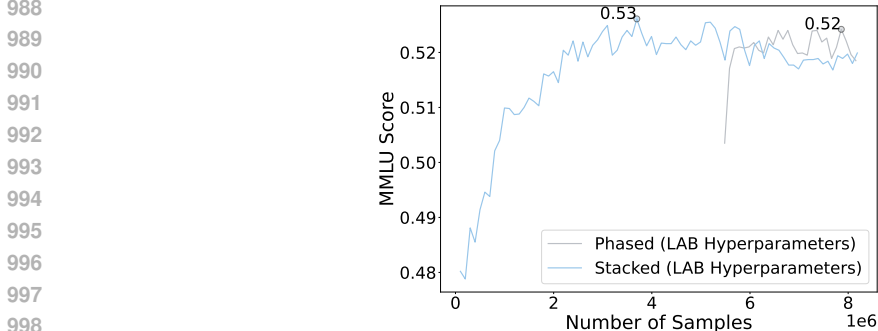


Figure 5: Final MMLU Performance comparison using LAB hyperparameters: stacked vs. phased training.

- **Phase I:** The bottom 50% of the data containing short sentences.
- **Phase II:** The top 50% of the data containing long sentences, plus a 1% subset of the short sentences as a replay buffer when transitioning to long sentences.

We fine-tuned the Granite 7B base model using the same hyperparameters—a batch size of 4,000 and a learning rate of  $2 \times 10^{-5}$ —under both the phased and stacked training strategies. Our results (Figure 7) showed no significant difference between phased and stacked training in this setting. Both performed similarly, with stacked training slightly outperforming phased training across all benchmarks. This suggests that even when the data is carefully partitioned based on difficulty, phased training does not improve model performance over stacked training. Moreover, phased training remains less sample-efficient due to the additional time and samples required to determine the optimal checkpoint for phase transitions.

### A.5.2 IMPACT OF BATCH SIZE

Figure 9 illustrates the performance of different batch sizes in both stacked and phased training on the MTBench benchmark. In addition to the MTBench results, we include the MMLU performance comparisons here. Figure 8 illustrates the impact of batch size on model performance in both stacked and phased training on the MMLU benchmark. The observations are consistent with those from MTBench: larger batch sizes lead to better final performance. The trends observed in the MMLU results reinforce our findings that larger batch sizes improve final performance, although they may require more training samples and computational resources. As with MTBench, smaller batch sizes initially reach higher performance levels more quickly but plateau earlier, while larger batch sizes surpass them with extended training.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

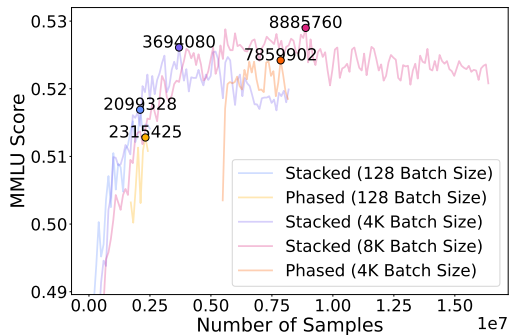


Figure 6: MMLU Sample efficiency comparison between stacked and phased training.

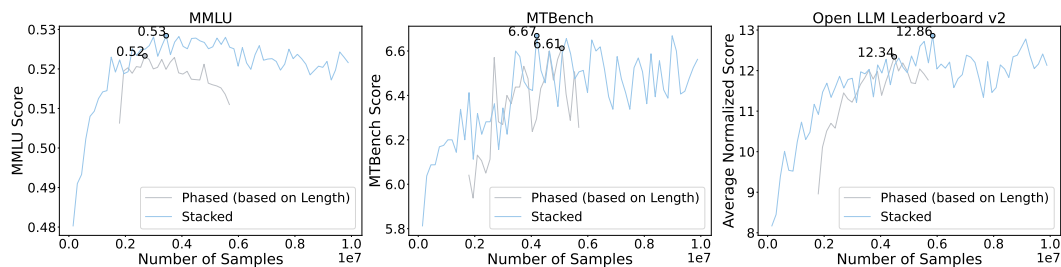


Figure 7: Performance comparison of stacked vs. phased training on difficulty-partitioned data (by answer length) across comprehensive benchmarks.

### A.5.3 EFFECT OF LEARNING RATE

As shown in Figure 10, the lowest learning rate of  $2 \times 10^{-5}$  yielded the best performance on the MTBench Benchmark.

We investigated whether larger batch sizes necessitate higher learning rates, based on the premise that with a larger batch size, the model processes more samples before each gradient step, potentially benefiting from a higher learning rate to make more significant updates and to maintain the variance of the gradient when compared to smaller batch sizes. Additionally, since larger batches result in fewer gradient steps over the same number of epochs, increasing the learning rate might improve training efficiency.

Our experiments compared models trained with different learning rates across batch sizes of 128, 3,840, and 7,680 samples. The runs included TULU hyperparameters at learning rates of  $2 \times 10^{-5}$  and  $3 \times 10^{-5}$ , and LAB hyperparameters with learning rates ranging from  $2 \times 10^{-5}$  to  $1 \times 10^{-4}$ . As shown in Figure 11, regardless of batch size, the lower learning rate of  $2 \times 10^{-5}$  consistently resulted in better or comparable performance on both MMLU and MTBench benchmarks. For instance, with a batch size of 128, performances were similar for both the learning rates. For larger batch sizes of 3,840 and 7,680, the  $2 \times 10^{-5}$  learning rate performed on par or better than higher learning rates.

An explanation for what we observe is that large batches provide more accurate gradient estimates due to averaging over more samples. Consequently, a lower learning rate suffices to make effective progress without introducing instability. Using a higher learning rate with large batches may lead to larger updates that overshoot minima.

### A.5.4 TULU vs. LAB

**Memorization and Generalization.** We focused on Phase 05 in the sequential phased training strategy, which is designed to augment the model’s foundational knowledge and memorization of facts. We evaluated the models using specific MMLU subjects related to memorization, including history, law, and science domains. We compared the performance of models starting from both the base Granite model and the best checkpoint obtained from Phase 00 training. Table 8 shows that

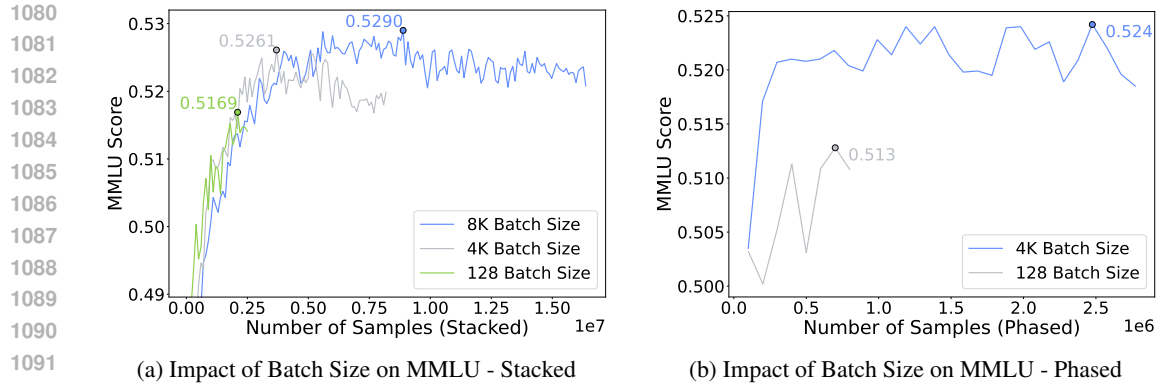


Figure 8: Impact of batch size on model performance in stacked and phased training on MMLU benchmark.

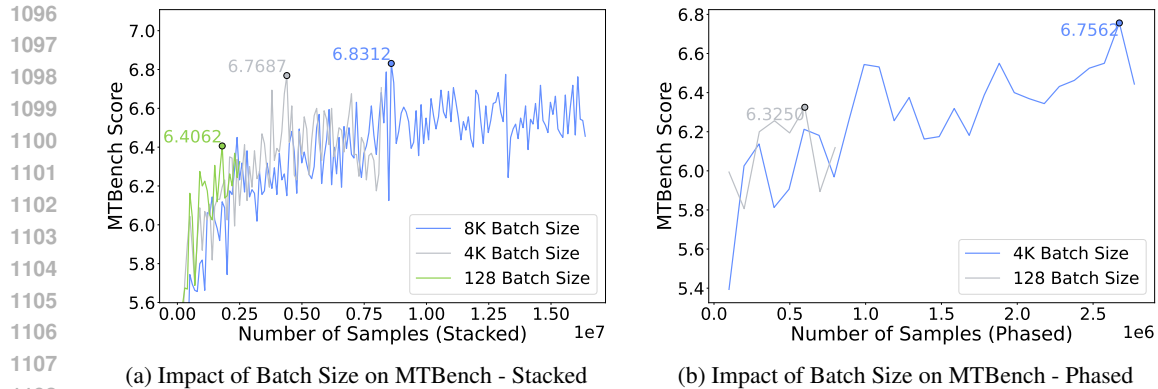


Figure 9: Impact of batch size on model performance in stacked and phased training on MTBench benchmark. MMLU results are in Appendix A.5.2.

models trained with LAB hyperparameters outperform those trained with TULU hyperparameters on the memorization-focused MMLU tasks. We evaluated the models' generalization abilities after Phase 10 in the sequential phased training strategy, which focuses on complex skills and compositional tasks. The MTBench benchmark was used to assess performance on tasks requiring reasoning, problem-solving, and adaptation to unseen scenarios. As shown in Table 8, the model trained with LAB hyperparameters significantly outperforms the one trained with TULU hyperparameters on MTBench.

Table 8: Comparison of TULU vs. LAB on MMLU Memorization and MTBench Generalization Scores. Cells highlighted in green indicate better scores, and blue indicates higher sample efficiency (fewer samples used).

Benchmark	Model	Score	Samples
MMLU (Memorization)	Granite Base	0.48	-
	TULU (Base)	0.59	199,936
	LAB (Base)	0.61	1,597,440
	TULU (Phase 00)	0.60	599,808
	LAB (Phase 00)	0.62	1,098,240
MTBench (Generalization)	TULU	6.33	599,808
	LAB	6.76	2,673,216

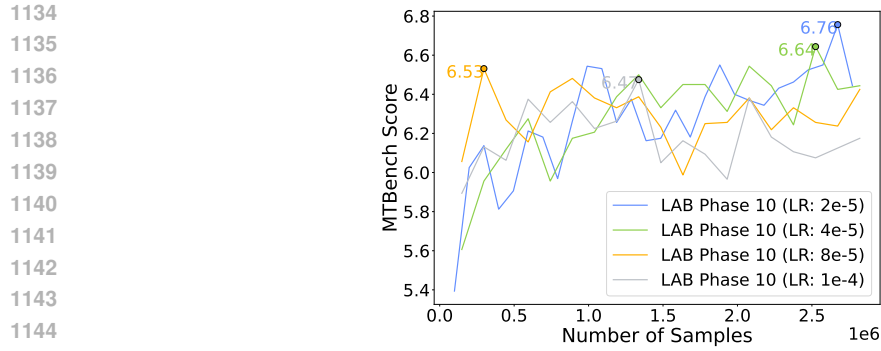
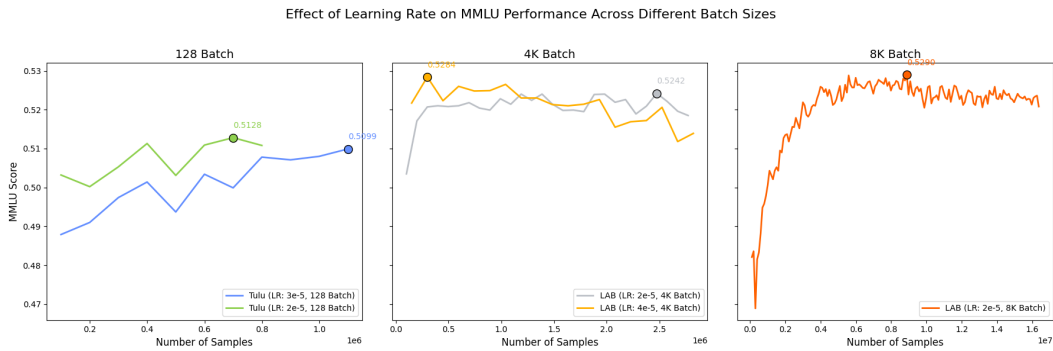
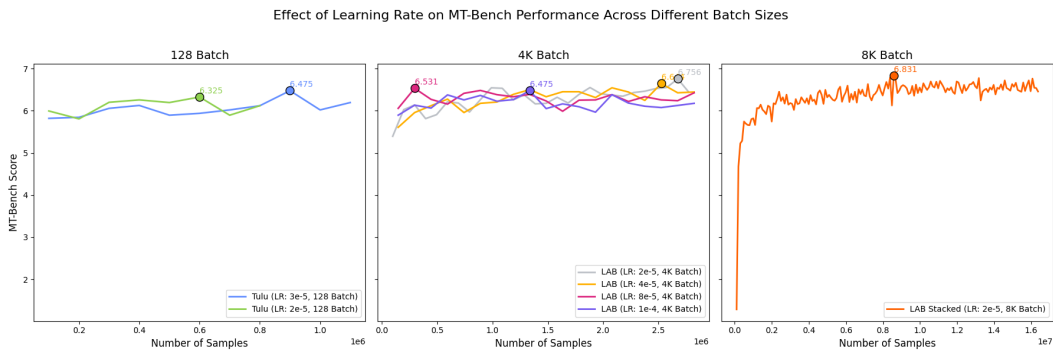


Figure 10: MTBench performance after Phase 10 training with different learning rates.



(a) MMLU Performance Across Learning Rates



(b) MTBench Performance Across Learning Rates

Figure 11: Performance comparison across different learning rates and batch sizes on MMLU and MTBench benchmarks.

To ensure a fair comparison, we ran each experiment for the same number of gradient steps, resulting in different number of samples due to different batch sizes, as seen in Figure 12. Figure 12a shows that models trained with LAB hyperparameters outperform those trained with TULU hyperparameters on the memorization-focused MMLU tasks. Additionally, as shown in Figure 12b, the model trained with LAB hyperparameters significantly outperforms the one trained with TULU hyperparameters on MTBench. Table 9 shows that LAB performs better than TULU across all benchmarks.

#### A.5.5 EFFECT OF LEARNING RATE SCHEDULES ON LARGE BATCH SIZES

As shown in Figure 13, the models trained with a constant learning rate (no decay) performed on par with those trained with a cosine decay schedule on both MMLU and MTBench, and in some cases even outperformed them, particularly on MTBench.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

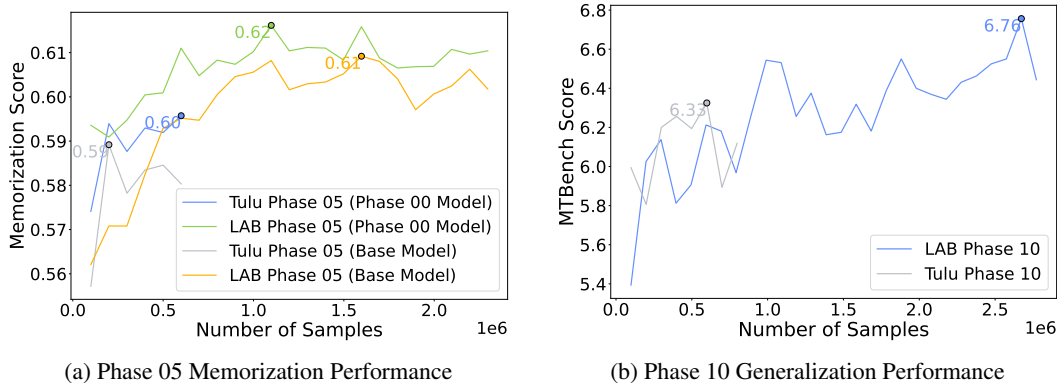


Figure 12: Comparison of TULU vs. LAB on memorization and generalization.

Table 9: Comparison of LAB vs. TULU Hyperparameter Configurations on Various Benchmarks. Cells highlighted in green indicate better scores, and blue indicates higher sample efficiency (fewer samples used).

Benchmark	Score			Samples	
	Granite Base	LAB	TULU	LAB	TULU
Leaderboard (BBH)	0.09	0.10	0.08	8,057,918	599,808
Leaderboard (MuSR)	0.01	0.07	0.03	8,057,918	599,808
ARC	0.78	0.75	0.74	8,057,918	599,808
GSM8K	0.11	0.37	0.36	8,057,918	599,808

### A.5.6 EFFECT OF WARMUP STEPS

Figures 14 and 15 show the performance comparison with different warmup steps: 0, 25, and 100 warmup steps, on the MMLU and MTBench benchmarks, respectively. The model trained without warmup steps achieved better performance on MMLU and similar performance on MTBench, indicating that warmup steps is not necessary for stable and effective training.

### A.5.7 ADAPTATION TO A DOMAIN-SPECIFIC DATASET

To evaluate the generalizability of our findings to domain-specific datasets, we conducted experiments using a Math, Reasoning, and Code (MRC) dataset. This dataset specializes in mathematical problem-solving, logical reasoning, and programming tasks, representing a focused domain compared to our original diverse dataset.

We evaluated the models on several benchmarks, including GSM8K (Cobbe et al., 2021), ARC (Clark et al., 2018), and the Open LLM Leaderboard v2 benchmarks including MATH and MuSR. We compare the LAB and TULU hyperparameter configurations on the MRC dataset. Using the LLaMA 3B model, we fine-tuned with both configurations: LAB used a batch size of 4,000 and a constant learning rate, while TULU used a batch size of 128 with warmup and linear decay. As shown in Table 10, the LAB configuration outperforms TULU across all evaluation metrics, reaffirming that larger batch sizes and simplified learning rate schedules are beneficial even when fine-tuning on domain-specific data.

Additionally, we fine-tuned the LLaMA 3B model using both the stacked and sequential phased training strategies with LAB hyperparameters. For phased training, as described in Appendix A.5.1, the dataset was partitioned into two phases based on response length: Phase I with shorter responses (bottom 50%) and Phase II with longer responses (top 50%). As shown in Table 11, stacked training demonstrates slightly higher performance and greater sample efficiency compared to phased training across all benchmarks.

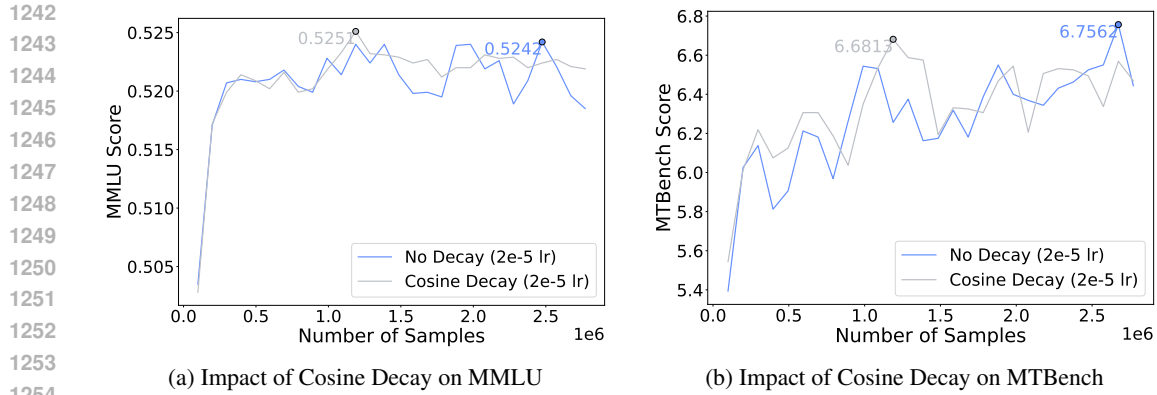


Figure 13: Comparison of learning rate schedules with a batch size of 3,840 samples on MMLU and MTBench benchmarks.

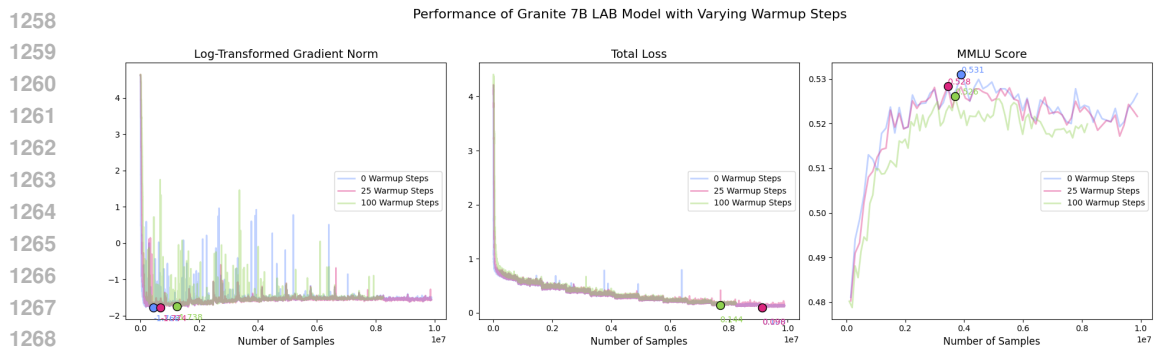


Figure 14: MMLU Performance with different warmup steps: 0, 25, and 100 warmup steps.

Table 10: Comparison of LAB vs. TULU Hyperparameter Configurations on the MRC Dataset Using the LLaMA 3B Model. Cells highlighted in green indicate better scores, and blue indicates higher sample efficiency (fewer samples used).

Benchmark	Score			Samples	
	LLaMA Base	LAB	TULU	LAB	TULU
Leaderboard (MATH Lvl 5)	0.02	0.04	0.04	9,980,259	3,468,664
Leaderboard (MuSR)	0.05	0.08	0.04	16,966,128	2,973,753
ARC	0.78	0.75	0.68	2,745,290	247,372
GSM8K	0.27	0.69	0.66	12,225,143	5,450,009

Table 11: Comparison of Stacked vs. Phased Training Strategies on the MRC Dataset Using the LLaMA 3B Model. Cells highlighted in green indicate better scores, and blue indicates higher sample efficiency (fewer samples used).

Benchmark	Score			Samples	
	LLaMA Base	Stacked	Phased	Stacked	Phased
Leaderboard (MATH Lvl 5)	0.02	0.04	0.03	9,980,259	18,455,850
Leaderboard (MuSR)	0.05	0.08	0.08	16,966,128	18,206,922
ARC	0.78	0.75	0.71	2,745,290	14,964,367
GSM8K	0.27	0.69	0.67	12,225,143	14,964,367

These findings demonstrate that our recommendations regarding training strategies and hyperparameters generalize to domain-specific datasets, supporting their applicability in specialized fine-tuning scenarios.



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307

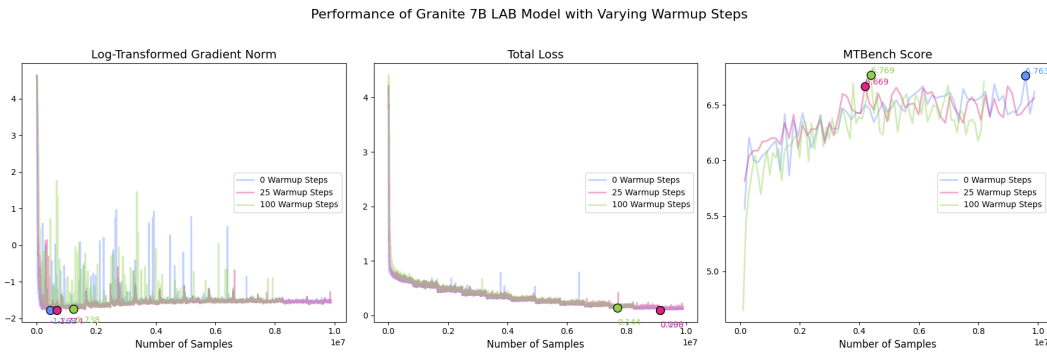


Figure 15: MTBench Performance with different warmup steps: 0, 25, and 100 warmup steps.

1308  
1309  
1310

#### A.5.8 ADAPTATIONS TO DIFFERENT MODEL SIZES AND ARCHITECTURES

1311  
1312  
1313  
1314  
1315

To assess the scalability and generality of our findings, we extended our experiments to different model families, architectures, and sizes, specifically testing the Mistral 7B model, the Granite 3B model, and the LLaMA 3B model.

1316  
1317  
1318  
1319

**Adaptation to a New Architecture.** We conducted a learning rate sweep for the Mistral 7B model to determine which learning rate yields the best final performance. Our objective was to apply the methodology used for finding the optimal learning rate with the Granite models to the Mistral architecture.

1320  
1321  
1322  
1323  
1324  
1325

This methodology involves starting with a low learning rate because lower learning rates have been found to stabilize training when fine-tuning language models. A low learning rate helps prevent large, destabilizing weight updates, allowing the model to fine-tune its parameters gradually and avoid overfitting. Additionally, lower learning rates facilitate more precise adjustments to the model weights, which is particularly important when adapting pre-trained models to new tasks or domains without forgetting previously learned information.

1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333

Then, once we begin with a low learning rate, the methodology/prescription we offer for finding the optimal hyperparameter—which in this case is the learning rate—is to test slightly higher and lower values to see if performance improves, and then keep following that path of lower or higher values until you don’t get anything better. That’s our prescription or methodology. We began the search around  $2 \times 10^{-5}$ , which worked best for Granite, and then tested slightly higher and lower values to see if performance improved. This approach allows us to prescribe a suitable range for the learning rate or any other hyperparameter and refine it based on empirical results. Through this sweep, we identified  $1 \times 10^{-6}$  as the best learning rate for Mistral.

1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344

Our results are presented in Figures 16a and 16b. As a baseline, the Mistral 7B pretrained model achieved a score of 0.60 on the MMLU benchmark in a 5-shot scenario, and the Mistral 7B Instruct model scored 6.84 on MTBench. We check if what we have observed before for Granite—that is, the general trend where lower gradient norms and higher loss are associated with better generalization and final performance—also applies to Mistral. Specifically, the lowest learning rate,  $1 \times 10^{-6}$ , produced the best overall performance on the MMLU benchmark. An interesting pattern emerged, similar to that observed with the Granite model: for the most effective learning rates, the gradient norm started at its lowest value and increased towards the end of training. Despite the higher gradient norm in the later stages, the associated loss remained higher throughout training (which is expected because lower learning rates typically result in higher loss values during training). This suggests that maintaining higher loss values may be necessary for the model’s optimization process to reach better generalization.

1345  
1346  
1347  
1348  
1349

These observations confirm that our methodology for selecting learning rates and our observation on the correlation between early training dynamics and final performance are applicable to different model architectures.

**Adaptation to Different Model Sizes.** We also examined whether our findings hold for smaller models by conducting experiments with the Granite 3B model. Specifically, we compared an 8k

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

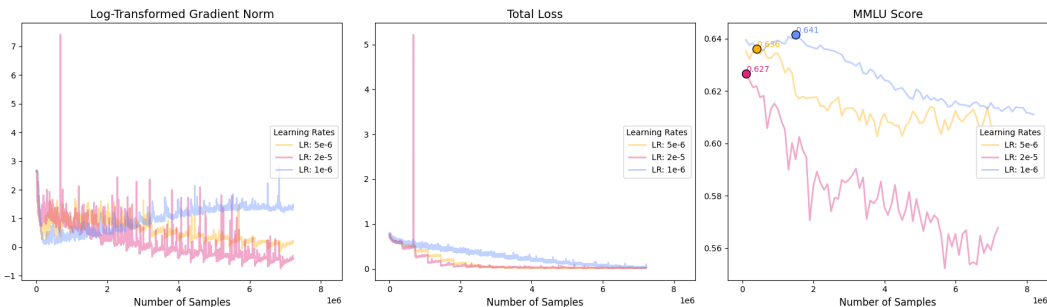
1400

1401

1402

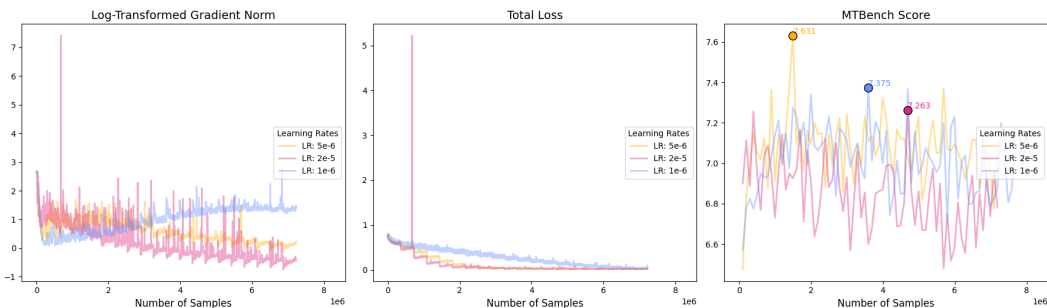
1403

Impact of Learning Rates on Training Dynamics (Grad Norm, Loss) and Final Model Performance (MMLU) for Mistral Model



(a) Training Dynamics and MMLU Performance

Impact of Learning Rates on Training Dynamics (Grad Norm, Loss) and Final Model Performance (MTBench) for Mistral Model



(b) Training Dynamics and MTBench Performance

Figure 16: Training dynamics for Mistral 7B with different learning rates, and their final performance on MMLU and MTBench benchmarks.

batch size with stacked training versus a 4k batch size with phased training. Our goal was to determine if the observations regarding batch size and training strategies for the Granite 7B model extend to the 3B model. In the 8k stacked setting for the Granite 3B model, we observed a lower gradient norm, higher loss, and improved MMLU performance compared to the 4k phased configuration. This trend is illustrated in Figures 17a and 17b.

The larger batch size likely benefits from increased data diversity within each batch, encompassing a wide variety of tasks, skills, and knowledge from the stacked phases. This diversity allows the model to learn more generalizable features, which likely contributes to the performance improvement. The lower gradient norm in the 8k stacked setting suggests more stable learning, while the higher loss indicates that the model is exploring a wider region of the loss landscape, potentially helping it avoid overfitting and converge to better generalization. This balance of stability (lower gradient norm) and exploration (higher loss) may explain why the 8k stacked configuration yields superior performance on the MMLU benchmark. These results suggest that the correlation between early training dynamics and final performance holds across different model sizes, indicating that early training dynamics are reliable predictors of final model performance regardless of model size.

These findings confirm that our earlier observations about training dynamics and performance correlations hold across different model sizes and architectures, reinforcing the importance of early training dynamics as indicators of final performance and the effectiveness of our methodology in selecting hyperparameters for fine-tuning language models.

**Generalization to a Different Model Family and Size.** To assess whether our findings extend to a different model architecture and size at the same time, we conducted experiments using the LLaMA 3B model (Touvron et al., 2023). We note that the Granite model shares the same architecture as the LLaMA model. Hence we believe that the findings in this paper can generalize across the LLaMA model family. We fine-tuned the model using both stacked and phased training strategies, as well as comparing the LAB and TULU hyperparameter configurations.



Figure 17: Training dynamics for Granite 3B with different batch sizes and training strategies, and their final performance on MMLU and MTBench benchmarks.

Table 12: Comparison of Stacked vs. Phased Training Strategies Using the LLaMA 3B Model. Cells highlighted in green indicate better scores, and blue indicates higher sample efficiency (fewer samples used).

Benchmark	Score			Samples	
	LLaMA Base	Stacked	Phased	Stacked	Phased
Leaderboard (BBH)	0.14	0.19	0.18	7,734,723	6,734,847
Leaderboard (MATH Lvl 5)	0.01	0.02	0.01	250,089	4,490,320
Leaderboard (MuSR)	0.05	0.22	0.11	10,979,309	4,988,983
MMLU	0.56	0.57	0.53	6,986,437	5,737,613
ARC	0.78	0.78	0.75	2,744,559	7,483,283
GSM8K	0.27	0.51	0.45	3,742,399	6,734,847
MTBench	-	5.00	4.30	9,232,227	6,734,847

For phased training, as described in Appendix A.5.1, the dataset was partitioned into two phases based on response length: Phase I with shorter responses (bottom 50%) and Phase II with longer responses (top 50%). We compared the LAB configuration (batch size of 4,000 with constant learning rate) to the TULU configuration (batch size of 128 with warmup and linear decay). The models were evaluated on benchmarks including MMLU, MTBench, GSM8K, ARC, and the Open LLM Leaderboard v2 benchmarks including BBH, MATH, and MuSR.

The results, depicted in Table 12, indicate that the stacked training strategy achieves better performance than phased training across all benchmarks. Results in Table 13 indicate that the LAB hyperparameter configuration consistently outperforms TULU, reinforcing our earlier conclusion that larger batch sizes and a constant learning rate schedule are advantageous. These findings suggest that our recommended training strategies and hyperparameters are effective across different model

Table 13: Comparison of LAB vs. TULU Hyperparameter Configurations on the LLaMA 3B Model. Cells highlighted in green indicate better scores, and blue indicates higher sample efficiency (fewer samples used).

Benchmark	Score			Samples	
	LLaMA Base	LAB	TULU	LAB	TULU
Leaderboard (BBH)	0.14	0.19	0.17	7,734,723	2,473,477
Leaderboard (MATH Lvl 5)	0.01	0.02	0.01	250,089	741,924
Leaderboard (MuSR)	0.05	0.22	0.15	10,979,309	1,731,217
MMLU	0.56	0.57	0.55	6,986,437	2,473,477
ARC	0.78	0.78	0.74	2,744,559	2,473,477
GSM8K	0.27	0.51	0.49	3,742,399	2,473,477
MTBench	-	5.00	4.97	9,232,227	2,473,477

architectures and sizes simultaneously, including the LLaMA family. Practitioners may consider applying these insights to fine-tune various small-sized LLMs, potentially achieving improvements in performance.

#### A.5.9 EARLY TRAINING DYNAMICS AS PREDICTORS OF FINAL PERFORMANCE

In addition to the MTBench results presented in the main paper, we include the MMLU performance comparison here.

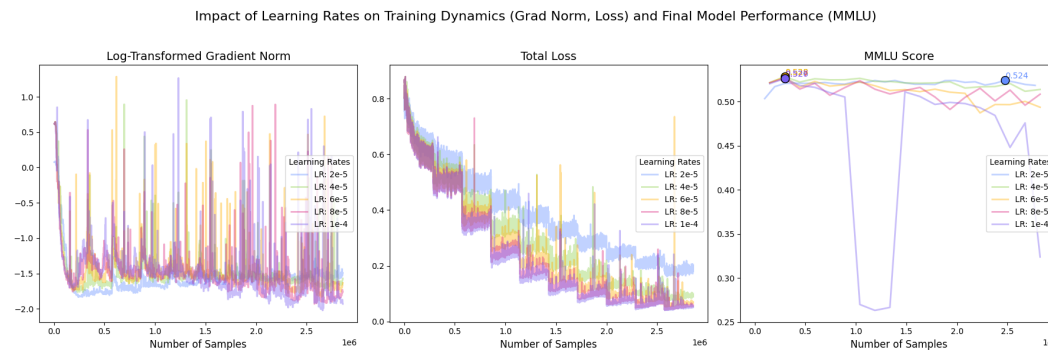


Figure 18: LAB Learning Rate Sweep: Training Dynamics and MMLU Performance.

Models trained with a learning rate of  $2 \times 10^{-5}$  demonstrated lower gradient norms initially, which increased toward the end of training, and higher loss throughout, ultimately resulting in superior final performance on MMLU compared to models trained with higher learning rates. This pattern mirrors the observations made for MTBench in the main paper, reinforcing the correlation between early training dynamics and final performance across different benchmarks.

Figures 19, 14, 15, and 20 illustrate the correlation between early training dynamics—gradient norms and loss values—and final benchmark performances across other experiments:

- **Batch Size Comparison in Stacked Training.** We further examined the impact of batch size on early training dynamics and final performance in stacked training by comparing batch sizes of 3,840 and 7,680 samples (denoted as 4k and 8k). Figures 19a and 19b present the gradient norms and loss values over training, along with the corresponding performances on MMLU and MTBench. We observed that the 8k batch size consistently exhibited lower gradient norms and higher loss throughout training compared to the 4k batch size. Ultimately, the 8k batch size achieved better final performance on both MMLU and MTBench benchmarks.

The larger batch size likely benefits from increased data diversity within each batch, encompassing a wide variety of tasks, skills, and knowledge from the stacked phases. This diversity allows the model to learn more holistically and make more informed optimization steps, enhancing its

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

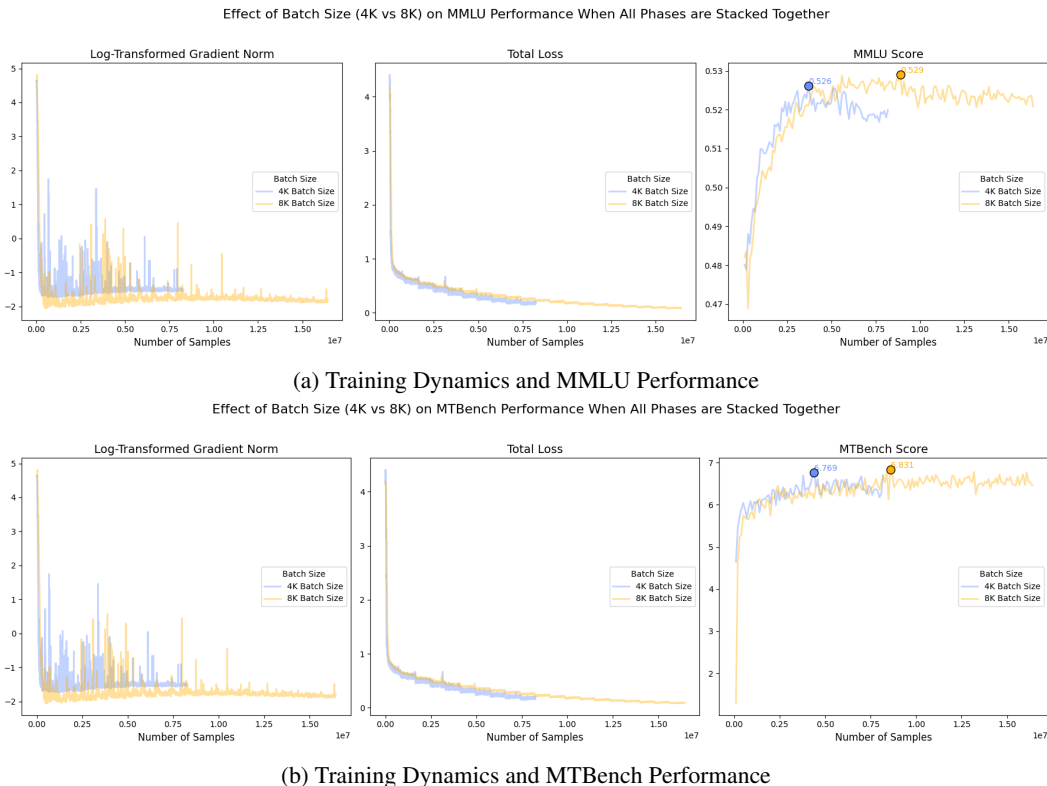


Figure 19: Training dynamics for batch sizes 4k and 8k in stacked training, and their final performance on MMLU and MTBench benchmarks.

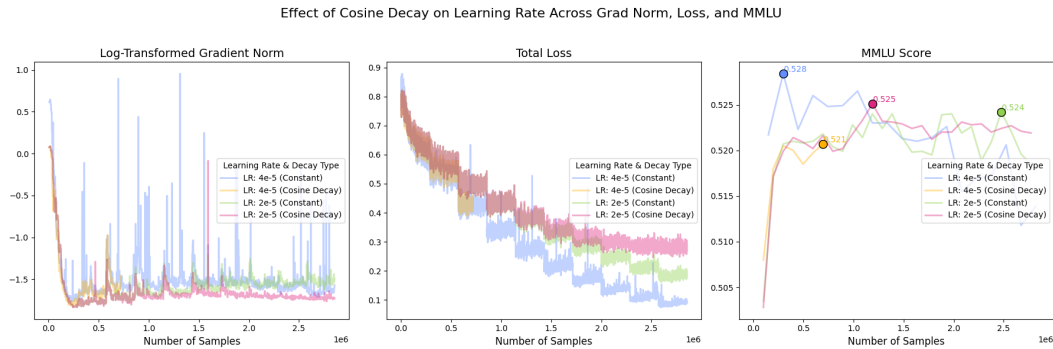
ability to generalize across different domains. However, we also noted that for smaller numbers of training samples, the 4k batch size achieved higher MMLU and MTBench scores, suggesting a trade-off. If computational resources or training time are limited, the 4k batch size may offer better performance in the early stages of training—up to approximately 3.75 million samples for MMLU and 8.5 million samples for MTBench. Beyond these points, the 8k batch size surpasses the 4k batch size in performance as observed in Figures 19a and 19b.

- **Warmup Steps Comparison.** We examined the impact of different warmup configurations (0, 25, and 100 steps) on early training dynamics and final performance. All models demonstrated very similar performance, loss curves, and gradient norms throughout training. The model trained without warmup steps (0 warmup) achieved slightly better performance on the MMLU benchmark and comparable performance on MTBench compared to models trained with 25 or 100 warmup steps.

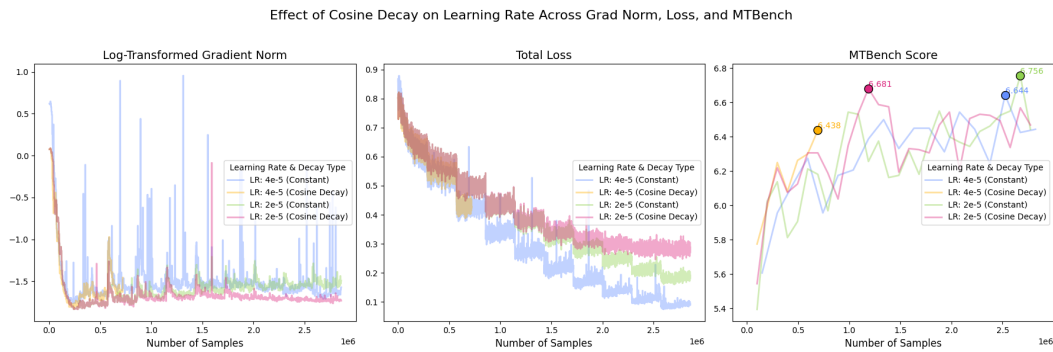
Gradient norm and loss curves serve as a proxy for the smoothness of the optimization process. Large fluctuations in early gradnorm values may indicate instability, which could negatively affect convergence, while more stable or lower gradnorm magnitudes suggest a smoother path toward optimal performance. Given that all warmup configurations resulted in similar final performance across both benchmarks and exhibited nearly identical loss and gradnorm curves, it indicates a strong correlation between training dynamics and final performance which can be seen in Figures 14 and 15.
- **Learning Rate Schedule Comparison.** We also analyzed the effect of using a constant learning rate versus a cosine decay schedule with learning rates of  $2 \times 10^{-5}$  and  $4 \times 10^{-5}$  on early training dynamics and final performance. Figures 20a and 20b present the gradient norms and loss values over training, along with the corresponding performances on MMLU and MTBench. Up until approximately 1 million samples, the gradient norms, loss values, and MMLU scores remain nearly identical for both the constant learning rate and cosine decay configurations, as the learning rate decay has not yet kicked in. However, after 1 million samples, the models with cosine decay

exhibit marginally better or on par MMLU/MTBench performance compared to those with a constant learning rate. We also observe that, beyond the 1 million sample mark, cosine decay results in lower gradient norms and higher loss values.

The lower gradient norms suggest more stable learning, while the higher loss may indicate that the model is exploring a broader region of the loss landscape, potentially enhancing generalization. This combination might explain the slight improvement in performance with cosine decay, as the balance between stability (lower grad norm) and exploration (higher loss) can help the model converge to a more optimal solution.



(a) Training Dynamics and MMLU Performance



(b) Training Dynamics and MTBench Performance

Figure 20: Training dynamics for constant learning rate vs. cosine decay, and their final performance on MMLU and MTBench benchmarks.

#### A.5.10 GRADIENT ACCUMULATION EQUIVALENCE TO FULL BATCH TRAINING

We investigated whether using gradient accumulation on a single node with a large batch size is equivalent to distributed training across multiple nodes with the same effective batch size. Theoretically, both methods should yield identical training dynamics and result in the same fine-tuned model if implemented correctly.

In our experiments, we compared two setups:

- **Single Node with Gradient Accumulation.** We utilized a single node with gradient accumulation to achieve an effective batch size corresponding to 60,000 tokens.
- **Multi-Node Distributed Training.** We employed distributed training across four nodes, maintaining the same effective batch size of 60,000 tokens without gradient accumulation.

We evaluated both setups by comparing their training loss curves, as well as performance on the MMLU and MTBench benchmarks. The results showed that the loss trajectories were virtually identical between the two methods. Additionally, the final performances on MMLU and MTBench were the same within experimental variance. These findings confirm that gradient accumulation on a

1620 single node can replicate the training dynamics and outcomes of full-batch distributed training across  
1621 multiple nodes. This equivalence provides flexibility for practitioners with limited computational  
1622 resources, allowing them to achieve the same model quality using gradient accumulation on fewer  
1623 GPUs.  
1624

1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673