

Pre-training Small Base LMs with Fewer Tokens

Sunny Sanyal*
UT Austin

Sujay Sanghavi
UT Austin

Alexandros G. Dimakis
UT Austin

Abstract

We study the effectiveness of a simple approach to develop a small base language model (LM) starting from an existing large base LM: first inherit a few transformer blocks from the larger LM, and then train this smaller model on a very small subset (0.1%) of the raw pretraining data of the larger model. We call our simple recipe Inheritune and first demonstrate it for building a small base LM with 1.5B parameters using 1B tokens (and a starting few layers of larger LM of 3B parameters); we do this using a single A6000 GPU for less than half a day. Across 9 diverse evaluation datasets as well as the MMLU benchmark, the resulting model compares favorably to publicly available base models of 1B-2B size, some of which have been trained using 50-1000 times more tokens.

We investigate Inheritune in a slightly different setting where we train small LMs utilizing larger LMs and their full pre-training dataset. Here we show that smaller LMs trained utilizing some of the layers of GPT2-medium (355M) and GPT2-large (770M) can effectively match the val loss of their bigger counterparts when trained from scratch for the same number of training steps on OpenWebText dataset with 9B tokens. We analyze Inheritune with extensive experiments and demonstrate its efficacy on diverse settings. Our code is available at <https://github.com/sanyalsunny111/LLM-Inheritune>.

1 Introduction

Scaling large language models (LLMs) is very challenging both in terms of the required computation resources but also in terms of the size and quality of the available datasets. The standard pre-training recipe involves working with multi-billion parameter models and training with trillion(s) of tokens (Kaplan et al., 2020; Hoffmann

et al., 2022; Peiyuan Zhang and Lu, 2023; Touvron et al., 2023b). In this paper we investigate what can be achieved when pre-training small base language models (LMs) of size 1.5B parameters with just 1 billion tokens of pre-training data. A base LM is a decoder style LLM that is solely trained for next token prediction without any subsequent instruction tuning, RLHF (reinforcement learning with human feedback), etc. In this paper, we focus on crafting/developing small base LMs from existing large reference LMs, when we have only a very small fraction of the pre-training data of the reference model available in public. This setting is important because there are many large models for which the weights are available, but the full pre-training datasets are not. Our method, is very simple: we keep the first few transformer blocks (layers) of the large reference LM and perform pre-training on the small available dataset, iterating for multiple epochs over the data. Our method is efficient in data and compute: the training was performed with 1 A6000 GPU in less than 12 hours (half the time (Geiping and Goldstein, 2022) used for BERT pre-training from scratch). Next we thoroughly analyzed our recipe in three different settings 1) where we start from large and more performant reference models, 2) we train model of different sizes with the available 1B data, 3) we trained our 1.5 billion parameter model using slightly larger subsets (1-5%) of pre-training data, ranging from 10 to 50 billion tokens. We then compared its performance with that of a model trained on 1B tokens of available data, but repeated for multiple epochs.

Summary. In summary, our key findings are as follows.

1. We introduce a pre-training technique called Inheritune. This method involves training a small size target model, extracted from a larger reference LM. We use OpenLLaMA-3B as our reference model and utilize few

*The corresponding author can be reached out at sanyal.sunny@utexas.edu.

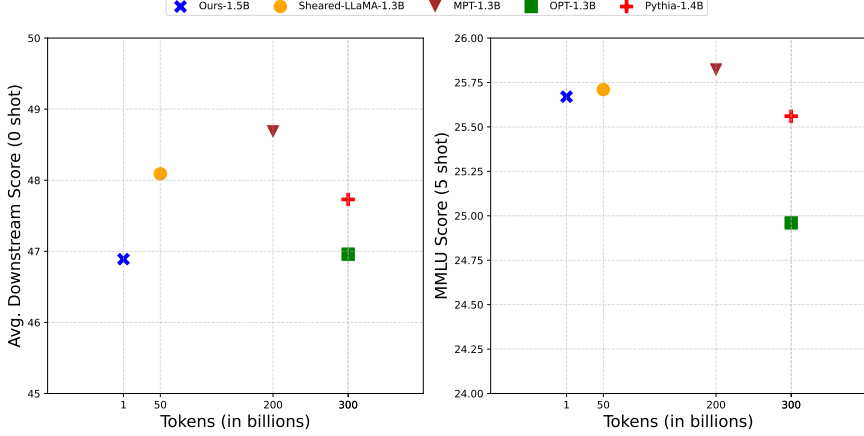


Figure 1: Performance of our 1.5B base LM derived using 1B data with Inheritune on an average of 9 different datasets (left) and MMLU benchmark (right) that evaluates commonsense, truthfulness, natural language inference and language understanding. We compare our model’s performance with reference model-OpenLLaMA-3B (2x size) and other small base LMs of size 1B-2B parameters such as MPT-1.3B, OPT-1.3B, Pythia-1.4B (pre-trained from scratch) and ShearLLaMA-1.5B (continually trained using existing large base LM).

of the transformer blocks (layers) to obtain our target model of size 1.5B. Subsequently the target model was further trained on 1B of training data for 8 epochs. Our small base LM achieves 89% of average downstream accuracy on 9 different datasets and 94% of MMLU (5-shot) score compared to the reference model which is double in size. Further, our small model achieves comparable performance on average downstream accuracy of 9 tasks and MMLU 5-shot score compared to publicly available baseline models pre-trained with 50-300 times more data as shown in Figure 1. We empirically analyze Inheritune with a 50B subset of pre-train data and larger reference models to demonstrate the generality of our approach.

2. We analyze Inheritune in a slightly different setting where we assume full access to the pre-training data. We observe that a much smaller target model can be extracted if we use the full pre-training data. We ran controlled experiments with GPT2-large and GPT2-medium LLMs. Utilizing Inheritune we show that for GPT2-large we can keep 50% of the layers and 45% of parameters while for a GPT2-medium, we keep 33% layers and 28% parameters without compromising the validation loss (log perplexity) as shown in Table 6. Intriguingly, we also observe that these smaller models derived with Inheritune exhibit lower validation loss to their same-sized counter-

parts are trained from scratch for 2x the number of training steps. Moreover, these smaller models derived with Inheritune exhibit a similar convergence pattern to their larger counterparts, as depicted in Figure 5.

Algorithm 1 Inheritune

Require: Reference model \mathcal{M}_{ref} , a subset $\hat{\mathcal{D}}_{train}$ from \mathcal{D}_{train} used to train \mathcal{M}_{ref} , the number of layers n we want in our final model, number of epochs E , and training hyper-parameters.

- 1: Let k be the number of transformer blocks in \mathcal{M}_{ref} , with the corresponding parameters being $\theta_0, \theta_1, \dots, \theta_{k-1}$.
 - 2: **Initialize:** target LM θ_{tgt} with the first n layers $\theta_0, \theta_1, \dots, \theta_{n-1}$, as well as the token embedding and prediction head layers of \mathcal{M}_{ref} .
 - 3: **Train** the target LM from this initialization for E epochs using the given hyper-parameters.
-

2 Method: Inheritune

In this section, we delve into the details of our proposed pre-training approach, Inheritune. Here we discuss Inheritune methodology in a setting where a small fraction of the pre-training data is available along with an existing large base LM. Later, we also investigate Inheritune where we assume full access to the pre-training data in Section 4.

Setup. We assume that there exists a pre-trained reference model \mathcal{M}_{ref} , comprising k layers, rep-

| Models | Layers | Hidden Size | Heads |
|-------------------|--------|-------------|-------|
| OpenLLaMA-3B | 26 | 3200 | 32 |
| OpenLLaMA-7B | 32 | 4096 | 32 |
| LLaMA2-7B | 32 | 4096 | 32 |
| GPT2-large(770M) | 36 | 1280 | 20 |
| GPT2-medium(355M) | 24 | 1024 | 16 |

Table 1: Overview of reference models used in this study and their architectural configurations. We obtain a pre-trained OpenLLaMA-3B and we trained all GPT2 models with OpenWebText consisting of 9B tokens.

| Model | Training Data (# tokens) |
|----------------------|--------------------------|
| OpenLLaMA-3B v1(ref) | RedPajama(1T) |
| Ours-1.5B* | RedPajama (1B) |
| Shear-LLaMA-1.3B* | RedPajama(50B) |
| MPT-1.3B | RedPajama(200B) |
| Pythia-1.4B | The Pile(300B) |
| OPT-1.3B | Custom data(300B) |

Table 2: Overview of the baseline models, pre-train data, and number of training tokens used to train these models.

resented by $\theta_{\text{ref}} = \{\theta_0, \theta_1, \dots, \theta_{k-1}\}$ trained with $\mathcal{D}_{\text{train}}$, however this full training data is unavailable and we only have a very small subset $\hat{\mathcal{D}}_{\text{train}}$. Moreover we assume that the distribution is preserved such that $\hat{\mathcal{D}}_{\text{train}} \sim \mathcal{D}_{\text{train}}$. We make no special efforts in selecting a high quality subset rather $\hat{\mathcal{D}}_{\text{train}}$ is just a random subset of all the domains of $\mathcal{D}_{\text{train}}$.

Step 1: Inherit the first n layers of \mathcal{M}_{ref} . We initialize the target model \mathcal{M}_{tgt} , with the first n layers of the \mathcal{M}_{ref} ; thus the weights of \mathcal{M}_{tgt} are $\{\theta_0, \theta_1, \dots, \theta_{n-1}\}$. The prediction head and token embeddings are also inherited from \mathcal{M}_{ref}

Step 2: Train \mathcal{M}_{tgt} for the given number E of epochs on the available training data $\hat{\mathcal{D}}_{\text{train}}$.

3 Developing a small base LM in low data regime using Inheritune

We now outline the setup and the main results for our first setting where using a very small fraction of pre-training data i.e. just 1B tokens (1000x small compared to original training data) to craft a 1.5B target LM which is competitive with existing base LMs both trained from scratch or derived using existing large base LMs of similar size.

3.1 Experimental Setup

Data. We utilize a 1B subset¹ of the Redpajama v1 dataset (Computer, 2023) which originally consists of 1T tokens from 7 different domains namely common crawl, c4, wikipedia, books, arxiv papers, github and stackexchnge. We combine the dataset in the proportions suggested in LLaMA by Touvron et al. (2023a), which is also the same set for proportions that was used for pretraining our reference LM: OpenLLaMA-3B. Our pre-training data thus constitutes a randomly selected 0.1% subset of OpenLLaMA-3B’s training dataset. Next for an extended evaluation to analyze the effects of Inheritune while using a larger proportion of the same data we randomly sample 50B subset (i.e. 5% subset) from Redpajama v1 dataset.

Model and Training. We use OpenLLaMA-3B version1 (Geng and Liu, 2023) as the reference model for our main results to efficacy of the Inheritune recipe. This model has $k = 26$ layers; we report comparisons for a choice of $n = 13$. That is, our model derived using Inheritune has 13 transformer blocks². We train for 8 epochs (where each epoch uses all the 1B tokens) with a batch size of 131K tokens per batch. The choice of training epochs are made based on analysis done in Figure 4. We use lit-gpt framework (AI, 2023) for training all small base LMs discussed in this paper. Further discussions on the training hyper-parameters can be found in the appendix. For an extended evaluation of our results we also employ larger reference LMs namely LLaMA2-7B (Touvron et al., 2023b) and OpenLLaMA-7B version 1 (Geng and Liu, 2023).

Baseline Models. Table 2 describes the models we compare against; all of them are publicly available checkpoints of similar size to our 1.5B derived model, following related prior work (Jawaheripi et al., 2023; Xia et al., 2023; Li et al., 2023; Peiyuan Zhang and Lu, 2023). However there are a few differences between them in how they are pre-trained.

We use pre-training data as a key criteria for selecting the baseline models. We pre-dominantly select publicly available baseline models based on their number of parameters and pre-training data i.e. redpajama. As the quality of the pre-training data plays a key role in model development. For

¹<https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T-Sample>

²We also include a study of the performance for other choices of n in Figure 2.

| Model | | Commonsense Reasoning | | | | |
|--------------------------|--------------|-----------------------|-------|--------------|------------|--------|
| Name (# train tokens) | Reference | Winograd | PIQA | Boolq | WinoGrande | Logiqa |
| OpenLLaMA-3B (1T) | n/a | 63.46 | 74.97 | 67.18 | 62.27 | 28.4 |
| OPT-1.3B (300B) | n/a | 38.46 | 71.82 | 57.83 | 59.51 | 27.04 |
| Pythia-1.4B (300B) | n/a | 36.54 | 70.89 | 63.12 | 56.99 | 27.65 |
| MPT-1.3B (200B) | n/a | 63.46 | 71.44 | 50.89 | 58.09 | 28.26 |
| Sheared LLaMA-1.3B (50B) | LLaMA2-7B | 36.54 | 73.45 | 62.02 | 58.17 | 27.34 |
| Ours-1.5B (1B) | OpenLLaMA-3B | 50.96 | 56.47 | 61.68 | 51.69 | 25.19 |

| Model | | Lang. Understanding & Inference | | | | Factuality |
|--------------------------|--------------|---------------------------------|--------------|--------------|--------------|--------------|
| Name (# train tokens) | Reference | MMLU(5) | WNLI | QNLI | MNLI | TruthfulQA |
| OpenLLaMA-3B (1T) | n/a | 27.21 | 50.7 | 51.3 | 37.3 | 35 |
| OPT-1.3B (300B) | n/a | 24.96 | 42.25 | 51.29 | 35.82 | 38.67 |
| Pythia-1.4B (300B) | n/a | 25.56 | 53.52 | 49.48 | 32.76 | 38.66 |
| MPT-1.3B (200B) | n/a | 25.82 | 40.85 | 50.52 | 35.93 | 38.68 |
| Sheared LLaMA-1.3B (50B) | LLaMA2-7B | 25.71 | 49.3 | 50.98 | 37.94 | 37.14 |
| Ours-1.5B (1B) | OpenLLaMA-3B | 25.67 | 43.66 | 49.41 | 34.42 | 48.61 |

Table 3: Comparison of our target model (\mathcal{M}_{tgt}) derived using Inheritune with reference model (\mathcal{M}_{ref}) and other baseline models of similar size when pre-trained from scratch and pre-trained with inherited weights and pruning. Our model although trained with fewer tokens achieves comparable performance compared to the baseline models. We have highlighted all the scores in **bold** where Our-1.5B model achieves at least 90% of the score compared it’s reference LM or it outperforms at least two of the publicly available baseline LMs. All the tasks are evaluated using 0 shot except MMLU which is 5-shot. The models against which n/a is mentioned are trained from scratch.

the purpose of fair comparison we choose MPT-1.3B (Author, s), Sheared LLaMA=1.3B (Xia et al., 2023) trained with redpajama. Next we also include models OPT-1.3B (Zhang et al., 2022) and Pythia-1.3B (Biderman et al., 2023) as these models are pre-trained with dataset which are known to be close with redpajama. Sheared LLaMA-1.3B is developed with targeted structural pruning with 0.4B tokens and continual pre-training using 50B data. Pruning with training has shown to have some added benefits towards overall generalization (Jin et al., 2022). Therefore it’s a bit unfair to directly compare a pruning and continual training method with our recipe which can be better classified as a model initialization technique. However both Shearing and Inheritune start training using an existing model and because of that we anyway make the comparison while treating it as a special case that we discuss separately in Section 3.2.

Evaluation. In this study, we use few-shot accuracy particularly 0-shot and 5-shot accuracy on multiple different downstream tasks to measure the quality of base LMs. This is a standard practice of evaluating pre-trained LLMs and has been done in several prior works (Brown et al., 2020; Touvron et al., 2023a; Javaheripi et al., 2023; Biderman et al., 2023). Our evaluation methodology

spans downstream tasks across four distinct genres: commonsense reasoning, natural language understanding, factuality, and natural language inference. We evaluate commonsense reasoning using 0-shot performance on 5 diverse datasets—PIQA (Bisk et al., 2020), BOOLQ (Clark et al., 2019), WINOGRANDE (Sakaguchi et al., 2020), WINOGRAD (Kocijan et al., 2020) and LOGIQA (Liu et al., 2020). We evaluate language understanding using 5-shot performance on massive multitask language understanding dataset (MMLU) (Hendrycks et al., 2020). Factuality is a key concern in LLMs and we evaluate factuality using 0-shot performance on TruthfulQA (Lin et al., 2022) data. Finally we evaluate language inference capabilities using 0-shot performance on MNLI (Bowman et al., 2015), QNLI (Wang et al., 2018) and WNLI (Wang et al., 2018) datasets. For the entire evaluation we use the lm eval harness framework (Gao et al., 2023).

3.2 Results: Inheritune with 1B data

We provide a visual summary of our main results in Figure 1, accompanied by a comprehensive evaluation of performance, broken down by tasks, in Table 3.

Comparison to baseline models available publicly. We compare against base LMs of size similar to ours but which have been trained from

scratch. As seen in Table 2, this lack of advantage is reflected in the remarkably larger number of tokens they are trained on. We show this comparison mainly to illustrate that when one is fortunate enough to have a large reference base model and a subset of its pre-training data, inheriting the target size base LM with a subset of the layers of the reference LM **makes the training remarkably more sample efficient than pre-training from scratch**. From Table 2, it is evident that our 1.5B model, developed using Inheritune, excels in 7 out of 10 individual tasks. It achieves a score of 90% or higher compared to the reference language model, which is twice its size and trained with 1000 times more data, or it outperforms at least two other base LMs of similar size trained with 50-300 times more data. Scores deemed favorable by us have been highlighted in bold. Next we compare our small LM with MPT-1.3B³ model trained from scratch with 200B tokens of redpajama and observe that we match 97% accuracy in all 9 downstream tasks and matches the MMLU (5-shot) score. To this end we compare with 2 small base LMs (OPT-1.3B and Pythia-1.3B) trained with large data corpses other than redpajama and show that we outperform both in MMLU (5-shot) score and perform comparably in other 9 datasets.

Comparison with a pruning technique ShearedLLaMA (Xia et al., 2023). Our training recipe can be better classified as an initialization technique than a pruning technique. There are some major differences of our recipe with pruning.

a) Pruning does not always require re-training (Li et al., 2020; Xia et al., 2022), whereas any initialization based recipe requires training. However, some prior works opt to continue train after pruning as an additional step, such as shearing (Xia et al., 2023). b) Typically, pruning techniques adhere to the "train large, then compress" paradigm (Li et al., 2020), which requires starting from larger models, whereas Inheritune is a training recipe that pre-trains a model from ground up, utilizing a better initialization than random. c) Targeted structural pruning is a carefully designed technique based on some weight update criteria hence it requires a lot of compute to perform, Shearing (Xia et al., 2023) took 8 GPUs, 0.4B tokens and 3k steps for just pruning. Later it also continually trains the pruned model for 50B

³<https://huggingface.co/mosaicml/mpt-1b-redpajama-200b>

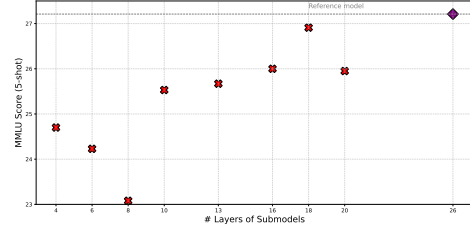


Figure 2: **Inheritune scales accross multiple different model sizes.** Utilizing the OpenLLaMA-3B as a reference large base LM we demonstrate that multiple performant small base LMs (submodels) of target size can be crafted using the Inheritune methodology with just 1B training tokens. Here we measure target size using number of layers used to initialize the submodels in step 2 of the Inheritune recipe and the performance using MMLU (5-shot) benchmark scores.

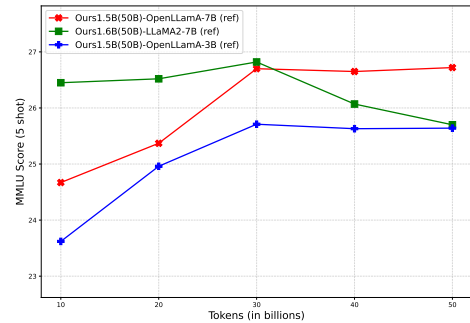


Figure 3: Performance of our 1.5B base LM derived using Inheritune on an average of 9 different datasets (left) and MMLU benchmark (right) that evaluates commonsense, truthfulness, natural language inference and language understanding. We present results with three different reference models OpenLLaMA-7B, LLaMA2-7B and OpenLLaMA-3B.

tokens using 16 GPUs. Our recipe performs training with 1 A6000 GPU for less than 12 hours using a total 1B data to develop a small base LM. Overall although it's not fair to compare a targeted structural pruning technique with an initialization recipe but for the sake of completeness we still perform the comparison. Here in Table 2 we observe that Inheritune shows competitive scores compared to Sheared LLaMA-1.3B in MMLU, Winograd, Boolq, QNLI and truthfulqa (i.e. 5/10 tasks) despite being trained with 50 times less data and much larger compute requirement.

3.3 Inheritune scales accross different model sizes

For the results in the previous section we only considered a single choice $n = k/2$ i.e. half the layers – for the size of the smaller model. In this section

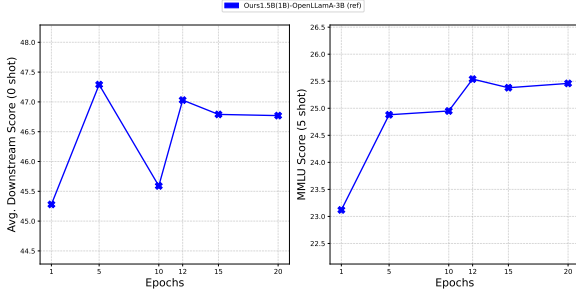


Figure 4: Performance of our 1.5B base LM derived using Inheritune based on existing OpenLLaMA-3B base model. Here we use 1B tokens and perform data repetition during continual training. We further evaluate our model on an average of 9 different datasets (left) and MMLU benchmark (right) that evaluates commonsense, truthfulness, natural language inference and language understanding.

we do a preliminary investigation into Inheritune with different choices of n (but the same 1B token dataset). All models in this section also use the OpenLLaMA-3B as the large pretrained reference model, and also share the same choices for all training hyperparameters – the only thing changing is the choice of n .

We developed 8 different submodels with $n = \{4, 6, 8, 10, 13, 16, 18, 20\}$. Figure 2 depicts the MMLU (5-shot) score as a function of n . As can be expected, we see a positive-sloping trend line. The submodel with 20 layers, exhibits a slight decrease in performance, potentially due to data overfitting issues due to increase in model size. The training details of all these submodels remain consistent with the target 1.5B small base LM and mentioned in appendix. We leave a more comprehensive investigation on the choice of n – including possibly varying both n and the number of training tokens jointly, and evaluating on a more exhaustive set of tasks – as future work.

3.4 Additional analysis with larger reference LMs and 50B data

We further analyze Inheritune to see the impact of its performance when more tokens are available. Initially for the main results we limited ourselves to 1B (i.e. 0.1%) tokens from the 1T pre-training data, here we use a 50B subset (i.e. 5%) of the pre-train data. Moreover we also extend this study to include larger base LMs of 7B parameters as reference models, employing OpenLLaMA-7B and LLaMA2-7B as reference models. For the purpose of this study we do not repeat the tokens from our

| Model (# tokens), ref | MMLU(5) |
|--------------------------------------|--------------|
| Ours-1.6B (1B), LLaMA2-7B | 24.27 |
| Ours-1.5B (1B), OpenLLaMA-3B | 25.67 |
| Ours-1.5B (50B), OpenLLaMA-3B | 25.71 |
| Ours-1.6B (50B), LLaMA2-7B | 26.07 |
| Ours-1.6B (50B), OpenLLaMA-7B | 26.72 |

Table 4: Performance comparison of models on the MMLU (5-shot) task. Our models, even when trained with fewer tokens, show competitive performance compared to benchmarks. We have highlighted the best MMLU 5-shot score in **bold**.

| Model (# tokens) | Data type | MMLU (5-shot) |
|------------------------|-----------|---------------|
| Ours-1.5B (1B) | 10 epochs | 24.95 |
| Ours-1.5B (50B) | 10B fresh | 23.62 |
| Ours-1.5B (1B) | 20 epochs | 25.46 |
| Ours-1.5B (50B) | 20B fresh | 24.96 |

Table 5: Zero-shot performance on MMLU (5-shot) task of Our-1.5B small base LM derived using 1B data for multiple data repetition–10 epochs and 20 epochs compared to the same model trained without data repetition for 10B and 20B fresh tokens. We derive all the variants of Our-1.5B small base using Inheritune with OpenLLaMA-3B as reference model. The models featured in this table correspond to those discussed in Figures 3 and 4.

50B subset. As shown in Figure 3, we observe that there is clear improvement in overall MMLU (5-shot) score with more data. Additionally it is interesting to see that 1.5B (or 1.6B models) developed with Inheritune using larger reference models show even greater improvements when fed with 50B subset of non repetitive data (i.e fresh tokens). We present a Table 4 using Figure 3 to show the best MMLU (5-shot) scores achieved using different reference LMs. For developing our small base LMs using larger reference LMs we use $n=7$ (i.e. 7 layers). The training details are mentioned in supplementary materials.

Ablations with number of epochs. We ran ablations (refer Figure 4) to choose the total number of epochs (multiple passes over the data) and observe that repetition when training our 1.5B (or 1.6B) LM is helpful particularly for MMLU. We also observe that the for an average of all the 9 other datasets (i.e. except MMLU) peaks its performance at 5 epochs and then deteriorates. Some prior works have studied this phenomenon that the scaling of downstream tasks with data is not always linear

(Biderman et al., 2023).

To repeat or not to repeat the tokens. Next we tackle the question – whether one should re-use 1B tokens for multiple epochs or use the same number of fresh tokens? Some prior works have recommended that if you have a reasonably large size dataset one can repeat it upto 4 epochs (Muenighoff et al., 2023). In our study we observe that one can safely re-use 1B tokens upto 10-20 epochs as shown in Table 5. We emphasize that this phenomenon needs a through investigation in itself and we defer this to future work. The models discussed in Table are saved checkpoints during a single training run and not the final model unless otherwise specified.

4 Exploratory Analysis of Inheritune in the presence of full pre-training data.

In this section we study the performance of Inheritune in a slightly different setting from that of the previous section. In particular, in this section we look at the case where we not only have the larger reference base model \mathcal{M}_{ref} , but we also have the entire dataset that was used to train \mathcal{M}_{ref} .

For this setup we investigate if the smaller base LMs we make can **match** the performance of the larger reference LM. Note that this is in contrast to what was observed in the previous (different) setup of Section 3, where the smaller models had significantly worse performance than their corresponding reference models.

4.1 Setup

For the experiments in this section, due to limited compute resources, we used less than billion size reference models GPT2-large, and GPT2-medium; and the 9B token OpenWebText as our dataset. Table 1 provides details on these models. For the description of the experiment, we will just focus on GPT-2 large; the same process was followed for GPT-2 Medium.

1. We first split this dataset into a training set \mathcal{D}_{train} and a smaller validation subset \mathcal{D}_{val} .
2. We then train the full 36-layer GPT-2 large model on the entire training set \mathcal{D}_{train} for certain number of training steps (T) and evaluate its validation loss on \mathcal{D}_{val} which we call benchmark val loss. Here by validation loss we mean the standard log-perplexity (refer Table 6).

3. We then apply Inheritune v2 (refer Algorithm 2) with the following inputs: the reference model given as input is the (now trained) 36-layer model, the smaller model’s intended number of layers was $n = 18$, the dataset is the full \mathcal{D}_{train} , the benchmark val loss and the number of steps of training was the same as that used in step (2) above for the reference model. This gives us our 18-layer model; we evaluate the validation loss of this model.
4. If the validation loss of the our model in step (3) is worse than that of the reference model in step (2), we re-do step (3) but now with the value of n increased by 2. That is, we progressively explore bigger sizes for the inherited model until we achieve parity in validation loss with the original reference model.
5. Separately and for comparison, we also trained a randomly initialized but other wise identical 18-layer model for the same number of epochs as the models in steps (2) and (3).

Algorithm 2 Inheritune v2

Require: Reference model \mathcal{M}_{ref} , split \mathcal{D} in \mathcal{D}_{train} and \mathcal{D}_{val} , benchmark val loss, number of steps T , and training hyper-parameters.

- 1: Let k be the number of transformer layers in \mathcal{M}_{ref} , with the corresponding parameters being $\theta_0, \theta_1, \dots, \theta_{k-1}$.
 - 2: **Initialize** target LM θ_{tgt} with the first $n = k/2$ layers $\theta_0, \theta_1, \dots, \theta_{n-1}$, as well as the token embedding and prediction head layers of \mathcal{M}_{ref} .
 - 3: **Train** the target LM from this initialization for T steps using the same hyper-parameters.
 - 4: **Grow** θ_{tgt} and retrain (step 3) until θ_{tgt} matches the benchmark val loss.
-

4.2 Results and Discussions

Table 6 illustrates that a smaller target LM with a few layers can be extracted without losing performance in the validation loss from the GPT2 family. Moreover, we also show that the target LM crafted using a larger reference LM through Inheritune demonstrates comparable zero-shot performance on two relevant downstream tasks. Next, our results show that the target LM developed using Inheritune surpasses a similar-sized model in performance, particularly when the target LM is

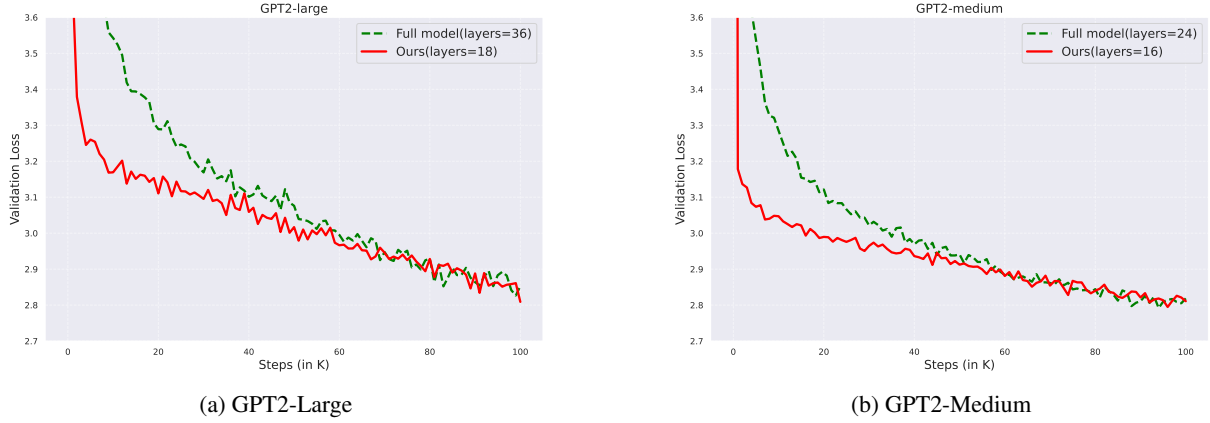


Figure 5: Pre-training GPT2-large and GPT2-medium models with OpenWebText-9B using vanilla recipe and our training recipe for 100K steps. Models derived using Inheritune despite being smaller in size converge faster while matching the final validation loss of the full sized models. See also Table 6

| Models | Layers | Initialization | Steps | Pre-train | Downstream (0-shot) | |
|---------------------------|--------|----------------|-------|---------------------------|---------------------------|---------|
| | | | | Val loss (\downarrow) | Wikitext (\downarrow) | Lambada |
| GPT-2 Large | 36 | rand init | 100K | 2.85 | 34.84 | 34.14 |
| | 18 | rand init | 100K | 2.97 | 37.63 | 30.97 |
| | 18 | rand init | 200K | 2.84 | — | — |
| | 18 | Ours | 100K | 2.80 | 35.38 | 34.64 |
| GPT-2 Medium | 24 | rand init | 100K | 2.81 | 31.93 | 36.54 |
| | 16 | rand init | 100K | 2.86 | 33.67 | 34.60 |
| | 16 | rand init | 200K | 2.83 | — | — |
| | 12 | Ours | 100K | 2.87 | — | — |
| Final Model \rightarrow | 14 | Ours | 100K | 2.84 | — | — |
| | 16 | Ours | 100K | 2.81 | 32.04 | 35.96 |

Table 6: Pre-training and downstream performance of GPT-2 medium and large LLMs evaluated using val loss, Wikitext, and Lambada Openai downstream tasks. Small LMs derived using our method perform comparably to their full-sized counterparts for GPT-2 large and GPT-2 medium. The model initialized with our method performs better than models with random initialization. For a GPT-2 Large model, we achieve the benchmark val loss in the first round of our method. For the GPT-2 medium model, we trained the model with 3 rounds of our method to achieve the benchmark val loss. We evaluate the downstream performance only for the final models.

| Model | Layers | Initialization | Steps | Pre-train Val loss (\downarrow) |
|--------------|--------|-----------------------------|-------|-------------------------------------|
| GPT-2 Medium | 24 | rand init | 100K | 2.81 |
| | 16 | rand init | 100K | 2.86 |
| | 16 | rand init | 200K | 2.83 |
| | 16 | attn+mlp (w/o layernorm) | 100K | 2.80 |
| | 16 | attn (w/ layernorm) | 100K | 2.84 |
| | 16 | mlp (w/ layernorm) | 100K | 2.85 |
| | 16 | Ours (w/o res. connections) | 42K | 7.5 |
| | 16 | Ours | 100K | 2.81 |

Table 7: Final pre-training performance of GPT2-medium (355M) models while initializing various sub modules in a transformer block. Here we first fixed the number of layers (to 16) to achieve the same val loss as the full model when trained with OpenWebText-9B dataset and then analyze the impact of initialization using an existing 24 layer GPT2-medium model. We observe that Inheritune which initializes attention (key, query, value, projection), mlp and layernorm weights has very similar gains compared to just initializing attention and mlp weights. There is little to no impact of layernorm initialization.

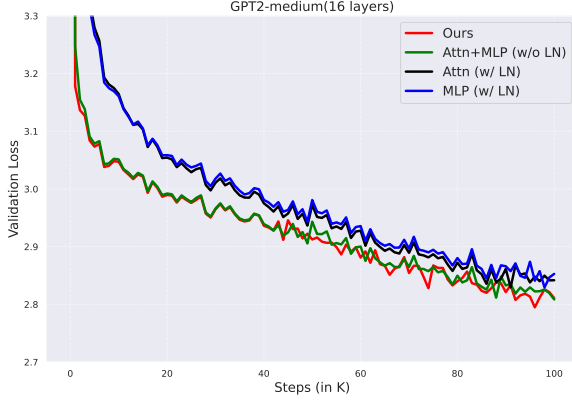


Figure 6: **Inheritune outperforms various other initialization baselines.** We analyze Inheritune approach while initializing some specific weights of a transformer blocks. Here we initialize each transformer block of a 16 layer GPT2 medium model in three different configurations. First we initialize the attention and mlp weights with randomly initializing the layernorms. Next we initialize only the attention and the mlp weights with all the respective layernorms.

trained from scratch. This result also aligns with our results in limited data setting (refer Section 3) where we show that small base LM derived with Inheritune matches the downstream performance of similar size models trained from scratch on large datasets. Interestingly our target model trained with Inheritune for 100K steps outperforms a same size model pre-trained from scratch for double the number of training steps i.e. 200K steps. Through the perspective of convergence, some prior works have made some connections of over-parameterization with faster convergence (Bengio et al., 2005; Vaswani et al., 2018). In Figure 5 we show that the LMs derived with Inheritune although smaller compared to their reference models still converge as fast as their large size reference models.

Next we ran extensive experiments with GPT2-medium (355M) to better understand the impact of initialization on transformer blocks on both generalization (val loss) and convergence speed during the pre-training phase. It is important to note that GPT2 models has a parameterized layernorm (LN). Previously in Table 6 we observed that our 16 layer GPT2-medium model with Inheritune matches the val loss of the 24 layer GPT2-medium. We include the 24 and 16 layer GPT2-medium model from scratch and also our 16 layer model. Next we trained 4 new GPT2-medium 16 layer variant model while initializing different parts of the trans-

former blocks with weights from 24 layer GPT2-medium trained from scratch. We initialize the transformer blocks with 1) attention ((key, query, value and projection) and the layernorm weights (attn w/ layernorm), 2) attention and mlp weights without the layernorm (attn+mlp w/o layernorm), 3) mlp weights with the layernorm (mlp w/ layernorm), and 4) Inheritune without residual connections. We emphasize that Inheritune performs initialization which inherits weights of attention and mlp weights with the layernorm (attn+mlp w/ layernorm). In Table 7 we observe that models trained with attention and mlp weights has shown the best performance agnostic to the layer norm initialization. The convergence patterns of these models can be seen in Figure 6. Therefore there is a clear advantage of initialization of both attention and mlp weights. Next we also make a surprising observation that initializing either the attention or mlp weights resulted in very similar improvements in both the speed of convergence and the final validation loss.

5 Related Works

Conventional pre-training of small base LMs.

Small base LMs are typically trained from scratch with trillion size data which takes a huge compute budget and time. For instance, tinyllama-1B (Peiyuan Zhang and Lu, 2023) a scaled down version of LLaMA2 models (Touvron et al., 2023b) is independently pre-trained from scratch with 3T tokens for 90 days utilizing 16 GPUs. Efforts have also been made to create small base LMs using high-quality synthetic data, as demonstrated in the training of the Phi models (Jawaheripi et al., 2023). Despite the use of high-quality synthetic data, these models still require training with trillions of tokens. In this paper we argue that small base LMs can be trained in a data efficient way (also faster) if trained with better initialization and with multiple passes over the data. Furthermore, our pre-training dataset consists of 1B tokens randomly sampled from redpajama, without any additional efforts in data curation.

Model compression. In settings outside of LLMs, such as in small neural networks (NNs) and convolutional neural networks (CNNs), pruning weights—especially from the width of pre-trained networks post training has proven to be highly effective (Han et al., 2015; Li et al., 2016) in reducing memory usage through sparsity. These models can

be pruned by up to approximately 90% (LeCun et al., 1989; Frankle and Carbin, 2019) without compromising their test performance. However, to our knowledge, there is no existing works in LLM scale that showed a similar result. This is likely because the functions these large base LMs aim to learn during pre-training are significantly more complex than those previously demonstrated in CNNs and NNs. Additionally we have provided a detailed discussion of Inheritune with some other pruning methodologies in Section 3.2.

Connection with knowledge distillation. Some prior knowledge distillation (KD) works have initialized the student network with few layers from the teacher network and then distilled the student using full pre-training data. It is important to note that these works (Turc et al., 2020; Sanh et al., 2019) employed this initialization has a trick embedded in their training details and not a main focus of their work. These studies typically explored the impact of KD in scenarios where the model was fine-tuned for a specific task, also they did not analyze how much of the improvement came from the initialization technique as opposed to the distillation process itself. Furthermore these works also used the full pre-train data used to train the teacher for distilling the small student model. In our work we provide Inheritune as algorithm and deeply analyzed it in a low data regime and full data regime. Next we also showed improvements in multi task pre-train setup. Finally KD is extremely compute intensive and Inheritune is very efficient in compute as we ground-up train a small base LM using 1 GPU for less than 12 hours.

Initialization. Some recent works (Xu et al., 2024) in transformers that uses weight initialization for faster fine-tuning for vision setting. However there is an evidence (Trockman and Kolter, 2023) that shows attention patterns of vision and language pre-trained models are different and hence initialization method that works for vision setting may not necessarily works in language setting as well. Next another recent work (Liu et al., 2023) achieves better domain specific performance using weights from an existing model. However this work have also used high quality synthetic data for finetuning. Our paper makes no claims of high quality data rather we sample a small random subset of an pre-training data used to train an existing Large base LM.

6 Implications

In this section we discuss some of the key implications of our work.

Cheap and easy development of small base LMs. Pre-training a small base LM of 1-2B parameters from scratch is extremely expensive. For instance mpt-1.3B base LM is pre-trained with 440 A100 GPUs for half a day, while the Pythia-1.4B base LM (Biderman et al., 2023) utilized 64 A100-40GB GPUs for 4830 GPU hours. Similarly, TinyLLaMA-1.1B model (Peiyuan Zhang and Lu, 2023) was pre-trained using 16 A100 GPUs for 3 months. Our 1.5B (1B data variant) LM shows competitive performance despite being trained with 1 A6000 GPU for less than 12 hours. Typically small base LMs are finetuned for a specific task before deployment and are not used in it’s base form. With Inheritune we present a really easy and cheap way for developing a small base LM to be later finetuned before deployment.

Naive baseline for pre-training a scaled down variant of large base LMs. Typically small variants of large base LMs are pre-trained using the same pre-training data (Peiyuan Zhang and Lu, 2023; Groeneveld et al., 2024). Our recipe introduces a new perspective of identifying sufficient depth without losing any generalization on the held out set (pre-train val loss). Next we also show that even with a small fraction of pre-train data (randomly sampled) and few initial layers of the large base LM one can develop a small base LM. Therefore our Inheritune recipe has the potential to become the naive baseline for any pre-training pipeline aiming to develop a smaller variant of a large base LM.

Sufficient Depth for Bigger LLMs. The architectural choices particularly the total number of layers used for an LLM is largely taken arbitrarily. In Section 4 we present Inheritune v2 (refer Algorithm 2) which can be utilized to identify sufficient depth of a particular model without loosing on pre-train validation loss. The usage of our recipe can go beyond depth.

7 Conclusion and Limitations

Limitation and Future Works. The Inheritune recipe has some limitations. The key limitation of our recipe is that it can only change the number of transformer blocks/layers and that restricts the

architectural design space of selecting the hidden size and attention head for a custom small base LM. Our method is likely to be very sensitive to the quality of the training dataset, since it is so small. Finally, the selection of which blocks to keep and how to select the training dataset and hyper-parameters have not been explored in our work but could offer some benefits.

Conclusion. In conclusion, this study demonstrates the efficacy of our Inheritune method for pre-training small base language models using significantly fewer data tokens and computational resources than traditional methods. We proposed a very simple way of creating a smaller LM using a large reference model and a small training set. We also analyzed a slightly different setting where Inheritune is employed with the full pre-training dataset. In that case, the validation loss achieved by GPT2-large with 100K training steps on OpenWebText 9B tokens can be matched using our GPT2 variant with half the layers. Furthermore, we presented similar results utilizing GPT2-medium with 24 layers and our variant of GPT2-medium with 16 layers. Our proposed recipe provides a new perspective on simple methods to make small base LMs from an existing large pre-trained model and a small subset of available pre-training data.

References

- Lightning AI. 2023. Lit-gpt. <https://github.com/Lightning-AI/lit-gpt>.
- Author(s). 2023. MPT-1.3B Model: A Large-Scale Language Model. <https://huggingface.co/mosaicml/mpt-1b-redpajama-200b>. Accessed: 2023-02-14.
- Yoshua Bengio, Nicolas Roux, Pascal Vincent, Olivier Delalleau, and Patrice Marcotte. 2005. *Convex neural networks*. In *Advances in Neural Information Processing Systems*, volume 18. MIT Press.
- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. Pythia: A suite for analyzing large language models across training and scaling. *arXiv preprint arXiv: 2304.01373*.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. *PIQA: reasoning about physical commonsense in natural language*. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7432–7439. AAAI Press.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. *A large annotated corpus for learning natural language inference*. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. *Language models are few-shot learners*. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. *Boolq: Exploring the surprising difficulty of natural yes/no questions*. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 2924–2936. Association for Computational Linguistics.
- Together Computer. 2023. *Redpajama-data: An open source recipe to reproduce llama training dataset*.
- Jonathan Frankle and Michael Carbin. 2019. *The lottery ticket hypothesis: Finding sparse, trainable neural networks*. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. *A framework for few-shot language model evaluation*.
- Jonas Geiping and T. Goldstein. 2022. *Cramming: Training a language model on a single gpu in one day*. *International Conference on Machine Learning*.

- Xinyang Geng and Hao Liu. 2023. [Openllama: An open reproduction of llama](#).
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. 2024. Olmo: Accelerating the science of language models. *arXiv preprint arXiv: 2402.00838*.
- Song Han, Jeff Pool, John Tran, and William J. Dally. 2015. Learning both weights and connections for efficient neural networks. *NEURIPS*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, D. Song, and J. Steinhardt. 2020. Measuring massive multitask language understanding. *International Conference On Learning Representations*.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *International Conference On Learning Representations*.
- Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, Jyoti Aneja, Sebastien Bubeck, Caio César Teodoro Mendes, Weizhu Chen, Allie DelGiorno, Ronen Eldan, Sivakanth Gopi, Suriya Gunasekar, Mojan Javaheripi, Piero Kauffmann, Yin Tat Lee, Yanzhi Li, Anh Nguyen, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Michael Santacrose, Harkirat Singh Behl, Adam Taumann Kalai, Xin Wang, Rachel Ward, Philipp Witte, Cyril Zhang, and Yi Zhang. 2023. Phi-2: The surprising power of small language models. <https://www.microsoft.com/en-us/research/blog/phi-2-the-surprising-power-of-small-language-models>.
- Tian Jin, Michael Carbin, Daniel M. Roy, Jonathan Franklin, and Gintare Karolina Dziugaite. 2022. [Pruning’s effect on generalization through the lens of training and regularization](#). In *Advances in Neural Information Processing Systems*.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Vid Kocijan, Thomas Lukasiewicz, Ernest Davis, Gary Marcus, and Leora Morgenstern. 2020. A review of winograd schema challenge datasets and approaches. *arXiv preprint arXiv: 2004.13831*.
- Yann LeCun, John Denker, and Sara Solla. 1989. [Optimal brain damage](#). In *Advances in Neural Information Processing Systems*, volume 2. Morgan-Kaufmann.
- Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. 2016. [Pruning filters for efficient convnets](#). *ArXiv*, abs/1608.08710.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie DelGiorno, Suriya Gunasekar, and Yin Tat Lee. 2023. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv: 2309.05463*.
- Zhuohan Li, Eric Wallace, Sheng Shen, Kevin Lin, K. Keutzer, D. Klein, and Joseph Gonzalez. 2020. Train large, then compress: Rethinking model size for efficient training and inference of transformers. *International Conference On Machine Learning*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. *ACL*.
- Bingbin Liu, Sebastien Bubeck, Ronen Eldan, Janardhan Kulkarni, Yanzhi Li, Anh Nguyen, Rachel Ward, and Yi Zhang. 2023. Tinygsm: achieving >80 *arXiv preprint arXiv: 2312.09241*.
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning. *arXiv preprint arXiv: 2007.08124*.
- Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D. Lee, Danqi Chen, and Sanjeev Arora. 2023. Fine-tuning language models with just forward passes. *NEURIPS*.
- Niklas Muennighoff, Alexander M. Rush, Boaz Barak, Teven Le Scao, Aleksandra Piktus, Nouamane Tazi, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. 2023. Scaling data-constrained language models. *NEURIPS*.
- Tianduo Wang Peiyuan Zhang, Guangtao Zeng and Wei Jia. 2023. [Tinyllama](#).
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. [Winogrande: An adversarial winograd schema challenge at scale](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New*

- York, NY, USA, February 7-12, 2020, pages 8732–8740. AAAI Press.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *NEURIPS*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models. *ARXIV*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv: 2307.09288*.
- Asher Trockman and J. Z. Kolter. 2023. [Mimetic initialization of self-attention layers](#). *International Conference on Machine Learning*.
- Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2020. Well-read students learn better: On the importance of pre-training compact models. *ICLR*.
- Sharan Vaswani, F. Bach, and Mark W. Schmidt. 2018. Fast and faster convergence of sgd for over-parameterized models and an accelerated perceptron. *International Conference on Artificial Intelligence and Statistics*.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2018. [Glue: A multi-task benchmark and analysis platform for natural language understanding](#). *BLACK-BOXNLP@EMNLP*.
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. 2023. Sheared llama: Accelerating language model pre-training via structured pruning. *arXiv preprint arXiv: 2310.06694*.
- Mengzhou Xia, Zexuan Zhong, and Danqi Chen. 2022. [Structured pruning learns compact and accurate models](#). *Annual Meeting of the Association for Computational Linguistics*.
- Zhiqiu Xu, Yanjie Chen, Kirill Vishniakov, Yida Yin, Zhiqiang Shen, Trevor Darrell, Lingjie Liu, and Zhuang Liu. 2024. [Weight selection for model initialization](#). In *The Twelfth International Conference on Learning Representations*.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. Opt: Open pre-trained transformer language models. *ARXIV.ORG*.

Supplementary Materials

Contents

- [A](#): Frequently Asked Questions.
- [B](#): Implementation Details

A Frequently Asked Questions

A.1 Why Inheritune focuses on small base LMs pre-trained using redpajama or similar datasets and not with other publicly available models such as Phi?

Phi-1.5 is trained on GPT-4 generated high quality synthetic data which neither full or even a small fraction is not made publicly available. Since data quality is paramount to LLM pre-training we restricted ourselves to model trained with redpajama-1T for fairness. We showed pythia because the data distribution of pile and redpajama is similar. The key idea we demonstrate here is that one can make a performant small base model utilizing an existing model and it's a small subset of pre-training data.

A.2 Inheritune has some resemblance with structural pruning why are you not comparing it with such pruning techniques?

We have compared Inheritune with targeted structural pruning recipe called Shearing (Xia et al., 2023) and also provided a detailed discussion in Section 3.2 about why it is bit unfair to directly compare any pruning methodology to Inheritune which we believe can be best described as a model initialization technique as we ground up pre-train a model. Let's understand this with an example, let us assume that we want a scaled down or smaller version of OpenLLaMA-7B model and as an academic researcher we can only manage 1 GPU with 20GB memory. The OpenLLaMA-7B model takes 17 GB of memory to load on that particular GPU. Now pruning is virtually impossible as it requires a lot of compute and memory (for iterations of pruning steps). In such setting Inheritune can be used to pre-train an scaled down version with very little compute requirements. Given one assumes a lot of available compute with access to the full pre-train data Inheritune may also help in determining the sufficient depth of a model as shown in Section 4.

A.3 I am not sure what is the usefulness of this method as it does not seem to be parameter efficient, not compute efficient and its benefits against existing pruning methods are yet to be explored.

We show that with a little training compute 1 GPU for less than 12 hours (half the time assumed by Geiping and Goldstein (2022) for BERT pre-training) one can develop a performant 1.5B base LM using just 1B data, which compares favorably across widely used pre-training datasets/benchmark(s). Next we show that less deeper base LMs can be pre-trained while using Inheritune (with upto 50% parameter reduction for GPT2-large model) without compromising the pre-train val loss. Some key implications of our work are discussed in Section 6.

Although we don't make any specific claim that Inheritune is parameter efficient. Some prior parameter efficient methods (Hu et al., 2021; Malladi et al., 2023) are designed to reduce the memory requirement during training. Unlike these prior works Inheritune reduces both trainable and total model parameters and thus makes training more parameter efficient. We defer a through investigation of Inheritune for parameter efficient adaptation to future.

B Implementation Details

Training details of Small base LMs trained with 1B data. We present our main results with Our-1.5B model trained with an existing OpenLLaMA version 1 (Geng and Liu, 2023) and 1 B tokens randomly sampled from 1T redpajama version1 data. The hyper-parameters related to this model is provided in Table 8. It is important to note that our claim that we only use 1 GPU for less than 12 hours to train Our-1.5 B model is specific to models derived using Inheritune with 1B data. Next we also train multiple sub-models as shown in Figure 2 the training details remains consistent with that of the initial model discussed earlier. However we observe that increasing the number of layers in a sub-model also increase the training time.

Training details of small base LMs with 50B data. We also trained our 1.5B model with larger subsets of data as shown in Figure 3. It is important to note that all the intermediate tokens until 50B are intermediate checkpoints of a single training run. Some of the key hyper-parameters of our training runs are discussed in Table 9. We have also

trained three variants of small base LMs utilizing 3 different reference base LMs namely OpenLLaMA-3B, OpenLLaMA-7B and LLaMA2-7B. For target LMs developed with OpenLLaMA-3B we use $n=13$ i.e. 13 layers. For target LMs developed using reference LMs of 7B parameters we use $n=7$ i.e. 7 layers. The training hyper-parameters remains consistent across all the models trained with 50B subset of the pre-train data.

Training details of GPT2 models. For an exploratory analysis we also trained GPT2-medium and GPT2-large models as shown in Table 6. The GPT2-medium model with 24 layers and also the inheritune variant with lesser layers are all trained with a 50K tokens per batch. The GPT2-large model and it’s variants developed using Inheritune are trained with 16K tokens per batch. The training details are given in Table 10. The GPT2 models are trained on a single node consisting of 3 A100 with 40 GB memory.

| Hyper-parameters | Value |
|------------------|--------------------|
| Training tokens | 1B |
| Training epochs | 8 |
| Training steps | 64K |
| Learning rate | 3×10^{-4} |
| Scheduler | Cosine |
| Weight decay | 0.1 |
| Optimizer | AdamW |
| Warm up | 1000 |
| Batch size | 131K |
| GPU count | 1 |
| GPU type | A6000 |
| GPU hours | ~8 hours |
| GPU hours/epoch | ~54 mins |

Table 8: Training hyper-parameters of our target 1.5B base LM using OpenLLaMA-3B as refernce LM.

| Hyper-parameters | Value |
|------------------|--------------------|
| Training tokens | 50B |
| Training epochs | ~1 |
| Training steps | 191K |
| Learning rate | 3×10^{-4} |
| Scheduler | Cosine |
| Weight decay | 0.1 |
| Optimizer | AdamW |
| Warm up | 1000 |
| Batch size | 131K |
| GPU count | 1 |
| GPU type | A100 |
| GPU hours | ~18 hours |

Table 9: Training hyper-parameters of our target 1.5B and 1.6B small base LMs developed with Inheritune as shown in Figure 3 using OpenLLaMA-3B, OpenLLaMA-7B and LLaMA2-7B as reference base LMs.

| Hyper-parameters | Value |
|------------------|--------------------|
| Data | OpenWebText |
| Training tokens | 9B |
| Learning rate | 3×10^{-4} |
| Scheduler | Cosine |
| Weight decay | 0.1 |
| Optimizer | AdamW |
| Warm up | 1000 |
| GPU count | 3 |
| GPU type | A100 |

Table 10: Training hyper-parameters of our target GPT2 models developed with Inheritune as shown in Table 3 using OpenLLaMA-3B, OpenLLaMA-7B and LLaMA2-7B as reference base LMs.