

# Honey, I Shrunk the Language: Language Model Behavior at Reduced Scale

Vijeta Deshpande<sup>1</sup>, Dan Pechi<sup>2</sup>, Shree Thatte<sup>1</sup>, Vladislav Lialin<sup>1</sup>, Anna Rumshisky<sup>1,3</sup>

<sup>1</sup>University of Massachusetts Lowell, Computer Science Department

<sup>2</sup>New York University, Center for Data Science

<sup>3</sup>Amazon Alexa AI

{vijeta\_deshpande, shree\_thatte}@student.uml.edu, danpechi@nyu.edu

{vlialin, arum}@cs.uml.edu

## Abstract

In recent years, language models have drastically grown in size, and the abilities of these models have been shown to improve with scale. The majority of recent scaling laws studies focused on high-compute high-parameter count settings, leaving the question of when these abilities begin to emerge largely unanswered. In this paper, we investigate whether the effects of pre-training can be observed when the problem size is reduced, modeling a smaller, reduced-vocabulary language. We show the benefits of pre-training with masked language modeling (MLM) objective in models as small as 1.25M parameters, and establish a strong correlation between pre-training perplexity and downstream performance (GLUE benchmark). We examine downscaling effects, extending scaling laws to models as small as 1M parameters. At this scale, we observe a break of the power law for compute-optimal models and show that the MLM loss does not scale smoothly with compute-cost (FLOPs) below  $2.2 \times 10^{15}$  FLOPs. We also find that adding layers does not always benefit downstream performance.<sup>1</sup>

## 1 Introduction

In the past few years, large language models (LLMs) have grown ever larger (Brown et al., 2020; Shoybi et al., 2019; Chowdhery et al., 2022; Fedus et al., 2022), and the emergent abilities of these models improve with scale. While several studies have looked at the relationship between model size, the amount of training, and performance for LLMs (Kaplan et al., 2020; Hoffmann et al., 2022), the main focus has been on scaling laws for high-compute settings. Very few studies have considered the effects of pre-training at a smaller scale (Turc et al., 2019; Huebner et al., 2021). Thus, the

<sup>1</sup>Our filtered pre-training data, reduced English vocabulary, and code are available at [https://github.com/text-machine-lab/mini\\_bert](https://github.com/text-machine-lab/mini_bert).

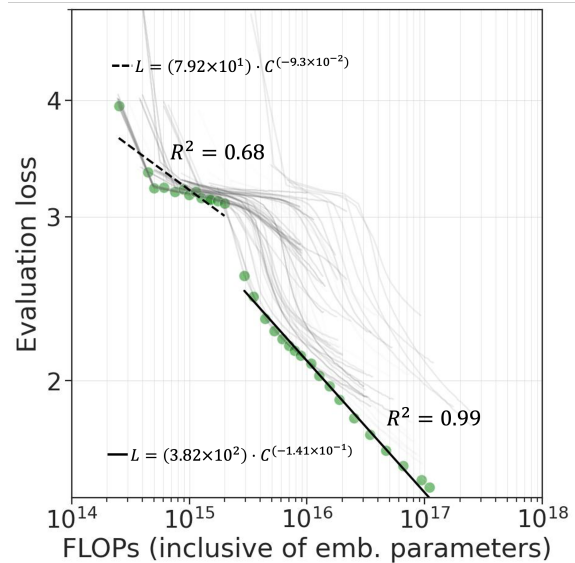


Figure 1: **Power law breaks in low FLOPs region.** Faded gray lines are unique model configurations. The green dots are compute-optimal instances i.e., the model instance with a minimum value of MLM loss for a given range of FLOPs. We observe that there exists a discontinuity in the power curve around  $2.2 \times 10^{15}$  FLOPs

question of when exactly model abilities begin to emerge remains largely unanswered.

In this study, we were interested in understanding whether the emergent phenomena can be observed at a drastically reduced scale, and what the relationship is between upstream and downstream performance at this scale. We also wanted to examine model shapes, configurations, and other factors that might affect whether we see the benefits of pre-training when downscaling a model.

Smaller models have been shown to do poorly when trained even on large volumes of data (Turc et al., 2019), which makes studying downscaling non-trivial. However, during language acquisition, humans are exposed to a reduced-size language before gradually expanding their vocabulary, yet they become fluent even when their vocabulary is limited. Taking our cue from humans, we explore

the hypothesis that reducing language size might allow us to observe the effects of pre-training in small models.

There has been one previous attempt to reduce language size (Huebner et al., 2021), but it was quite limited: one reduced-size Transformer encoder was trained with a non-standard version of masked language modeling (MLM) loss on a relatively small corpus of child-directed speech and evaluated for its ability to pass linguistic tests from a custom grammar test suite.

We use a vocabulary of 21,000 words derived from AO-CHILDES (Huebner and Willits, 2021), a corpus of child-directed speech, to create a filtered corpus containing a subset of standard pre-training corpora: C4 (Raffel et al., 2020), Wikipedia, Book Corpus (Zhu et al., 2015), and others. We pre-train over 70 Transformer encoder models in the 1-100M parameter range, varying model shape, and configuration, including the number of layers, hidden size, the number of attention heads, and the feed-forward layer dimension (intermediate size). We fine-tune and evaluate a series of checkpoints at different FLOPs count on the subset of GLUE filtered with the same vocabulary.

We present evidence that for a realistically downscaled language, the benefits of pre-training are observable even in smaller models with as few as 1.25M parameters. Our results also indicate that models with fewer layers achieve better performance on GLUE tasks. In contrast to Tay et al. (2022), we find a strong correlation between upstream and downstream performance (here, model perplexity and GLUE score). However, pre-training compute optimality does not appear to be crucial for downstream results. We also show that for compute-optimal models at this scale, parameter count does not reliably predict MLM loss, suggesting a limitation to scaling laws. We observe a departure from the FLOPs-Perplexity law, characterized by a sudden shift in the exponent value in the low-compute region of FLOPs  $\leq 2.2 \times 10^{15}$  (cf. Figure 1). This represents a divergence from previous observations regarding scaling laws (Kaplan et al., 2020; Hoffmann et al., 2022).

## 2 Related Work

**Scaling Laws** Kaplan et al. (2020) has demonstrated and popularized power law dependency between language model parameter count and per-

plexity. This further motivated the existing trend for increasing model size (Brown et al., 2020; Chowdhery et al., 2022; Fedus et al., 2022). Investigation of smaller models has mostly focused on distillation (Sanh et al., 2019; Turc et al., 2019) and achieving the best performance given the parameter count. In contrast, we are interested in understanding at what scale pre-training works and the emergence of language model abilities to help downstream tasks.

Changing pre-training data size via training token volume (Pérez-Mayos et al., 2021; Zhang et al., 2020), vocabulary size (Gowda and May, 2020), or the number of epochs (Voloshina et al., 2022) has also been explored for effects on language acquisition. These studies have generally demonstrated low-level linguistic tasks involving syntax require relatively little data volume compared to more complex tasks like WinoGrande (Sakaguchi et al., 2019). The relationship between model size, input data, and downstream performance remains the subject of much debate as to the nature of scaling laws. Hoffmann et al. (2022) concluded data size and model size should be scaled equally to achieve compute-optimal training for LLM’s. Further credence to reducing LLM compute is lent by Sorscher et al. (2022), who found input data pruning can improve models to scale exponentially, and Chan et al. (2022), who showed LLM in-context learning derives from the Zipfian distribution of pre-training data.

Most relevant to this work is Huebner et al. (2021) who found a small language model trained on child-directed speech can achieve comparable performance to larger LMs on a set of probing tasks. In contrast to them, we train multiple models, explore scaling in the low-compute region and evaluate on a filtered version of GLUE instead of a set of linguistic tests.

## 3 Methodology

This section discusses the language simplification process, pre-training data, development of the data tokenizer, language model configuration, and pre-training objective in detail.

### 3.1 Simplifying language

To create a corpus of reduced English, we filter large text corpora based on a word vocabulary from AO-CHILDES (Huebner and Willits, 2021). The AO-CHILDES corpus contains English transcripts

Corpus name	Sentences (mil.)	Tokens (mil.)
C4 <sup>3</sup>	3	427
C4	27	428
Book Corpus	12	190
Wikipedia	4.8	76
Simplified Wikipedia	0.19	3
Children’s Book Test	0.08	1
Total	47.07	1,125

Table 1: Size of the filtered pre-training data used in this study.

of child-directed speech. With the transcripts, we generate a vocabulary by removing special characters and tokenizing words by spaces. We also remove gibberish words present in the transcripts e.g. “bababa”. With this process, we construct a set of 21 thousand unique words.

### 3.2 Pre-training data

We filter data from five text corpora: Wikipedia, Simple English Wikipedia<sup>2</sup>, Book Corpus (Zhu et al., 2015), Children’s Book Test (CBT) (Hill et al., 2015), and Common Crawl (C4) (Raffel et al., 2020), to obtain pre-training data. We filter C4 two ways: span level (110 words span size, 30 words stride) and sentence level. We select a text span (or sentence) to include in the pre-training data if and only if there are no words out of the vocabulary of interest, ignoring any numeric characters. For sentence-level filtration, we process text data on all five corpora. With sentence-level filtration, we collect approximately 44 million sentences which we concatenate to construct six million spans. The combination of both span- and sentence-level data filtration provided us with over nine million pre-training sequences of an average length of 127 BPE tokens. Finally, we split the filtered data into three sets: train, development, and test, of sizes nine million, 100 thousand, and 100 thousand, respectively. We provide the amount of filtered data from each text corpus in Table 1.

For the rest of the paper, we use the word “vocabulary” to refer to the number of unique tokens instead of unique whitespace-separated words, unless otherwise mentioned.

### 3.3 Tokenizer

Since we are working with a reduced language, commonly used subword vocabulary sizes for En-

glish might be suboptimal. We want to achieve a reasonable balance between over-splitting (splitting the words into units smaller than the smallest meaningful morphological units, e.g., splitting into characters) and under-splitting (e.g., retaining full words instead of splitting them into meaningful subwords).

We conducted a series of experiments in which we tuned the vocabulary size to find the right balance. While varying the vocabulary size, we track two metrics for the tokenized text: *word-split ratio* and another metric we define, the *exact sub-token matching score* (ESMS). Word-split ratio is the number of tokens into which a word is split, where words are separated by whitespace. For example, if the word “cooking” is converted to “cook” and “ing”, then the word-split ratio value is two. We measure and report the word-split ratio value for 5,000 examples sampled from the set of collected pre-training data without replacement.

To measure ESMS, we compare the tokenizer performance with morpheme-based subword tokens. For example, in case of the word “cooking”, we check whether the tokenizer is splitting the word into two tokens, ‘cook’ and ‘ing’. For this purpose, we used a manually-curated list of 127 words with their corresponding morpheme-based sub-tokens, (see Table 5 in the Appendix for some examples). ESMS is computed as an exact match to the reference tokenization. For one example, it is equal to 1 if the word is tokenized exactly as in the reference and 0 in any other case.

We experiment with three types of tokenizers, Byte-Pair Encoding (BPE) (Radford et al., 2019), WordPiece (Devlin et al., 2018), and SentencePiece (Raffel et al., 2020). Similar to the study conducted by FitzGerald et al. (2022), we select vocabulary size for each type of tokenizer by minimizing the absolute difference of word-split ratio compared to the reference tokenizer. We consider separate reference tokenizers for each tokenizer type. For BPE, WordPiece, and SentencePiece, we select pre-trained tokenizers published by (Liu et al., 2019), (Devlin et al., 2018), and (Raffel et al., 2020), respectively, as our reference tokenizers. After selecting the vocabulary size for each tokenizer type, we select the tokenizer with the highest value of ESMS as our final choice.

With the above-mentioned selection process, we find that the BPE tokenizer with a vocabulary size of 19,000 and ESMS of 0.2604, is the best-suited

<sup>2</sup><https://simple.wikipedia.org>

<sup>3</sup>Span-level filtering.

tokenizer for our study. We provide the results of our tokenizer selection experiments in Appendix B.

### 3.4 Model architecture and configuration

The models we pre-train in our experiments closely follow the configuration setting of RoBERTa (Liu et al., 2019). We scale down and vary the model’s hidden size, intermediate size (FFN hidden dimension size), number of hidden layers, and number of attention heads such that the total number of trainable parameters does not exceed 20 million. To separately control model hidden size and embedding size, we also add a linear layer, followed by a normalization layer (Ba et al., 2016) between the embedding block and the first Transformer layer.

### 3.5 Pre-training objective

In our study, we pre-train models on a Masked Language Modeling (MLM) task (Devlin et al., 2018). We chose MLM instead of regular (causal) language modeling, because of its effectiveness for natural language understanding tasks at a smaller scale as demonstrated by BERT.

We conducted an exploratory set of experiments to observe the effect of various MLM objective settings on validation perplexity. We found that using a random word replacement strategy and same-word replacement strategy doesn’t improve the model at a small scale. Hence, to enable considerable learning in the limited parameter setting, we do not use random replacement and same-word replacement of the token selected for masking. In other words, we always replace the token selected for masking with the mask token <mask> before inputting it into the model. Otherwise, we adopt the same strategy as BERT pre-training by masking 15% of tokens.

## 4 Experimental Setup

In our experiments, we explore the relationship between training data size, model size (number of parameters), model shape, cost of training (FLOPs), and performance during pre-training and downstream. In the following subsections, we will discuss our strategy for exploring various model shapes followed by a discussion on hyperparameter settings in detail.

### 4.1 Exploration of model configuration

To investigate the impact of reduced model size, we start by scaling down the base configuration of

RoBERTa (Liu et al., 2019) from its initial hidden size of 768 to 256, and the number of hidden layers and attention heads from 12 to 8. For intermediate layer size, we follow the common practice of setting it to a value four times that of the hidden size. We refer to this configuration as the anchor configuration. We pre-train a model with the anchor configuration and explore three values for embedding size, hidden size, intermediate size, number of hidden layers, and number of attention heads, varying one hyperparameter at a time. With such unidirectional exploration, we pre-train 16 models. We refer to this set of 16 models as **set-1**. To explore more model configurations, we randomly sample 16 configurations that are not included in **set-1**. For random sampling, we only explore values that are powers of two and are upper-bounded by 256, 256, 1024, 8, and 8, for the embedding size, hidden size, intermediate size, number of attention heads, and number of hidden layers, respectively. We refer to this set of 16 models as **set-2**. Furthermore, we pre-train 30 more models by performing unidirectional explorations of hidden size and the number of hidden layer values by anchoring other hyperparameter values. We refer to this set of 30 models as **set-3**.

### 4.2 Pre-training

For every model configuration, we keep the input sequence length fixed at 128 tokens. All models are initialized with a fixed random seed value of zero.

Once initialized, we train the model for one epoch with a batch size of 256 for 35,000 weight updates. We use an inverse square root learning rate scheduler with 5% warmup.

We conducted a few trials guided to decide the peak learning rate value for our experiments. We started with values higher than  $6e-4$ , based on findings published by Liu et al. (2019) and Kaplan et al. (2020), and kept on reducing the learning rate until we observed a stable training loss curve. We observed  $1 \times 10^{-1}$  to be suitable for models with more than 18 million parameters and  $5 \times 10^{-1}$  otherwise.

For optimization, we use the AdamW optimizer (Loshchilov and Hutter, 2017) with  $\beta_1$ ,  $\beta_2$ , and  $\epsilon$  values set to 0.9, 0.95, and  $10^{-8}$ , respectively. After preliminary experiments we set the weight decay parameter to 0.01. Besides the learning rate, we keep all optimizer-related hyperparameters

constant across all model configurations. For all dropout layers in the model, we adopt the same value of 10% as that of the RoBERTa model (Liu et al., 2019).

### 4.3 Fine-tuning

We evaluate pre-trained models on GLUE (Wang et al., 2018). Because our pre-trained data consists of a limited vocabulary, we fine-tune and test on GLUE task datasets with the same vocabulary filtering, in addition to unfiltered variants.

For all tasks, we fine-tune our pre-trained models for 5 epochs and report the performance of the best performance value on the validation set averaged over three seed values. For all fine-tuning experiments, we keep the batch size fixed at 32. Over the five training epochs, we vary the learning rate value with a linear scheduler with a warmup of 5%. We set the peak-learning rate within the range from  $2e-5$  to  $2e-4$  value, according to the task. In addition to these pre-trained models, we fine-tune and evaluate GLUE for randomly-initialized versions of pre-trained models as well.

### 4.4 Evaluation metrics

**Pre-training** For pre-training results, we measure and report the cross-entropy loss and perplexity on the test split of the data. We use the cross-entropy loss and perplexity calculated on the development set for curve fitting. In both cases, we calculate the cross-entropy loss only for the masked tokens, and the perplexity value is calculated by exponentiating the cross-entropy loss value. We also calculate the FLOPs (compute cost) as defined by Hoffmann et al. (2022). We first calculate FLOPs per training sequence based on the model parameters (including the embedding parameters) and multiply it by the amount of training data seen by the model to get a total number of FLOPs. We provide a detailed formula of FLOPs calculation in Appendix A.

**Cost-effectiveness analysis** We use the Incremental Cost-Effectiveness Ratio (ICER) (Bambha and Kim, 2004) to conduct cost-effectiveness analysis of different model configuration hyperparameters. We treat the FLOPs and model perplexity values as proxies for the expenditure and the outcome of expenditure, respectively. Therefore,

$$ICER = \frac{\Delta Outcomes}{\Delta Cost} = \frac{\Delta Perplexity}{\Delta FLOPs} \quad (1)$$

We calculate the difference (the  $\Delta$  values) by comparing a model configuration with the next cheaper option (e.g., we compare the model with a hidden size of  $2^8$  to the model with a hidden size of  $2^7$ ). For the specific case of increasing hidden size from  $2^7$  to  $2^8$ , ICER represents performance gain (reduction in perplexity) per additional FLOPs spent on increasing the hidden size value from  $2^7$  to  $2^8$ . We calculate ICER values for four hyperparameters namely, embedding size, hidden size, intermediate size, and the number of hidden layers.

**Fine-tuning** We use standard metrics for the GLUE benchmark: accuracy, Matthew’s correlation score, and combined correlations score depending on the task. For our conducted experiments, we report the average value of all performance metrics across tasks.

## 5 Results and Discussion

### 5.1 Curve fitting

To assess the empirical relationship between model performance and data size, model size, and FLOP values, we fit a power curve of the form  $y = C \cdot x^e$ , separately, to model size, data size, and FLOP values. We only consider the compute-optimal instances for curve fitting. To find compute optimal instances, we first divide the FLOPs values into over 30 bins and fetch the checkpoint corresponding to the minimum value MLM loss for each FLOPs bin. We use an implementation of the Levenberg-Marquardt (Moré, 1978) algorithm provided under the *SciPy* library for curve fitting.

We observe that the optimal values of the exponents for data size and model size are,  $-0.2459$  and  $-0.2805$ . Note that these values are expected to be different from those in Kaplan et al. (2020), since we work with a different loss and a reduced-vocabulary language. The small difference between both exponent values suggests that the MLM loss reduces with a similar pace for data and model scaling. Hence, in our downscaled problem, we find data and model scaling equally important for compute-optimality. Although, we find  $R^{24}$  values for both curves i.e., loss vs. data size and loss vs. model size, are low (Figures 2 and 3). Moreover, we observe that for FLOPs values greater than  $2^{15}$ , a power curve nearly perfectly predicts the compute optimal MLM loss value for a given compute

<sup>4</sup> $R^2$  is the coefficient of determination and we adopt the default definition of  $R^2$  in the Scipy Python library. Please refer to *Scipy* for more details.

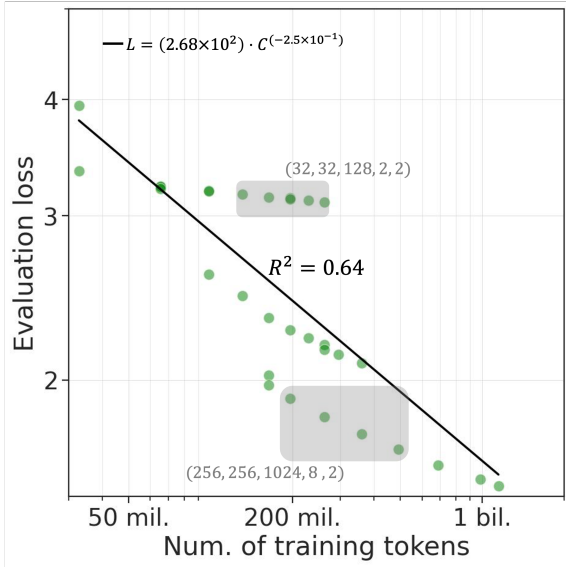


Figure 2: **Reduction of MLM loss for increasing size of pre-training tokens.** The green dots represent the compute-optimal instances found in our experiments. The solid black line represents the fitted power curve with the exact equation in the top left corner. The gray highlighted regions are examples of regions where we find a model configuration (of the format (embedding size, hidden size, intermediate size, number of layers, number of attention heads)) consistently being compute-optimal.

budget. In this region, we find the exponent for the FLOPs values to be  $-0.1412$  with an  $R^2$  value of  $0.9888$ .

In our experiments, we find few model configurations to be effective for a long range of FLOP values. We highlight a couple of examples of such configurations in Figures 2 and 3. The occurrence of such a configuration causes a discontinuous transition of MLM loss with respect to data size, model size, and FLOP values. We observe this effect to be more pronounced in the lower FLOPs region. Consequentially, a power curve does not fit in the region of lower FLOP values ( $\leq 2 \times 10^{15}$ ). With the exponent value of  $-0.0929$ , we observe the best value of  $R^2$  to be  $0.68$ .

In order to illustrate the effects of parameter count increase in the downscaled setting, we extended our pre-training experiments to include models of larger size in Figures 2 and 3. We observed that increasing the parameter count up to 100 million does not allow the model to beat the perplexity achieved by a smaller 16 million parameter model (See Appendix D).

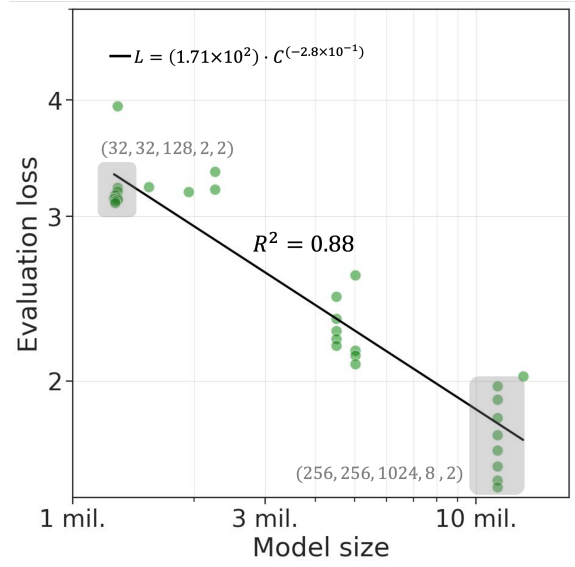


Figure 3: **Reduction in the MLM loss for increasing model size.** The green dots represent the compute-optimal instances found in our experiments. The solid black line represents the fitted power curve with the exact equation in the top left corner. The gray highlighted regions are examples of regions where we find a model configuration (of the format (embedding size, hidden size, intermediate size, number of layers, number of attention heads)) consistently being compute-optimal.

## 5.2 Incremental cost-effectiveness analysis

For the cost-effectiveness analysis, we focus on the pre-trained model configurations in **set-1**. We calculate the ICER values separately for four hyperparameters: embedding size, hidden size, intermediate size, and the number of hidden layers. We arrange the models in increasing order of FLOPs and calculate ICERs by comparing each model to the next cheapest option. The ICER values presented in Table 2 represent performance gain (reduction in perplexity) per an additional expenditure of a billion FLOPs, scaling only one hyperparameter at a time.

We observe the highest ICER value of  $3.0075$  for scaling the hidden size from  $32$  to  $64$ , refer to Table 2. For further scaling of hidden size, from  $64$  to  $128$ , and from  $128$  to  $256$ , ICER values drop at least by  $3x$  for each increment. Besides rapidly reducing values, ICERs for hidden size were always the highest, making it the most cost-effective choice. Comparably high ICERs were observed for scaling the model by increasing the number of hidden layers. For increasing the hidden layers from one to two, we record an ICER of  $2.6271$ . This value reduces to  $0.6810$  and  $0.2089$ , when scaling

Model config. (E, H, I, L, A)	ICER
(256, 32, 1024, 8, 8)	–
(256, 64, 1024, 8, 8)	3.0075
(256, 128, 1024, 8, 8)	0.8316
(256, 256, 1024, 8, 8)	0.2411
<hr/>	
(256, 256, 1024, 1, 8)	–
(256, 256, 1024, 2, 8)	2.6271
(256, 256, 1024, 4, 8)	0.6810
(256, 256, 1024, 8, 8)	0.2089
<hr/>	
(32, 256, 1024, 8, 8)	–
(64, 256, 1024, 8, 8)	0.6277
(128, 256, 1024, 8, 8)	0.2105
(256, 256, 1024, 8, 8)	0.1669
<hr/>	
(256, 256, 128, 8, 8)	–
(256, 256, 256, 8, 8)	0.4002
(256, 256, 512, 8, 8)	0.4127
(256, 256, 1024, 8, 8)	0.1970

Table 2: Incremental cost-effectiveness ratio values for increasing values of various hyperparameters. Model hyperparameters map accordingly, E: embedding size, H: hidden size, I: intermediate size, H: number of hidden layers, A: number of attention heads. Because we calculate ICER values relative to the model with the next lowest value of hyperparameter, ICERs are not calculated for the lowest values of hyperparameters.

the model with two layers to have four layers, and a model with four layers to have eight layers, respectively.

We find the ICER values for embedding size and intermediate size significantly lower than the values for hidden size and the number of hidden layers. The differences between ICERs were higher for the lower values of each hyperparameter. Although, when all hyperparameter values reached the corresponding highest values, the difference in ICERs diminished. A comparison between ICER values for embedding size and intermediate size shows that increasing embedding size from 32 to 64 brings 0.2275 more improvement in the perplexity per million FLOPs, compared to increasing intermediate size from 128 to 256. However, for all further increments in the hyperparameter values, increasing intermediate size results in at least 0.03 more ICER value than for embedding size.

### 5.3 Downstream evaluation

We report fine-tuning performance on vocabulary-filtered GLUE benchmark in Table 3 (cf. Section 3.1). For reference, we also report performance on unfiltered GLUE. We find that GLUE performance peaks for models with 2 and 4 hidden layers with average scores of 59.39 and 60.67, respectively. Interestingly, we find the average GLUE

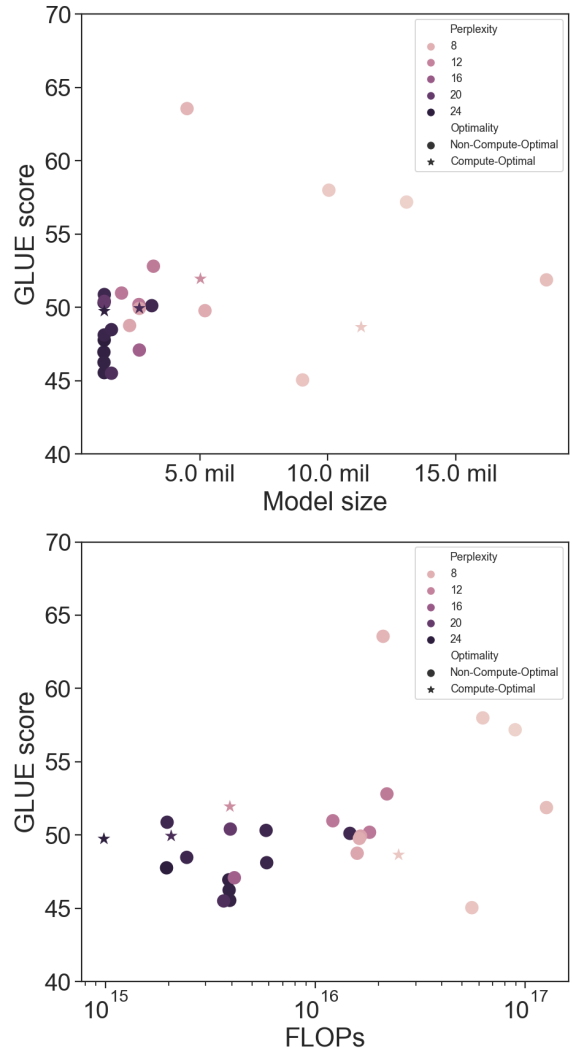


Figure 4: GLUE scores by model size and FLOPs for compute-optimal and non-compute-optimal models

score decreases for the model with 8 hidden layers to 56.29. Such reduction is not observed when increasing the hidden size or the embedding size. As expected, models consistently demonstrate better performance on vocabulary-filtered GLUE. We also see that model performance is strongest for models with 2 and 4 hidden layers assessed on vocabulary-filtered GLUE.

To assess whether pre-training effects are beneficial in the downscaled setting, we compare the average GLUE score of each pre-trained model with the score of the same model fine-tuned without pre-training. Table 3 shows that for each model shape, the fine-tuned model outperforms its respective randomly initialized counterpart. Our results show that in a reduced-vocabulary setting, the advantages of pre-training are observable even for smaller models, starting with a 1.25M parameter

Model config. (E, H, I, L, A)	Model size (mil. parameters)	FLOPs ( $\times 10^{15}$ )	Perplexity	GLUE score (unfiltered)	GLUE score (filtered)	GLUE score (filtered, w/o PT)
(256, 256, 1024, 8, 8)	16.24	110	4.80	51.73	56.29	40.24
<b>(32, 256, 1024, 8, 8)</b>	11.89	80	5.60	46.99	48.39	49.80
<b>(64, 256, 1024, 8, 8)</b>	12.51	84	5.31	45.12	51.09	51.99
<b>(128, 256, 1024, 8, 8)</b>	13.75	92	5.11	48.07	52.73	<b>52.09</b>
(256, <b>32</b> , 1024, 8, 8)	6.10	42	10.42	47.98	50.23	45.93
(256, <b>64</b> , 1024, 8, 8)	7.34	50	7.56	49.34	53.78	50.67
(256, <b>128</b> , 1024, 8, 8)	10.04	69	5.88	51.16	55.63	50.60
(256, 256, <b>128</b> , 8, 8)	7.63	85	5.61	57.65	57.18	41.62
(256, 256, <b>256</b> , 8, 8)	8.15	88	5.45	54.91	56.48	41.28
(256, 256, <b>512</b> , 8, 8)	9.20	96	5.12	51.60	57.36	40.55
(256, 256, 1024, <b>1</b> , 8)	10.71	73	7.60	50.87	55.99	50.53
(256, 256, 1024, <b>2</b> , 8)	11.50	79	6.07	54.31	59.39	51.17
(256, 256, 1024, <b>4</b> , 8)	13.08	89	5.28	53.85	<b>60.67</b>	47.15
(256, 256, 1024, 8, <b>1</b> )	16.24	110	4.74	50.47	53.68	40.30
(256, 256, 1024, 8, <b>2</b> )	16.24	110	<b>4.67</b>	50.39	54.98	40.10
(256, 256, 1024, 8, <b>4</b> )	16.24	110	4.75	49.55	54.21	39.66
(32, 32, 128, 2, 2)	1.27	8.57	20.07	46.25	49.03	48.68
(32, 32, 128, 1, 1)	1.25	8.60	23.40	44.97	48.22	47.98
(32, 32, 64, 1, 1)	1.25	8.71	23.42	44.91	47.20	48.69

Table 3: Performance on vocabulary-filtered and unfiltered GLUE benchmarks. These models have access to the same number of tokens during pre-training. Model hyperparameters map accordingly: Embedding Size: E, Hidden Size: H, Intermediate Size: I, Hidden Layers: L, Attention Heads: A. In the last column, we report the fine-tuning results for the randomly initialized (without pre-training) counterpart of the model.

count.

We further fine-tune a set of 27 models, comprising a mix of compute-optimal and non-compute-optimal checkpoints, to better understand the relation between upstream and downstream performance. In Figure 4, we plot each model’s GLUE score against its size and number of FLOPs, with color indicating the test perplexity of each model. We observe that for a given parameter count, compute-optimal models do not necessarily outperform the undertrained models on the GLUE benchmark.

Lastly, considering all fine-tuning results together, we conduct a test to measure the correlation between perplexity (upstream performance) and average GLUE score (downstream performance). We find that the correlation between the average GLUE score for unfiltered GLUE datasets and model perplexity is inconclusive with the Spearman coefficient value of -0.17 and a p-value of 0.28. On the other hand, we find average GLUE score calculated for filtered GLUE datasets highly correlates with model perplexity with the Spearman coefficient value of -0.67 and a p-value  $\leq 0.01$ .<sup>5</sup>

<sup>5</sup>Our reported values of the test metric are calculated over a sample size of 32. The exact p-value is  $6.25 \times 10^{-7}$ .

## 5.4 Comparison with unconstrained text

To highlight the ability of smaller models to benefit from pre-training when the language size is reduced (rather than unconstrained), we also pre-trained 10 models on unconstrained language (i.e., without any vocabulary reduction). We provide the details of the data collection, tokenizer training, and the experimental setup in Appendix C.

Table 4 shows the relative performance figures for the models trained on unconstrained language. We fix the model configuration and report the change in performance relative to the constrained (limited-vocabulary) case. The training data size and hyperparameter setting for pre-training and fine-tuning are kept the same.

Note that for the unconstrained case, there is a considerable increase in the model size due to an increase in the Byte-BPE vocabulary (from 19,000 to 29,000). The increased model size also increases the compute cost by at least 32.87%. Despite the increased model size and compute cost, no corresponding improvement in pre-training performance is observed, as shown in Table 4. In fact, although the perplexity on reduced-vocabulary data decreases with model size, none of the model configurations studied reach the MLM perplexity of



Model config. (E, H, I, L, A)	$\Delta$ Model size (mil. parameters)	% $\Delta$ FLOPs	% $\Delta$ Perplexity	% $\Delta$ GLUE score (filtered)	% $\Delta$ GLUE score (unfiltered)
(256, 256, 1024, 8, 8)	5.13	32.87	23.96	-18.28	-11.79
(256, <b>32</b> , 1024, 8, 8)	2.89	47.59	17.86	-3.58	2.15
(256, <b>64</b> , 1024, 8, 8)	3.21	44.01	16.16	-3.32	1.60
(256, <b>128</b> , 1024, 8, 8)	3.85	38.99	20.85	-4.97	-0.08
(256, 256, 1024, <b>1</b> , 8)	5.13	49.44	18.90	-15.39	-1.90
(256, 256, 1024, <b>2</b> , 8)	5.13	46.12	17.24	-18.99	-7.80
(256, 256, 1024, <b>4</b> , 8)	5.13	40.66	16.93	-19.93	-7.61
(32, 32, 128, 2, 2)	0.65	51.23	18.90	-8.69	-1.00
(32, 32, 128, 1, 1)	0.65	51.89	25.79	-9.06	1.19
(32, 32, 64, 1, 1)	0.65	52.07	14.78	-8.94	-2.05

Table 4: Relative performance of the models trained on text data without any vocabulary filtration (unconstrained case). In this table, we present the results relative to the same configuration model trained on limited-vocabulary language data (constrained case). The  $\Delta$  symbol represents the difference  $\Delta = value_{unconstrained} - value_{constrained}$  in the results and  $\% \Delta$  represents the percentage difference  $\Delta = 100 \times \Delta / value_{constrained}$ . Perplexity values are calculated on the test split of constrained data. Positive values for  $\% \Delta$  perplexity mean a degradation of performance for the unconstrained case. Note that for the unfiltered GLUE, we do not expect limited-vocabulary models to do well. However, we observe that besides adding a compute overhead, pre-training language models on unconstrained data degrade the pre-training and downstream performance in the models with  $\leq 22$  million parameters.

the reduced-scale model, when evaluated on the test split of the data.

Since perplexity values are directly impacted by the increase in the Byte-BPE vocabulary size, we also evaluate the unconstrained data models on GLUE benchmarks for fairer evaluation. Similar to pre-training results, we observe a consistent degradation of the performance for the filtered versions of the GLUE benchmark. For all model configurations considered, the average GLUE score for filtered datasets reduces by up to  $\approx 20\%$  due to pre-training on unconstrained data. On the other hand, for the unfiltered version of the GLUE datasets, we do not expect models trained in limited-vocabulary data to do well. However, we find that the average GLUE score on unfiltered datasets improves only in three of the 10 model configurations we considered. These results further confirm that limiting the vocabulary benefits the models with  $\leq 22$  million parameters.

One possible explanation for this is the relative contribution of embedding parameters to model size. In an unconstrained setting, vocabulary embedding parameters account for most of the model, with no parameters left for transformer blocks.<sup>6</sup> There is prior work showing that pre-training loss improves only minimally for models with zero transformer blocks (Kaplan et al., 2020). Thus, con-

<sup>6</sup>For example, for a model with 10M parameters, if the embedding size is 200 with a vocabulary of 50,000, no transformer block can be added. But if the vocabulary size is 20,000, 6M parameters can be used in transformer blocks.

straining vocabulary allows one to increase transformer block capacity while otherwise maintaining a small parameter count.

## 6 Conclusions & Future Work

In this study, we investigated whether reducing language size allows the benefits of pre-training to be observed in a downscaled setting for models with  $\leq 20$ M parameters. We evaluated a range of model configurations and found that the advantages of pre-training are observable even for models with as few as 1.25M parameters, with a strong correlation between upstream and downstream performance. However, we also observed that compute-optimal training does not appear to be crucial for downstream results and that parameter count does not reliably predict upstream performance. Furthermore, we observed a break of the FLOP-Perplexity power law at the  $2.2 \times 10^{15}$  FLOP region, which shows the limited applicability of scaling laws.

Overall, our experiments provide insight into the behavior of small language models in a downscaled language setting. The next logical steps as a follow-on to this work would be to check whether generative models would demonstrate any emergent abilities in a downscaled setting.

## 7 Limitations

While we do explore a range of models in the 1-20M parameter space, our work does not constitute a complete study of downscaling. In this work,

we aimed to explore the more fundamental components of model shape, model size, and input data.

However, our findings may not generalize to other models with alternative applications of downscaling methods. Considering it to be out of scope for this study’s assessment of pre-training effects, we did not compare our results to knowledge distillation methods of similar model shape and size. Furthermore, our exploration of model shape and size was limited to a model’s hidden size, number of hidden layers, embedding size, and intermediate size, and number of attention heads as these are the most commonly-tuned hyperparameters.

Our usage of vocabulary filtration as a means of downscaling input data size may not be the most effective means of limiting input data. While shown to be effective, alternative approaches for input data manipulation such as curriculum learning, and data pruning merit study beyond the scope of this paper.

## Ethics Statement

Our exploration of smaller language models presents a number of implications for accessibility, environmental impact, and cost. By exploring models in the space of 1-20M parameters, our findings can inform language modeling work for those without access to large, GPU-enabled environments. This is important as it can encourage further research work in this space by those who are otherwise unable to work with SoTA LLMs. We acknowledge that our resources enabled the breadth of study in this paper; most of this study was conducted using a single GPU. This consideration underscores our commitment to improving accessibility for under-resourced technologists throughout the world. Furthermore, in working with downscaled LLMs, we hope to encourage methods that reduce overall carbon footprint and bolster sustainable practices in NLP. These considerations are especially important given the particular burden placed on those with limited access to electricity and technology. The cost of running and experimenting with these models may prove quite costly in terms of person-hours and compute resources. As such, we hope our work at smaller scale can help lessen these burdens, and positively impact the lives of technologists, and others. Any model from the study can be trained in less than a day on a single consumer-grade GPU.

## References

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450*.
- Kiran Bambha and W Ray Kim. 2004. Cost-effectiveness analysis and incremental cost-effectiveness ratios: uses and pitfalls. *European journal of gastroenterology & hepatology*, 16(6):519–526.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D 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 Ziegler, Jeffrey Wu, Clemens Winter, Chris 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*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Stephanie CY Chan, Adam Santoro, Andrew Kyle Lampinen, Jane X Wang, Aaditya K Singh, Pierre Harvey Richemond, James McClelland, and Felix Hill. 2022. Data distributional properties drive emergent in-context learning in transformers. In *Advances in Neural Information Processing Systems*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pilla, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. [Palm: Scaling language modeling with pathways](#). *arxiv:2204.02311*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- William Fedus, Barret Zoph, and Noam Shazeer. 2022. [Switch transformers: Scaling to trillion parameter](#)

- models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39.
- Jack FitzGerald, Shankar Ananthkrishnan, Konstantine Arkoudas, Davide Bernardi, Abhishek Bhagia, Claudio Delli Bovi, Jin Cao, Rakesh Chada, Amit Chauhan, Luoxin Chen, et al. 2022. Alexa teacher model: Pretraining and distilling multi-billion-parameter encoders for natural language understanding systems.
- Thamme Gowda and Jonathan May. 2020. Finding the optimal vocabulary size for neural machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3955–3964. Online. Association for Computational Linguistics.
- Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. 2015. The goldilocks principle: Reading children’s books with explicit memory representations. *arXiv preprint arXiv:1511.02301*.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. Training compute-optimal large language models.
- Philip A Huebner, Elior Sulem, Fisher Cynthia, and Dan Roth. 2021. Babyberta: Learning more grammar with small-scale child-directed language. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 624–646.
- Philip A Huebner and Jon A Willits. 2021. Using lexical context to discover the noun category: Younger children have it easier. In *Psychology of learning and motivation*, volume 75, pages 279–331. Elsevier.
- 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*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Jorge J Moré. 1978. The levenberg-marquardt algorithm: implementation and theory. In *Numerical analysis*, pages 105–116. Springer.
- Laura Pérez-Mayos, Miguel Ballesteros, and Leo Wanner. 2021. How much pretraining data do language models need to learn syntax?
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Winogrande: An adversarial winograd schema challenge at scale. *Commun. ACM*, 64:99–106.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2019. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*.
- Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S. Morcos. 2022. Beyond neural scaling laws: beating power law scaling via data pruning. In *Advances in Neural Information Processing Systems*.
- Yi Tay, Mostafa Dehghani, Jinfeng Rao, William Fedus, Samira Abnar, Hyung Won Chung, Sharan Narang, Dani Yogatama, Ashish Vaswani, and Donald Metzler. 2022. Scale efficiently: Insights from pretraining and finetuning transformers. In *International Conference on Learning Representations*.
- Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Well-read students learn better: On the importance of pre-training compact models.
- Ekaterina Voloshina, , Oleg Serikov, Tatiana Shavrina, and and. 2022. Is neural language acquisition similar to natural? a chronological probing study. In *Computational Linguistics and Intellectual Technologies. RSUH*.
- 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. *arXiv preprint arXiv:1804.07461*.
- Yian Zhang, Alex Warstadt, Haau-Sing Li, and Samuel R. Bowman. 2020. When do you need billions of words of pretraining data?
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *The IEEE International Conference on Computer Vision (ICCV)*.

## A Calculation of compute cost (FLOPs)

We adopt the same approach for calculating the compute cost (FLOPs) as presented by Hoffmann et al. (2022). For notational convenience, we denote the sequence length, vocabulary size, embedding size, hidden size, intermediate size (hidden dimension of the feed-forward network in the transformer block), number of attention heads, key size (for the attention block), and number of layers by  $S$ ,  $V$ ,  $E$ ,  $H$ ,  $I$ ,  $A$ ,  $K$  and  $L$ , respectively.

The FLOPs for the forward pass are calculated as follows.

- Single embedding block:

$$C_{emb} = 2 \times S \times (VE + EH)$$

- Single attention block:

- Cost of the key, query, and value projections

$$C_{att} = 2 \times 3 \times SH \times (KA)$$

- Cost of the dot product operation of key and query

$$C_{att} += 2 \times SS \times (KA)$$

- Cost of the softmax operation

$$C_{att} += 3 \times SS \times (A)$$

- Cost of the query reduction

$$C_{att} += 2 \times SS \times (KA)$$

- Cost of the final linear layer

$$C_{att} += 2 \times SH \times (KA)$$

- Single feed-forward layer:

$$C_{int} = 2 \times (HI + IH)$$

- Single language model head:

$$C_{lmh} = 2 \times SHV$$

- Single forward pass:

$$C_{forward} = C_{emb} + C_{lmh} + L \times (C_{att} + C_{int})$$

- Single backward pass:

$$C_{backward} = 2 \times C_{forward}$$

- Cost per training sequence:

$$C_{seq} = C_{forward} + C_{backward}$$

Therefore, we calculate the total compute cost as  $Number\ of\ parameter\ updates \times batch\ size \times C_{seq}$ .

## B Tokenizer selection

### B.1 List of reference words for ESMS

In Table 5 we provide a list of words we used for calculating the Exact Sub-token Matching Score (ESMS). We also provide the morpheme-based tokens per word and the maximum value of exact matches per word.

### B.2 Comparison of different vocabulary sizes and tokenizer types

We provide the results of our experiments to determine the best-suited tokenizer in this section. Table 6 provides the vocabulary size, corresponding word-split ratio, and ESMS value for the three types of tokenizers we evaluated. The bolded row is the final tokenizer we used in our pre-training experiments.

## C Comparison with unconstrained language

Our main experiments were conducted on language that is constrained by a predefined vocabulary. To study the effect of the applied vocabulary constraint in comparison with free text, we conduct a set of experiments on unconstrained language i.e., without any vocabulary-based filtering. In the following subsections, we provide details of the data collection, tokenizer training, and pre-training process adopted.

### C.1 Pre-training data

Our objective in curating constrained language was solely to impose a vocabulary constraint. However, our filtering method (Section 3.2) resulted in constrained language comprised of non-consecutive text sequences. To address differences beyond vocabulary, we conducted unconstrained language collection using the following approach. We divided all instances in a specific corpus into spans of

Reference word	Morpheme sub-tokens	Maximum exact matches per word
Cooking	cook, ing	2
Dangerous	danger, ous	2
Pretext	pre, text	2
Fitness	fit, ness	2
Antisocial	anti, social	2
Podium	pod, ium	2
Universe	uni, verse	2
European	europ, ean	2
Decode	de, code	2
Subvert	sub, vert	2
Proactive	pro, active	2
Concentric	con, centr, ic	3
Octopus	octo, pus	2

Table 5: Reference words and corresponding sub-word tokens for calculating Exact Sub-token Matching Score (ESMS) for the tokenized text

Tokenizer name	Vocabulary size	Word-split ratio	ESMS
BPE (Radford et al., 2019)	18,000	1.34	0.2868
BPE (Radford et al., 2019)	<b>19,000</b>	<b>1.32</b>	<b>0.2604</b>
BPE (Radford et al., 2019)	20,000	1.31	0.2490
BPE (Radford et al., 2019)	Pre-trained	1.32	0.1547
WordPiece (Devlin et al., 2018)	<u>16,000</u>	<u>1.17</u>	<u>0.0339</u>
WordPiece (Devlin et al., 2018)	<u>17,000</u>	1.17	0.0264
WordPiece (Devlin et al., 2018)	18,000	1.16	0.0188
WordPiece (Devlin et al., 2018)	Pre-trained	1.17	0.0339
SentencePiece (Raffel et al., 2020)	9,000	1.32	0.0301
SentencePiece (Raffel et al., 2020)	<u>10,000</u>	<u>1.29</u>	<u>0.0226</u>
SentencePiece (Raffel et al., 2020)	11,000	1.26	0.0188
SentencePiece (Raffel et al., 2020)	Pre-trained	1.29	0.0339

Table 6: Values of word-split ratio and ESMS for various tokenizer and vocabulary size settings.

110 words and randomly sampled spans. The number of randomly sampled spans was determined to maintain the same data distribution across different corpora, as indicated in Table 1. This method aimed to minimize the impact of data features other than vocabulary. We gathered an equivalent number of training sequences (approximately nine million) as in the constrained pre-training data. Finally, we ensured a fair comparison of pre-training performance by using the same evaluation and test split for both pre-training datasets.

## C.2 Tokenizer

After data curation, we conduct experiments with various tokenizers to finalize the tokenizer type and size of the token vocabulary for the language model. These experiments were conducted in the same manner described in Section 3.3. The final tokenizer we select is the Byte-BPE tokenizer (Radford et al., 2019) with a vocabulary of 29,000 tokens ( $1.6\times$  that of the vocabulary size for the constrained language). The word-split ratio and the ESMS (exact sub-token matching score) for the final tokenizer were 1.53 and 0.2339.

## C.3 Experimental setup

After finalizing the token vocabulary, we measure the pre-training as well as the downstream performance of the models trained on unconstrained language. We focus on the model configuration explored in the **set-1** (refer to section 4.1). Furthermore, guided by our results in the ICER analysis (refer to section 5.2), we only consider the model configurations that either perturb the hidden size or the number of layers in the model. With such a selection, we pre-train seven language models. In addition, we pre-train the smallest model configuration that highlighted the benefits of pre-training in our main experiment (refer to Table 3). Overall, we pre-train 10 models on the collected unconstrained language. We keep all the hyperparameter values the same as in our main experiments (refer to Sections 4.2 and 4.3). For the comparison of the pre-training performance, we measure the MLM loss and perplexity values calculated on the test split of the limited-vocabulary pre-training data. For comparison of the downstream performance, we finetune the final checkpoint of all pre-trained models on GLUE tasks and record the average GLUE score, separately for filtered and unfiltered versions of the GLUE datasets.

## D Training larger models

We continued our pre-training experiments with the constrained language (limited vocabulary) data to include larger models i.e., models with more than 20 million parameters. We first set the anchor configuration to have embedding size, hidden size, intermediate size, number of layers, and number of attention heads equal to 512, 512, 2048, 8, and 8, respectively. After defining the anchor configuration, we follow the same approach of varying each configuration feature to explore and pre-train various models. However, for the training of larger models, we only focused on the hidden size and number of layers.

We present the pre-training results of the larger models in Table 7. In the larger model configurations, we observe that the additional model parameters do not reduce the model perplexity below 4.67. However, within the set of larger models, we observe an expected reduction in the perplexity values with an increase in the size of the model.

For our main experiments, we calculated the training data size for the expected largest model size i.e., 20 million parameters based on the findings provided by Hoffmann et al. (2022). The data size value was between 500 to 600 million tokens. Note that this data size is the size required to train the 20 million parameter model ‘compute-optimally’. Hence, we collected more data to observe the effect after the compute-optimal point. Finally, we collected approximately double the quantity of data ( $\geq 1100$  million tokens). Findings provided by (Hoffmann et al., 2022) are based on decoder-only models but, at the time of the experimentation, this was the best available guide for us to make a decision. Hence, we speculate that the size of our filtered pre-training data is not sufficient for the larger models that we consider in this set of experiments where we pre-train considerably larger models. Therefore, we do not include the model configuration considered in this set of experiments for our main results and figures.

<b>Model config.</b> <b>(E, H, I, L, A)</b>	<b>Model size</b> <b>(mil. parameters)</b>	<b>FLOPs</b> <b>(<math>\times 10^{15}</math>)</b>	<b>Perplexity</b>
(256, 256, 1024, 8, 2)	16.24	110	4.67
(512, 512, 2048, 8, 8)	45.30	302	5.57
(512, <b>256</b> , 2048, 8, 8)	25.41	174	6.47
(512, <b>768</b> , 2048, 8, 8)	69.52	452	5.22
(512, <b>1024</b> , 2048, 8, 8)	98.06	624	4.94
(512, 512, 2048, <b>1</b> , 8)	23.23	158	13.27
(512, 512, 2048, <b>2</b> , 8)	26.38	179	6.67
(512, 512, 2048, <b>4</b> , 8)	32.69	220	5.75

Table 7: Pre-training results of larger models trained on the limited-vocabulary text data. In this table, we provide model configuration, model size, compute-cost (FLOPs), and model perplexity for the pre-training experiments we conducted on language models larger than the models pre-trained in our main experiments (set-1). We provide the best pre-training performance found in our main experiments (set-1) as a baseline (the first row). We find that additional parameters do not reduce model perplexity.

## ACL 2023 Responsible NLP Checklist

---

### A For every submission:

- A1. Did you describe the limitations of your work?

6

- A2. Did you discuss any potential risks of your work?

*Our study works with understanding models at a very small scale (<10M params). These models, unlike LLMs do not present harm in terms of environmental impact or misuse as they are not as capable.*

- A3. Do the abstract and introduction summarize the paper's main claims?

1

- A4. Have you used AI writing assistants when working on this paper?

*We used ChatGPT and GPT-3 to help us rephrase our text to make it more naturally looking and fluent. After rephrasing we could carefully read and edit the output of the system if necessary. In some cases we would completely discard the generated text in favor of manual text editing. Additionally, we used Grammarly to check the typos, grammar, and phrasing. These methods were used throughout the paper.*

### B Did you use or create scientific artifacts?

*Results*

- B1. Did you cite the creators of artifacts you used?

*I don't understand the question. If this refers to pre-trained models and datasets, then the methods and results section and a big chunk of the rest of the paper.*

- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?

*We will provide the licence when the dataset is published. Also, legally, the license must include the names of the authors, which contradicts anonymity.*

- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?

*We use open-access datasets that are not restricted.*

- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?

*Not applicable. Left blank.*

- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

*It will be provided upon dataset publication*

- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.

*Left blank.*

---

*The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.*



**C  Did you run computational experiments?**

*Results*

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?

*Results*

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

*Experimental setup*

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

*Results*

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

*It will be available in the published code, we can't publish code now according to the anonymity policy*

**D  Did you use human annotators (e.g., crowdworkers) or research with human participants?**

*Left blank.*

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

*No response.*

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

*No response.*

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

*No response.*

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

*No response.*

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

*No response.*