# BTLM-3B-8K: 7B Parameter Performance in a 3B Parameter Model

**Nolan Dey**[*1], **Daria Soboleva**[*1], **Faisal Al-Khateeb**[1], **Bowen Yang**[1], **Ribhu Pathria**[1], **Hemant Khachane**[1], **Shaheer Muhammad**[1], **Zhiming (Charles) Chen**[1], **Robert Myers**[2], **Jacob Robert Steeves**[2], **Natalia Vassilieva**[1], **Marvin Tom**[1], **Joel Hestness**[1]

1. Cerebras Systems    2. OpenTensor Foundation

*{nolan,daria.soboleva}@cerebras.net*

## Abstract

We introduce the Bittensor Language Model, called "BTLM-3B-8K", a new state-of-the-art 3 billion parameter open-source language model. BTLM-3B-8K was trained on 627B tokens from the SlimPajama dataset with a mixture of 2,048 and 8,192 context lengths. BTLM-3B-8K outperforms all existing 3B parameter models by 2-5.5% across downstream tasks. BTLM-3B-8K is even competitive with some 7B parameter models. Additionally, BTLM-3B-8K provides excellent long context performance, outperforming MPT-7B-8K and XGen-7B-8K on tasks up to 8,192 context length. We trained the model on a cleaned and deduplicated SlimPajama dataset; aggressively tuned the μP hyperparameters and schedule; used ALiBi position embeddings; and adopted the SwiGLU nonlinearity.

On Hugging Face, the most popular models have 7B parameters, indicating that users prefer the quality-size ratio of 7B models. Compacting the 7B parameter model to one with 3B parameters, with little performance impact, is an important milestone. BTLM-3B-8K needs only 3GB of memory with 4-bit precision and takes 2.5x less inference compute than 7B models, helping to open up access to a powerful language model on mobile and edge devices. BTLM-3B-8K is available under an Apache 2.0 license on Hugging Face: `https://huggingface.co/cerebras/btlm-3b-8k-base`.

## 1 Introduction

Large language models (LLMs) can perform a diverse collection of text-based tasks with brief instructions (Brown et al., 2020a), making them useful in many settings. Applications include natural language understanding, content generation, and computer programming. With the ability to generate coherent text, answer questions, translate languages, and summarize long documents, LLMs are transforming the way we interact with and leverage information.

With LLaMa Touvron et al. (2023a) it became possible to inefficiently train LLMs (Hoffmann et al., 2022) on trillions of tokens to achieve state of the art parameter efficiency. The resulting LLaMA models introduced the community to powerful open-source LLMs that can be deployed on a high-end laptop[1]. Since then, there have been many reproductions and extensions of LLaMA models (Together.ai, 2023; Geng & Liu, 2023; Tow, 2023; Almazrouei et al., 2023; Penedo et al., 2023; Nijkamp et al., 2023; Team, 2023b;a; Touvron et al., 2023b) with the 7B parameter size being the most popular due to its performance and portability.

But while users want the quality of 7B models, such models have memory and compute requirements that are prohibitively costly in many settings. Even with compression techniques such as quantization (Frantar et al., 2022), edge devices such as mobile phones and laptops generally do not have enough memory capacity to hold 7B model weights, and inference tends to be slow.

---

[*]Equal contribution
[1]`https://github.com/ggerganov/llama.cpp`

1

Another shortcoming of existing LLMs is that they don't support long contexts. The ability to model long-range contextual dependencies is essential for tasks such as summarizing or answering questions about long-form text, processing entire codebases, predicting DNA sequences, engaging in multi-turn dialogues, or generating content for articles.

In this work, we introduce the Bittensor Language Model "BTLM-3B-8K", a new state-of-the-art 3B parameter, open-source language model. Our model is competitive with 7B parameter models that were trained with 3.3× more compute, 2.5× more parameters, and 1.6× more tokens. BTLM-3B-8K can fit on devices with 3GB of RAM and requires 2.5x less inference compute than 7B models, enabling access to the performance of 7B models on billions of edge devices worldwide. BTLM-3B-8K uses ALiBi position embedding (Press et al., 2021) and is trained with up to 8,192 context length, enabling long context performance competitive with existing 7B parameter models.

Our contributions are as follows:

- **Training Procedure:** We detail the procedure we used to train BTLM-3B-8K on one epoch of the SlimPajama dataset using CG-1, a cluster of 64 Cerebras CS-2 Systems.

- **Model Evaluation:**

  - We provide extensive comparison of existing 3B and 7B parameter models on 22 benchmarks, evaluating common sense reasoning, world knowledge, reading comprehension, code generation, long sequence interpolation, long sequence extrapolation, bias, toxicity, and misinformation.
  - We demonstrate that BTLM-3B-8K sets the standard for 3B parameter models and often outperforms 7B models.

- **Training Improvement Ablations:** We perform ablations of the architectural changes and training methods that drove BTLM's superior performance, achieving a 5.36% improvement in loss over the baseline.

- **Releases and Availability:** We release the BTLM-3B-8K weights and the SlimPajama dataset we used to train BTLM with an Apache 2.0 license on Hugging Face: `https://huggingface.co/cerebras/`. We trust that these contributions can be of significant value to the open-source community.

## 2 BTLM Architecture and Training

### 2.1 Model Architecture

BTLM-3B-8K is an autoregressive transformer decoder model (Brown et al., 2020a) based on the GPT-3 architecture with fully dense attention. We make three architectural changes motivated by the experiments described in Section 4:

- SwiGLU nonlinearity (Shazeer (2020)) instead of GELU.

- ALiBi position embeddings (Press et al. (2021)) instead of learned position embeddings. This enables improved extrapolation to longer sequence lengths not seen during training.

- Maximal update parameterization ($\mu$P, Yang et al. (2021)) instead of the standard parameterization (SP). This involves applying scalar multipliers to the learning rate, output, and initialization of certain layers to counteract activation scales growing with width.

BTLM-3B-8K has the following model shape parameters: $d_{model}$=2560, $n_{layers}$=32, $d_{head}$=80, $d_{ffn}$=6826. This yields 2.6B model parameters. which we round to 3B as is now conventional.

| Data source | RedPajama Doc Filtr. % | RedPajama Byte Dupl. % | SlimPajama Proportion % |
|---|---|---|---|
| Commoncrawl | 0.02 | 63.76 | 52.20 |
| C4 | 4.70 | 6.85 | 26.70 |
| GitHub | 0.00 | 46.16 | 5.20 |
| Books | 0.00 | 2.01 | 4.20 |
| ArXiv | 0.62 | 0.06 | 4.60 |
| Wikipedia | 0.00 | 2.24 | 3.80 |
| StackExchange | 0.32 | 0.20 | 3.30 |
| Total | 1.86 | 49.60 | 100.00 |

Table 1: Document low-length filter rates and data source byte duplication rates found in RedPajama, in addition to final SlimPajama data source proportions.

## 2.2 Pretraining Data

Aspects of data quality, for example the data source mix, filtering methodology, and duplication rate, can have a significant impact on LLM performance. To bolster BTLM's performance, we create a high quality 627B token dataset called SlimPajama ((Soboleva et al., 2023)). Starting from the 1.21T token RedPajama dataset Computer (2023), we apply filtering and deduplication to improve data quality. First, we remove documents containing fewer than 200 characters, as we find these typically contain only metadata. Next, we perform global deduplication using MinHashLSH (Leskovec et al., 2014) to extensively remove documents with significant overlapping text. Table 1 shows a breakdown of the filter rate and deduplication rate for each of SlimPajama's data sources. Finally, we tokenize the data with byte-pair encoding using the the GPT-2 vocabulary with 50257 tokens ((Sennrich et al., 2016; Radford et al., 2019)). Overall, SlimPajama contains 627B tokens after tokenization. The SlimPajama dataset is available on https://huggingface.co/datasets/cerebras/SlimPajama-627B. We also release our preprocessing code under https://github.com/Cerebras/modelzoo/transformers/data_processing/slimpajama.

## 2.3 Training Procedure

BTLM-3B-8K was trained in two phases while holding batch size constant in terms of number of tokens:

1. 470B tokens with a sequence length of 2,048 and a batch size of 1920 sequences

2. 157B tokens with a sequence length of 8,192 and a batch size of 480 sequences

We used the AdamW optimizer Loshchilov & Hutter (2017) with $\beta_1 = 0.9, \beta_2 = 0.95, \epsilon = 10^{-8}$, weight decay of 0.1, and gradient clipping to a maximum norm of 1.0. Since we are using $\mu$P, the learning rate for each layer is derived from the base learning rate. We use a maximum base learning rate of 1.2e-2. We use a linear warmup length of 375M tokens, followed by a linear decay from the maximum base learning rate of 1.2e-2 down to 1.0198e-04. The base initialization standard deviation used was 0.073. In addition, we introduce two tunable scalar multipliers for the embedding output and output logits ((Yang et al., 2021)). We used an embedding multiplier of 14.6 and an output logit multiplier of 2.22. The base learning rate, base initialization standard deviation, embedding multiplier, and output logit multiplier were found through a random hyperparameter search with a 40M parameter proxy model following Yang et al. (2021); Dey et al. (2023).

## 2.4 Training Loss Stability

It is common for LLMs to encounter loss instability which can lead to loss divergence and require careful manual interventions to recover training ((Zhang et al., 2022; Chowdhery et al., 2022)). Figure 1 shows that BTLM training progressed with excellent loss stability, especially given how large our learning rate is relative to other models. We attribute this stability to the maximal update parameterization which controls
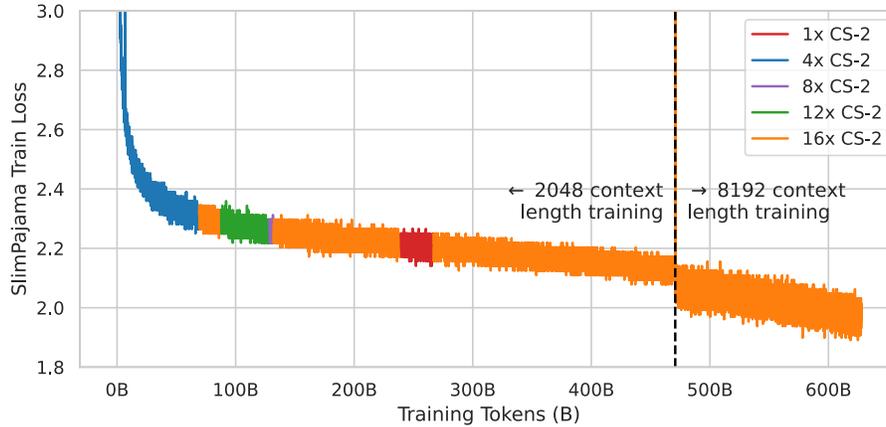
Figure 1: SlimPajama train cross-entropy loss versus training tokens. Training was scaled between different numbers of CS-2 systems depending on cluster availability.

activation explosion as model width is scaled up. BTLM only experienced two loss spikes: one at step 15 (59M tokens) and another at the transition to 8,192 context length training as the model adapts to longer sequences. The training fully recovered from both spikes, and they did not seem to impact the overall trajectory of the loss curve.

### 2.5 Hardware

BTLM was trained on the Condor Galaxy 1 (CG-1) AI supercomputer, a cluster of 64 Cerebras CS-2 systems built by G42 and Cerebras. Unlike GPU or TPU clusters, CS-2 clusters exclusively use data parallelism to scale to multiple systems during training[2], eliminating the complexity of splitting models into smaller chunks using tensor or pipeline model parallelism. During training, we needed to interleave training with other high priority jobs on the cluster. Thanks to the simplicity of data parallelism, we easily scaled up and down our training to different numbers of CS-2 nodes with near-linear speedups and without any code or configuration changes. Figure 1 shows the training loss curve for the training run, and different portions of the run colored according to the number of CS-2 machines on which that phase was trained. We encountered no hardware failures during training, demonstrating the reliability of the Cerebras wafer-scale cluster.

## 3 Model Evaluation

In this section, we compare BTLM-3B-8K model with 3B and 7B parameters open-source foundation models: RedPajama-INCITE (Together.ai, 2023), OpenLLaMA (Geng & Liu, 2023), StableLM-v2 (Tow, 2023), Falcon (Almazrouei et al., 2023), Falcon-RW (Penedo et al., 2023), XGen (Nijkamp et al., 2023), MPT (Team, 2023b;a), LLaMA (Touvron et al., 2023a), and LLaMA-2 (Touvron et al., 2023b).

Following Brown et al. (2020b), we evaluate models on zero-shot and few-shot tasks using the Eleuther AI evaluation harness framework Gao et al. (2021). To provide a more holistic view, we measure model capability across a wide variety of task domains: common sense reasoning (CSR), world knowledge (WK), reading comprehension (RC), massive multitask language understanding (MMLU), and coding abilities. In Table 2, we show the average accuracy within each task domain for 3B and 7B open-source base models. By reporting average accuracy across tasks within a domain we hope to provide a more accurate picture by smoothing out the high variability that individual tasks might introduce.

BTLM-3B-8K achieves state-of-the-art performance among 3B parameter models, outperforming others by a substantial margin while using the least pretraining compute and data. BTLM-3B-8K was trained on 627B tokens, significantly less than RedPajama-INCITE-3B at 800B tokens and OpenLLaMA 3Bv2 at 1T

---

[2]https://www.cerebras.net/blog/linear-scaling-made-possible-with-weight-streaming

tokens. BTLM-3B is even competitive with 7B models, outperforming RedPajama-INCITE-7B (Together.ai, 2023), OpenLLaMA-7B Geng & Liu (2023), and StableLM-Alpha-7B-v2 (Tow, 2023) in various task domains, despite 7B models using more than 3x the training compute and being trained on 1.6x more data.

| Model | | Pre-training ($\downarrow$) | | Downstream task accuracy ($\uparrow$) | | | | |
|---|---|---|---|---|---|---|---|---|
| | | Tokens | FLOPs | CSR | WK | RC | MMLU | Code |
| StableLM-Alpha-3B-v2 | 2.7B | 1.1T | 2.10e22 | 58.0 | 31.7 | 48.1 | 26.6 | 9.7 |
| RedPajama-INCITE-3B | **2.6B** | 800B | 1.50e22 | 56.7 | 34.6 | 48.4 | 27.0 | 5.0 |
| OpenLLaMA 3Bv2 | 3.3B | 1T | 2.20e22 | 57.7 | 33.7 | 47.7 | 26.6 | 9.5 |
| BTLM-3B-8K | **2.6B** | **627B** | **1.3e22** | **59.9** | **36.6** | **50.0** | **28.1** | **9.9** |
| StableLM-Alpha-7B-v2 | 6.7B | 1.1T | 4.90e22 | 61.2 | 38.3 | 48.1 | 26.6 | 15.0 |
| Falcon-7B | 6.9B | 1.5T | 7.00e22 | 63.4 | 45.0 | 51.1 | 26.3 | 0.0 |
| RedPajama-INCITE-7B | 6.7B | 1T | 4.40e22 | 59.5 | 40.1 | 50 | 27.5 | 5.2 |
| Falcon-RW-7B | **6.3B** | **350B** | **1.5e22** | 61.0 | 39.1 | 49.8 | 26.2 | N/A |
| OpenLLaMA 7B | 6.6B | 1T | 4.30e22 | 58.6 | 41.7 | 50.2 | 30.1 | 7.7 |
| MPT-7B | 6.7B | 1T | 4.40e22 | 63.2 | 42.7 | 50.7 | 28.5 | **15.4** |
| XGen-7B-8K | 6.7B | 1.5T | 7.10e22 | 60.7 | 40.0 | 51.5 | 35.9 | 14.2 |
| OpenLLaMA 7Bv2 | 6.6B | 1T | 4.30e22 | 60.5 | 40.7 | 50.7 | 40.4 | 14.7 |
| LLaMA-7B | 6.6B | 1T | 4.30e22 | **63.7** | 45.3 | 52.1 | 35.2 | 12.1 |
| LLaMA-2-7B | 6.6B | 2T | 9.30e22 | 63.4 | **47.5** | **53.2** | 45.8 | 13.7 |

Table 2: Average accuracy on common sense reasoning (CSR), world knowledge (WK), reading comprehension (RC), massive multitask language understanding (MMLU), and code tasks. All tasks are using 0-shot evaluation, except MMLU which is 5-shot. Code accuracy refers to HumanEval pass@1 accuracy.

In addition, we also evaluate model performance on long-context tasks in Section 3.6. BTLM-3B-8K outperforms MPT-7B-8K (Team, 2023a) and XGen-7B-8K Nijkamp et al. (2023) in QMSum and GovReports, two 8,192 context length summarization tasks (Zhong et al., 2021; Huang et al., 2021). In long range retrieval tasks, BTLM-3B-8K performs comparably to MPT-7B-8K and outperforms XGen-7B-8K. BTLM-3B-8K matches or outperforms these 7B models using 5.6x less pretraining FLOPs and 2.4x less pretraining data.

We attribute BTLM's competitive performance to the high quality SlimPajama dataset (Soboleva et al., 2023) and the training improvements described in Section 4. We provide detailed results within each category in subsequent sections. Detailed task descriptions and further evaluation results are available in Appendix A, C.

## 3.1   Common Sense Reasoning

To evaluate common sense reasoning (CSR) capability, we report zero-shot results on the following tasks: PIQA (Bisk et al., 2020), SIQA (Sap et al., 2019), HellaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2021), and OpenBookQA (OBQA) (Mihaylov et al., 2018). These tasks involve multiple choice questions that test general reasoning, understanding of the physical world, emotional and social intelligence, and skills in pronoun resolution problems.

Table 3 shows that BTLM-3B outperforms all other 3B models on common sense reasoning tasks by a significant margin. In addition, BTLM-3B achieves a higher average accuracy on common sense reasoning tasks than OpenLLaMA-7B and RedPajama-INCITE-7B while using far less compute.

## 3.2   Reading Comprehension

We measure reading comprehension (RC) abilities with zero-shot evaluation on RACE-middle (R-m), RACE-high (R-h) (Lai et al., 2017), and BoolQ (Clark et al., 2019). The RACE dataset is sourced from English exams in China for middle and high school students. The RACE questions are written by human experts and test word matching, paraphrasing, and reasoning. BoolQ involves answering yes or no questions about passages from Wikipedia articles.

| Model | Common Sense Reasoning (↑) | | | | | |
| | PIQA | SIQA | HellaSwag | WinoGrande | OBQA | Avg. |
|---|---|---|---|---|---|---|
| RedPajama-INCITE-Base-3B-v1 | 73.8 | 44.9 | 63.2 | 63.6 | 37.8 | 56.7 |
| OpenLLaMA 3Bv2 | 76.2 | 44.8 | 65.2 | 63.3 | 39.2 | 57.7 |
| StableLM-Base-Alpha-3B-v2 | **77.2** | 44.1 | 65.8 | 62.3 | **40.8** | 58.0 |
| BTLM-3B-8K | **77.2** | **46.5** | **69.8** | **65.8** | 40.4 | **59.9** |
| OpenLLaMA 7B | 74.5 | 46.9 | 64.7 | 66.8 | 40.0 | 58.6 |
| RedPajama-INCITE-7B-Base | 77.4 | 45.1 | 70.4 | 64.0 | 40.4 | 59.5 |
| OpenLLaMA 7Bv2 | 78.2 | 47.0 | 69.6 | 65.8 | 42.0 | 60.5 |
| XGen-7B-8K-Base | 75.9 | 47.9 | 74.2 | 65.5 | 40.2 | 60.7 |
| Falcon-RW-7B | 79.1 | 46.6 | 72.1 | 65.7 | 41.4 | 61.0 |
| StableLM-Base-Alpha-7B-v2 | 79.8 | 44.1 | 71.7 | 69.1 | 41.2 | 61.2 |
| MPT-7B | **80.6** | 48.1 | 76.2 | 68.1 | 42.8 | 63.2 |
| Falcon-7B | 80.5 | **49.1** | **76.3** | 67.1 | 44. | 63.4 |
| LLaMA-2-7B | 79.0 | 49.0 | 76.0 | 68.9 | 44.2 | 63.4 |
| LLaMA-7B | 79.2 | 48.5 | 76.2 | **70.0** | **44.4** | **63.7** |

Table 3: Zero-shot validation accuracy on each common sense reasoning task, except for OpenBookQA which uses the test split.

| Model | Reading Comprehension (↑) | | | | World Knowledge (↑) | | | | |
| | R-m | R-h | BoolQ | Avg. | ARC-e | ARC-c | NQ | TQA | Avg. |
|---|---|---|---|---|---|---|---|---|---|
| StableLM-Alpha-3B-v2 | 41.2 | 38.9 | 64.3 | 48.1 | 53.8 | 32.9 | 5.5 | 34.5 | 31.7 |
| OpenLLaMA 3Bv2 | 40.6 | 36.8 | 65.6 | 47.7 | 61.9 | 35.1 | 6.3 | 31.5 | 33.7 |
| RedPajama-INCITE-3B | 40.1 | 37.9 | 67.4 | 48.5 | 61.6 | 34.4 | 6.4 | **36.0** | 34.6 |
| BTLM-3B-8K | **40.6** | **39.4** | **70.0** | **50.0** | **66.9** | **37.6** | **6.9** | 34.9 | **36.6** |
| StableLM-Alpha-7B-v2 | 42.3 | 38.8 | 70.2 | 50.4 | 59.4 | 38.1 | 9.1 | 46.5 | 38.3 |
| Falcon-RW-7B | 41.7 | 38.6 | 69.1 | 49.8 | 67.9 | 38.7 | 9.8 | 39.9 | 39.1 |
| RedPajama-INCITE-7B | 41.2 | 38.2 | 70.8 | 50.1 | 69.3 | 39.2 | 5.5 | 46.2 | 40.1 |
| OpenLLaMA 7B | 42.3 | 37.7 | 70.5 | 50.2 | 67.1 | 37.1 | 12.2 | 50.3 | 41.7 |
| MPT-7B | 40.3 | 38.0 | 73.7 | 50.7 | 70.0 | 41.9 | 11.9 | 47.1 | 42.7 |
| OpenLLaMA 7Bv2 | 41.2 | 38.7 | 72.3 | 50.7 | 68.0 | 40.2 | 7.9 | 46.9 | 40.7 |
| Falcon-7B | 42.3 | 37.2 | 73.8 | 51.1 | 70.8 | 43.5 | **14.6** | 50.9 | 45.0 |
| XGen-7B-8K | 41.2 | 39.0 | 74.2 | 51.5 | 66.9 | 41.1 | 07.2 | 44.6 | 40.0 |
| LLaMA-7B | 40.9 | **40.3** | 75.0 | 52.1 | 72.9 | 44.7 | 11.7 | 52.1 | 45.3 |
| LLaMA-2-7B | **42.3** | 39.5 | **77.8** | **53.2** | **74.6** | **46.3** | 12.5 | **56.6** | **47.5** |

Table 4: Zero-shot accuracy on reading comprehension and world knowledge tasks. We report test accuracy except for BoolQ, where we report validation accuracy.

Table 4 shows BTLM-3B-8K achieves a significantly higher average reading comprehension accuracy than other 3B models and Falcon-RW-7B. RACE-middle is an exception where StableLM-Alpha-v2-3B surpasses BTLM. On the RACE-high task, BTLM-3B outperforms all 3B and 7B models except for LLaMA-7B and LLaMA-2-7B.

## 3.3 World Knowledge

To assess the depth of knowledge acquired by models during training and their proficiency in recalling it upon prompting, we use four closed-book question answering tasks: ARC-easy (ARC-e), ARC-challenge (ARC-c), NaturalQuestions (NQ), and TriviaQA (TQA) (Clark et al., 2018; Kwiatkowski et al., 2019; Joshi et al., 2017). In these tasks, models are presented questions and do not have access to documents containing evidence for the answer. ARC contains multiple choice grade school science questions, NaturalQuestions

contains short questions from Google search engine users, and TriviaQA contains questions from trivia quiz-league websites.

In Table 4 we show that BTLM-3B achieves the highest average accuracy on world knowledge tasks amongst 3B models. In contrast to all other task types, BTLM-3B underperforms every 7B model in average world knowledge accuracy. We hypothesize this is because world knowledge tasks evaluate what knowledge has been compressed into model parameters, which puts smaller models at a disadvantage. BTLM-3B performs comparably to 7B models in other task types where questions are presented in an open-book manner, testing language understanding. This interpretation suggests that smaller models are better suited to tasks where plenty of context is provided.

### 3.4 Massive Multitask Language Understanding

To evaluate models' performance on multiple choice questions from 57 subjects, spanning STEM to humanities, we measure performance on the massive multilingual language understanding (MMLU) benchmark (Hendrycks et al., 2020). This collection of tasks mirrors human evaluations, making it more challenging. The difficulty varies from elementary to professional levels while examining both general knowledge and problem-solving skills. Following Touvron et al. (2023a) we report 5-shot performance on humanities (Hum.), STEM, social sciences (Soc. Sci.), and "Other" task categories, as well as the overall average in Table 5. BTLM-3B not only performs better than all the 3B models but also outperforms Falcon-7B, Falcon-RW-7B, and RedPajama-INCITE-7B.

| Model | MMLU (↑) | | | | | Code (↑) | |
|---|---|---|---|---|---|---|---|
| | Hum. | STEM | Soc. Sci. | Other | Avg. | HE@1 | HE@100 |
| StableLM-Alpha-3B-v2 | 27.1 | 26.2 | 24.9 | 28.2 | 26.6 | 9.7 | **33.3** |
| OpenLLaMA 3Bv2 | 25.7 | 26.0 | 26.6 | 28.5 | 26.7 | 9.5 | 32.9 |
| RedPajama-INCITE-3B | 26.2 | 26.6 | 29.6 | 25.9 | 27.1 | 5.0 | 13.3 |
| BTLM-3B-8K | **27.6** | **27.1** | **27.9** | **29.8** | **28.1** | **9.9** | 29.7 |
| Falcon-RW-7B | 27.3 | 23.2 | 25.6 | 27.7 | 26.0 | N/A | N/A |
| Falcon-7B | 26.9 | 25.9 | 24.4 | 27.6 | 26.2 | 0.0 | 1.8 |
| RedPajama-INCITE-7B | 26.2 | 27.4 | 30.6 | 26.4 | 27.7 | 5.2 | 19.2 |
| MPT-7B | 27.4 | 28.1 | 29.2 | 29.7 | 28.6 | **15.4** | **54.2** |
| OpenLLaMA 7B | 28.4 | 28.4 | 31.3 | 32.9 | 30.3 | 7.7 | 24.9 |
| LLaMA-7B | 34.0 | 30.6 | 38.4 | 38.2 | 35.3 | 12.1 | 35.9 |
| XGen-7B-8K | 33.6 | 29.8 | 39.5 | 41.6 | 36.1 | 14.2 | 41.5 |
| OpenLLaMA 7Bv2 | 37.0 | 33.4 | 45.4 | 47.0 | 40.7 | 14.7 | 47.3 |
| StableLM-Alpha-7B-v2 | 42.6 | 36.6 | 49.3 | 51.2 | 44.9 | 15.0 | 44.9 |
| LLaMA-2-7B | **43.1** | **36.9** | **51.7** | **52.6** | **46.1** | 13.7 | 43.6 |

Table 5: Five-shot accuracy on the Massive Multitask Language Understanding (MMLU) benchmark and zero-shot performance on HumanEval (HE) with pass@1 and pass@100 on the test splits.

### 3.5 Code

To evaluate BTLM-3B-8K' coding ability, we use the HumanEval (HE) (Chen et al., 2021) task. In this task, models are presented with a concise program description, a function signature, and several valid input-output test cases. The objective is to generate a Python program that satisfies the test cases and program description. We adhered to the original HumanEval task Chen et al. (2021) settings (pass@1 uses 0.2 temperature sampling and for pass@100 we use 0.8).

Table 5 shows BTLM-3B-8K outperforms all 3B models, as well as Falcon-7B, RedPajama-INCITE-7B, and OpenLLaMA-7B. Performance on coding benchmarks is correlated with the amount of code present in a model's pretraining data. For example, MPT-7B contains the largest proportion of code related tokens (13.5%) and tops the 7B category. For BTLM, 5% of the training tokens were code related. We exclude Falcon-RW-7B since it was not trained on any code data.

### 3.6    Long Context Evaluation

The ability to perform inference on long context lengths is essential for applications such as document summarization, document question answering, content generation, or supporting a long chat history. In this section we evaluate the long context capability of BTLM-3B-8K against MPT-7B-8K (Team, 2023a) and XGen-7B-8K (Nijkamp et al., 2023), two foundation models trained to perform long context inference. In Section 3.6.1, we evaluate interpolation up to 8,192 context lengths using QMSum (Zhong et al., 2021) and GovReports (Huang et al., 2021), two text summarization tasks. Then, in Section 3.6.2, we evaluate extrapolation to longer contexts than were seen during training with the LongEval tasks (Li* et al., 2023). Finally, we more thoroughly examine BTLM-3B-8K's extrapolation capability in Section 3.6.3.

#### 3.6.1    Long Context Interpolation

| Model | | Pretraining (↓) | | QMSum (↑) | | | GovReports (↑) | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Tokens | FLOPs | R-1 | R-2 | R-L | R-1 | R-2 | R-L |
| XGen-7B-8K | 6.7B | 1.5T | 7.0e22 | 11.8 | 3.0 | 9.1 | 11.8 | 5.6 | 8.3 |
| MPT-7B-8K | 6.7B | 1.5T | 7.1e22 | 14.8 | **5.2** | 11.3 | 8.5 | 3.9 | 6.2 |
| BTLM-3B-8K | **2.7B** | **627B** | **1.3e22** | **16.3** | 2.5 | **12.4** | **15.5** | **5.8** | **10.2** |

Table 6: ROUGE scores on the QMSum and GovReports long text summarization tasks. To test the interpolation regime for models, we only evaluate samples less than 8,192 tokens in length.

Table 6 reveals that BTLM-3B-8K surpasses both MPT-7B-8K and XGen-7B-8K on QMSum and GovReports tasks. Notably, it achieves this using only 40% of the parameters, 41.7% of the pretraining tokens, and 17.9% of the pretraining FLOPs compared to the other models. Notably, MPT-7B-8K achieves a greater ROUGE-2 score for QMSum while BTLM-3B-8K achieves higher ROUGE-1 and ROUGE-L scores.

#### 3.6.2    Long Context Extrapolation

To measure extrapolation to longer context lengths than were seen during training, we perform evaluation on the two tasks from the LongEval benchmark Li* et al. (2023). The "Coarse-grained Topic Retrieval" task, which we abbreviate to "LongEval-Topics", requires models to retrieve the first discussed topic from a long conversation that spans multiple topics. The "Fine-grained Line Retrieval" task which we abbreviate to "LongEval-Lines", requires models to precisely retrieve a number from a long document. With our tokenizer, LongEval-Topics and LongEval-Lines contain examples up to 14.2K and 12.1K context length respectively. We present results in terms of number of topics or lines to remain agnostic to tokenizer differences.
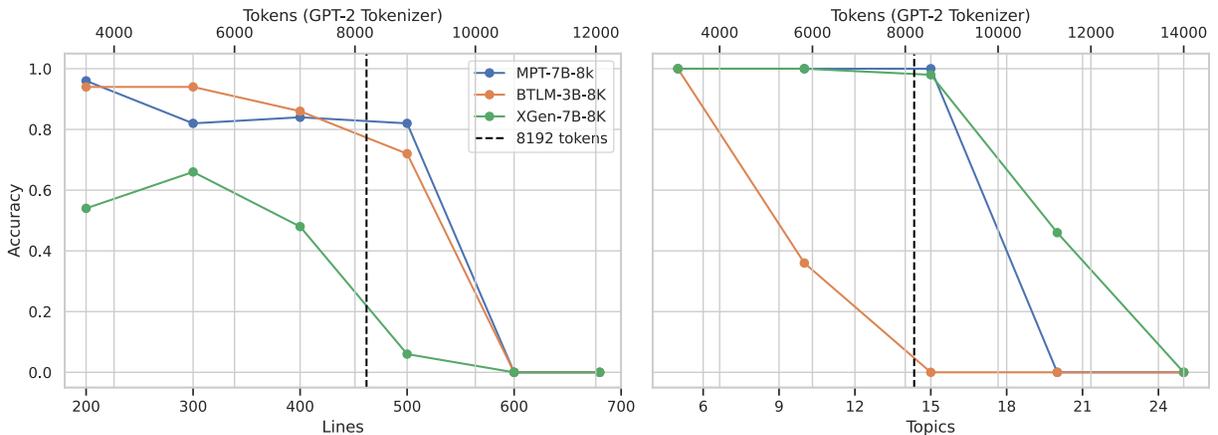


Figure 2: Accuracy on the LongEval-Lines and LongEval-Topics long-range retrieval tasks.

Figure 2 shows BTLM-3B-8K and MPT-7B-8K significantly outperform XGen-7B-8K on both LongEval tasks. This is because both use ALiBi position embeddings while XGen-7B-8K uses rotary position embeddings which do not extrapolate well without additional techniques (Chen et al., 2023; Pal et al., 2023). BTLM-3B-8K is comparable to MPT-7B-8K on LongEval-Lines but MPT-7B-8K extrapolates to slightly longer context lengths on LongEval-Topics, which we believe happens due to MPT model trained on 3.2x more tokens with 8,192 context length.

### 3.6.3 BTLM-3B-8K SlimPajama Extrapolation

To further assess BTLM's extrapolation capability, we evaluate on the SlimPajama test set with 32768 context length and plot loss at each token position in Figure 3. We evaluate checkpoints at different points during training to gain insight into how extrapolation capability evolves.
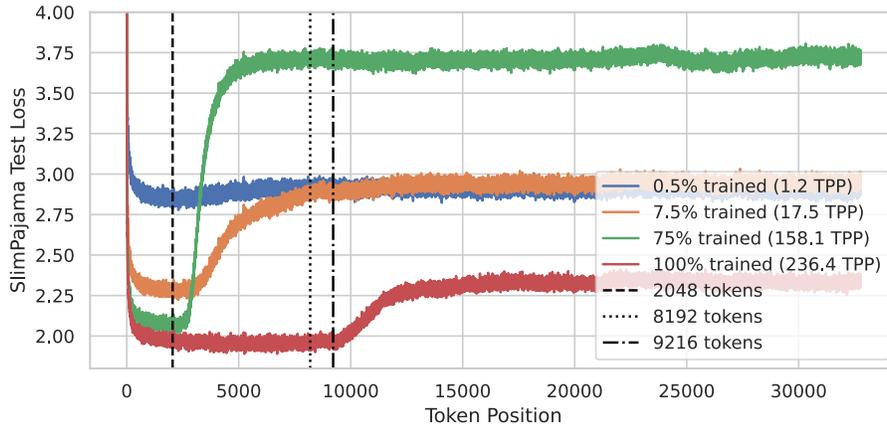


Figure 3: SlimPajama test set cross-entropy loss for various BTLM checkpoints at each token position. Inference is performed on examples packed to 32768 tokens in length.

Press et al. (2021) report that ALiBi grants impressive extrapolation properties with a 255M parameter model trained on 103M tokens. This corresponds to just 0.4 tokens per parameter (TPP), well below the 20 TPP recommendation from Hoffmann et al. (2022). With our 1.2 TPP checkpoint we observe similar extrapolation performance as Press et al. (2021) but this result appears to only be possible due the overall loss being quite poor quite early in training. As training progresses, our model learns to "overfit" to the current context length. We observe that the final checkpoint from the 2,048 context length training phase (75% complete) cannot extrapolate well beyond 2,048 context length. This demonstrates that ALiBi alone does not provide competitive extrapolation capability, and we suggest using variable context length training schedules to improve performance. The final BTLM-3B-8K model trained on 8,192 with a context length can extrapolate well up to ≈9,216 context length but suffers loss degradation beyond this.

### 3.7 Bias, Toxicity, and Truthfulness

Language models have been found to inherit biases present in their training data (Sheng et al., 2019; Kurita et al., 2019) and implicated in generating offensive and toxic content (Gehman et al., 2020). Therefore to quantify the potential harm BTLM could cause in deployment settings without additional mitigations, we compare the bias, toxicity, and truthfulness of BTLM with with OpenLLaMA-3B-v2 (Geng & Liu, 2023), RedPajama-INCITE-7B (Together.ai, 2023), Falcon-7B (Almazrouei et al., 2023) and LLaMA-2-7B (Touvron et al., 2023b) in Table 7.

The TruthfulQA task evaluates how well models can distinguish factually correct statements from incorrect ones (Lin et al., 2021). BTLM produces more reliable outputs than all tested models except for LLaMA-2-7B.

The WinoGender task presents models with sentences containing two subjects and a pronoun that requires models to correctly guess which subject the pronoun refers to (Rudinger et al., 2018). Subjects are people

who are referred to by their occupation, for example "the paramedic". "Gotcha" examples contain sentences where the pronoun gender does not match the occupation's majority gender based on the US Bureau of Labor Statistics. When we compare WinoGender accuracy for specific pronoun categories, we primarily assess a model's capability in common-sense reasoning. To evaluate bias, we look into the difference in accuracy that a model achieves on different pronoun categories. We observe BTLM is better at resolving gendered pronouns than gender neutral pronouns, indicating bias. BTLM also performs worse than random on the gotcha categories, indicating the model has internalized the gender biases associated with occupations.

In the ToxiGen task, models are asked to classify a sentence mentioning minority groups as toxic or non-toxic (Hartvigsen et al., 2022). There appears to be an inverse correlation between overall model performance and the probability of producing toxic outputs. BTLM model produces more toxic content than OpenLLaMA-3B-v2 and RedPajama-INCITE-7B, but less than Falcon-7B and LLaMA-2-7B.

The CrowS-Pairs task evaluates the bias of models on 9 different categories (Nangia et al., 2020). BTLM's bias in this task like RedPajama-INCITE-7B and Falcon-7B models which achieve comparable performance across a range of downstream tasks.

Overall BTLM exhibits bias, toxicity, and truthfulness like existing models. Nevertheless, we recommend exploring harm mitigation strategies in deployment contexts (OpenAI, 2023). Additionally, more careful dataset curation techniques such as filtering not-safe-for-work URLs (Penedo et al., 2023) showed to be helpful in reducing model harmfulness.

| Task | Subtask | BTLM-3B-8K | OpenLLaMA 3Bv2 | RedPajama-INCITE-7B | Falcon-7B | LLaMA-2-7B |
|---|---|---|---|---|---|---|
| TruthfulQA ↑ | Multiple choice | **35.9** | 34.8 | 33.0 | 34.2 | **39.0** |
| WinoGender ↑ | hers/her/she | **60.0** | 56.7 | 63.3 | 60.0 | **69.2** |
| | his/him/he | **60.0** | 56.7 | 60.0 | 55.0 | **62.5** |
| | their/them/someone | 57.5 | **60.0** | **72.5** | 56.7 | 69.2 |
| | hers/her/she (gotcha) | **48.3** | 37.9 | 48.3 | 41.4 | **62.1** |
| | his/him/he (gotcha) | 29.0 | **35.5** | 51.6 | 45.2 | **67.7** |
| | All | **59.2** | 57.8 | 65.3 | 57.2 | **66.9** |
| ToxiGen ↓ | Multiple choice | 50.7 | **44.6** | 45.3 | 52.7 | 57.8 |
| CrowS-Pairs ↓ | Age | 75.8 | **53.9** | **71.4** | **71.4** | 74.7 |
| | Disability | 69.2 | **64.6** | 76.9 | **67.7** | **67.7** |
| | Gender | 67.2 | **53.8** | 68.4 | 66.9 | **62.5** |
| | Nationality | 60.2 | **52.3** | 62.5 | 61.1 | **59.7** |
| | Physical Appearance | 77.8 | **66.7** | 79.2 | 76.4 | 77.8 |
| | Race/Color | 54.1 | **49.6** | 59.7 | 56.7 | 61.6 |
| | Religion | 74.8 | **71.2** | 76.6 | 73.9 | 81.1 |
| | Sexual Orientation | 86.0 | **69.9** | 88.2 | 86.0 | **78.5** |
| | Socioeconomic Status | 69.0 | **59.5** | **69.5** | **69.5** | 74.2 |
| | Average | 65.1 | **56.0** | **65.6** | 67.8 | 66.9 |

Table 7: Zero-shot evaluations on bias, toxicity, and truthfulness benchmarks: TruthfulQA, WinoGender, ToxiGen, and CrowS-Pairs.

# 4 Training Improvement Ablations

To arrive at the final training setup for BTLM, we test various architectural modifications and training techniques. In this section, we present an ablation study for each training improvement starting from a GPT-3-style training setup. By combining all the changes, we improve pretraining loss by 5.36% over the baseline training setup. As we show in Section 3, this combination of features results in BTLM outperforming all other 3B parameter foundation models and even surpassing some 7B models.

### 4.1 Baseline Training Setup

We begin from a GPT-3 style (Brown et al., 2020a) autoregressive transformer decoder model used in the Cerebras-GPT $\mu$P models (Dey et al., 2023; Yang et al., 2021). We train models with 20 tokens per parameter (TPP) (Hoffmann et al., 2022), the GELU activation function, linear learning rate decay to 10% of the maximum, learned position embeddings, and the following $\mu$P tuned hyperparameters:

- Base learning rate = 6e-3

- Base weight initialization standard deviation = 0.08

- Embedding multiplier = 10

- Output logit multiplier = 1

We use a 111M parameter model size with $d_{model}$=768, $n_{layers}$=10, and $d_{head}$=64 to perform each ablation. As Section 3 shows, the insights we gain from the 111M size transfer well to the final 2.6B size. Our ablations were performed with the Pile dataset (Gao et al., 2020), but we find our results generalize well to the SlimPajama dataset used in the final training run. Unless otherwise specified, models are trained with a 2,048 context length.

### 4.2 Architecture and Training Hyperparameter Improvements

First, we perform several ablations with different model architectures and training hyperparameters, then, we measure how they affect training efficiency. These ablations are summarized in Table 8 and Figure 4. We ablate each improvement in an additive fashion and measure the effect relative to the baseline described in Section 4.1. In Figure 4, we fit a scaling law to 111M, and 256M parameter baseline models trained with 236.4 TPP. This allows us to estimate the performance of the baseline setup at every FLOP budget. We use this property to estimate the iso-FLOP loss improvement, iso-loss FLOP reduction and iso-loss parameter reduction that each variant achieves over the 236.4 TPP scaling law. Through all the improvements presented in this section, we decrease loss by 5.36% relative to the 236.4 TPP or achieve the same loss using 35% of the FLOPs.
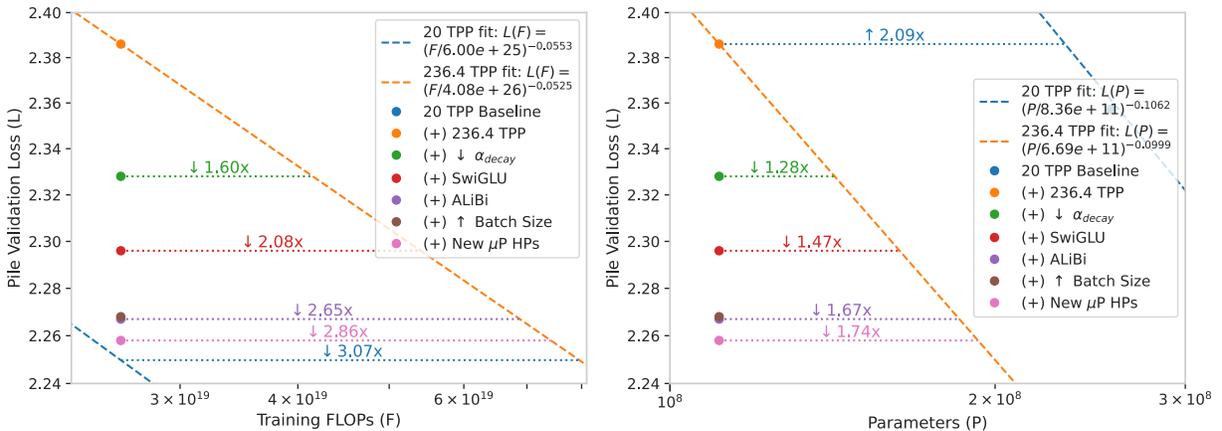


Figure 4: Overview of each architecture and training hyperparameter improvement ablated starting from a CerebrasGPT-$\mu$P baseline (Dey et al., 2023). Power law fits are included for 20 TPP and 236.4 TPP baselines. Relative to these power laws we illustrate the FLOP and parameter differences at the same loss.

#### 4.2.1 Increased Tokens per Parameter

BTLM-3B-8K is a 2.6B parameter model trained for 627B tokens or 236.4 tokens per parameter (TPP). Starting from a 111M parameter compute-optimal 20 TPP baseline (Hoffmann et al., 2022) described in

| Variant | TPP | $\alpha_{decay}$ | Activation Function | Position Embed. | Batch Size | $\mu$P HPs | Pile Valid. Loss | Iso-FLOP $\Delta$ Loss | Iso-Loss $\Delta$ FLOP | Iso-Loss $\Delta$ Param. |
|---|---|---|---|---|---|---|---|---|---|---|
| (+) 20 TPP | 20 | 0.1 | GeLU | Learned | N/A | Old | 2.247* | -5.82% | ↓3.07x | ↑2.09x |
| 236.4 TPP Baseline | **236.4** | 0.1 | GeLU | Learned | 120 | Old | 2.386 | 0% | 1x | 1x |
| (+) ↓ $\alpha_{decay}$ | **236.4** | **0.0085** | GeLU | Learned | 120 | Old | 2.328 | -2.43% | ↓1.60x | ↓1.28x |
| (+) SwiGLU | **236.4** | **0.0085** | **SwiGLU** | Learned | 120 | Old | 2.296 | -3.77% | ↓2.08x | ↓1.47x |
| (+) RoPE | **236.4** | **0.0085** | **SwiGLU** | RoPE | 120 | Old | 2.259 | -5.32% | ↓2.84x | ↓1.73x |
| (+) ALiBi | **236.4** | **0.0085** | **SwiGLU** | **ALiBi** | 120 | Old | 2.267 | -4.99% | ↓2.65x | ↓1.67x |
| (+) ↑ Batch Size | **236.4** | **0.0085** | **SwiGLU** | **ALiBi** | **420** | Old | 2.268 | -4.95% | ↓2.63x | ↓1.66x |
| (+) New $\mu$P HPs | **236.4** | **0.0085** | **SwiGLU** | **ALiBi** | **420** | **New** | **2.258** | **-5.36%** | ↓**2.86x** | ↓**1.74x** |

Table 8: Ablation of different training configurations. Settings used in the final BTLM setup are bolded. (*) Projected based on 20 TPP power law at 236.4 TPP FLOP budget.

Section 4.1, we increase TPP to 236.4 to more closely mimic BTLM's training conditions. Due to the training inefficiency introduced by this regime, Table 8 shows the 20 TPP setup achieves 5.82% lower loss than the 236.4 TPP baseline with the same compute budget. In other words, 236.4 TPP training requires 3.07x more compute to reach the same loss as a 20 TPP model. However, the 20 TPP setup requires 2.09x more parameter to reach the same loss as the 236.4 TPP baseline, demonstrating the inference benefit of over-training. This 236.4 TPP model serves as a baseline for the ablations that follow.

### 4.2.2 Increased Learning Rate Decay Ratio

For LLM training, it is most common to include a short learning rate warmup followed by a cosine or linear decay to 10% of the maximum learning rate. Hoffmann et al. (2022) found this decay to 10% of the maximum learning rate to be optimal for the 20 TPP setting. We hypothesize that in higher TPP ($\tau$) settings the learning rate decay fraction ($\alpha_{decay}$) should be increased to encourage finer grained weight updates later in training. Equation 1 proposes a simple heuristic: in higher TPP settings increase $\alpha_{decay}$ proportional to the $\alpha_{decay} = 0.1, \text{TPP} = 20$ setting.

$$\alpha_{decay} = 0.1 \cdot (20/\text{TPP}) \tag{1}$$

In Figure E we sweep $\alpha_{decay}$ for 370 TPP and find this rule of thumb to provide good prediction of $\alpha_{decay}$. For 236.4 TPP, Equation 1 suggests decaying to 0.85% of the maximum learning rate. Table 8 shows $\alpha_{decay} = 0.0085$ decreases loss by 2.43% relative to $r_{decay} = 0.1$ or requires 1.60x less FLOPs to achieve the same loss.

### 4.2.3 SwiGLU Activation Function

Shazeer (2020) showed that activation functions with gated linear units (GLU) improve transformer training. Then, Scao et al. (2022) demonstrated the SwiGLU activation function outperforms the GELU activation function. We repeat this ablation and show SwiGLU decreases loss by 1.37% relative to GELU (Table 8). To keep compute comparable to GELU models with $d_{ffn} = 4d_{model}$, we use $d_{ffn} = \frac{8}{3}d_{model}$ to account for the additional projection.

### 4.2.4 ALiBi and RoPE Position Embedding

Scao et al. (2022) showed the Attention with Linear Biases (ALiBi) position embedding Press et al. (2021) outperforms both learned and rotary position embeddings (RoPE) (Su et al., 2022). In Table 8 we observe the opposite: RoPE outperforms ALiBi at 2,048 context length training. Despite this we selected ALiBi for the BTLM model due to the superior extrapolation capability. 8 shows ALiBi decreases loss by 1.26% relative to learned position embeddings.

### 4.2.5 Increased Batch Size and Improved $\mu$P Hyperparameters

The maximal update parameterization ($\mu$P) enables the transfer of optimal hyperparameters (HPs) from a small proxy model up to a very large target model Yang et al. (2021). However, we should not ignore the effect of batch size on the optimal learning rate. If the proxy model is trained with a batch size smaller than the critical batch size McCandlish et al. (2018), learning rate transfer to a large model trained at or above the critical batch size will be sub-optimal. We perform a random search on a 40M parameter proxy model, ensuring a large enough batch size, and arrive at the following hyperparameters:

- Base learning rate = 1.2e-2

- Base weight initialization standard deviation = 0.073

- Embedding multiplier = 14.6

- Output logit multiplier = 2.22

With a 111M parameter model, we show increasing batch size from 120 to 420 has a negligible effect on the loss. Then using batch size 420, we transfer the optimal hyperparameters from our 40M parameter proxy model and show a 5.36% loss decrease or achieve the same loss with 2.86x fewer FLOPs relative to the 236.4 TPP baseline (Figure 4).

### 4.3 Variable Context Length Training

In this section, our goal is to find an efficient process for training a model which can perform high quality inference up to at least 8,192 context length. The naive approach to achieve this goal is to train a model entirely on data with 8,192 context length. Purely training this way would result in 1.53x more FLOPs than 2,048 context training. To save compute while still achieving long context performance, Devlin et al. (2019) introduced a simple strategy of training 90% of steps on 128 context length, then the final 10% on 512 context length. We extend this methodology by training a 111M parameter model on 75% of tokens at 2,048 context length followed by 25% of tokens at 8,192 context length using ALiBi position embeddings (Press et al., 2021). We compare this variable context length strategy against pure 2,048 context length and pure 8,192 context length training. To assess long sequence capability, we evaluate on the Pile validation set with 32,768 context length and plot the loss at each token position.

Figure 5 shows that the variable context length strategy achieves comparable long sequence loss to pure 8,192 context length training While using 74% of the FLOPs. Although ALiBi was designed to improve models' extrapolation to sequence lengths longer than seen during training, we observe a clear loss degradation at token positions slightly beyond than the training sequence length. The long sequence performance of pure 2,048 context length training shows ALiBi alone is not a sufficient substitute for long context training.

## 5 Related Work

**Parameter-efficient "over-trained" language models**. Touvron et al. (2023a) made a landmark contribution to open-source LLMs by releasing the weights of the LLaMA models. The LLaMA models are trained for many more tokens than would be optimal for training compute efficiency (Hoffmann et al., 2022) but this is done in service of inference-time efficiency. Authors did not release their training dataset however, prompting several groups to reproduce and extend the LLaMA training methodology. Some of these works include RedPajama-INCITE (Together.ai, 2023), OpenLLaMA (Geng & Liu, 2023), StableLM-v2 (Tow, 2023), Falcon (Almazrouei et al., 2023), Falcon-RW (Penedo et al., 2023), XGen (Nijkamp et al., 2023), MPT (Team, 2023b;a), LLaMA (Touvron et al., 2023a), and LLaMA-2 (Touvron et al., 2023b). Our work also extends the LLaMA methodology by training BTLM-3B-8K for 627B tokens, many more than the 54B tokens that would be optimal for training compute efficiency (Hoffmann et al., 2022).

**Long context LLMs**. Many LLM use cases involve performing inference on long context windows such as information retrieval, document summarization, or generating long-form content. Two notable avenues for
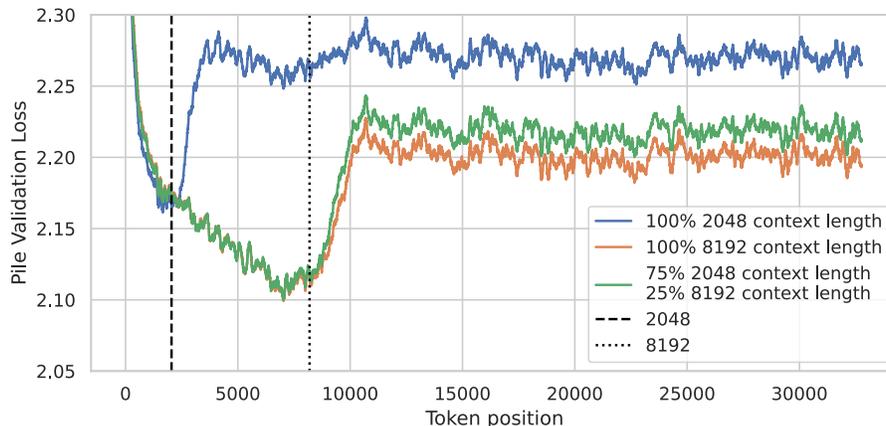
Figure 5: Loss versus token position for various sequence length schedules. Loss is plotted with a 100 value moving average to improve plot readability.

improving the long context performance of LLMs are to train with longer context lengths or use position embeddings designed for extrapolation. One can either train on a long context length for the entire training (Touvron et al., 2023b) or use an increasing sequence length schedule (Devlin et al., 2019) to improve compute efficiency. BTLM-3B-8K along with XGen-7B-8K (Nijkamp et al., 2023) and MPT-7B-8K (Team, 2023a) adopt a variable context length training schedule. In addition to the training sequence length, the choice of position embedding can also affect the long context LLM performance. Rotary position embeddings (RoPE) (Su et al., 2022), attention with linear bias (ALiBi) (Press et al., 2021), and xPos (Sun et al., 2023) were all designed to improve extrapolation to longer context length over baselines such as learned position embeddings (Radford et al., 2019). In this work we adopt ALiBi position embeddings, following the MPT models.

## 6 Conclusion

As LLMs become more ubiquitous, the amount of compute spent on both training and inference is rapidly increasing. In this work we present BTLM-3B-8K, a state-of-the-art 3B parameter language model with performance surpassing even some 7B parameter models while requiring only 40% of the inference compute. With 4-bit quantization our model can run within 3GB of RAM, enabling deployment on billions of mobile devices. BTLM-3B-8K can also perform high quality inference up to 8192 context length, making it suitable for useful applications such as text summarization or document question answering. Finally, the training improvements and extensively deduplicated SlimPajama dataset we present in this work are widely applicable and can significantly improve training, inference, and data efficiency for LLMs. Both the BTLM-3B-8K weights and SlimPajama training dataset are available with a permissible Apache 2.0 license on Hugging Face: `https://huggingface.co/cerebras`.

## Acknowledgements

manuscript provided by Gurpreet Gosal, Gavia Gray, Anshul Samar, William Marshall, and Rob Schreiber. Finally, we acknowledge the contributions of the many Cerebras engineers who made this work possible.

## References

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, et al. Falcon-40B: an open large language model with state-of-the-art performance, 2023.

Yonatan Bisk, Rowan Zellers, Ronan Le bras, Jianfeng Gao, and Yejin Choi. PIQA: Reasoning about Physical Commonsense in Natural Language. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020. URL https://ojs.aaai.org/index.php/AAAI/article/view/6239.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, 2020a.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language Models are Few-Shot Learners, 2020b. URL https://arxiv.org/abs/2005.14165.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating Large Language Models Trained on Code, 2021. URL https://arxiv.org/abs/2107.03374.

Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending Context Window of Large Language Models via Positional Interpolation, 2023. URL https://arxiv.org/abs/2306.15595.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. PaLM: Scaling Language Modeling with Pathways. 2022.

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. BoolQ: Exploring the Surprising Difficulty of Natural Yes/No Questions, 2019. URL https://arxiv.org/abs/1905.10044.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have Solved Question Answering? Try ARC, the AI2 Reasoning Challenge, 2018. URL https://arxiv.org/abs/1803.05457.

Together Computer. RedPajama: An Open Source Recipe to Reproduce LLaMA training dataset, 2023. URL https://github.com/togethercomputer/RedPajama-Data.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT*, 2019.

Nolan Dey, Gurpreet Gosal, Zhiming, Chen, Hemant Khachane, William Marshall, Ribhu Pathria, Marvin Tom, and Joel Hestness. Cerebras-GPT: Open Compute-Optimal Language Models Trained on the Cerebras Wafer-Scale Cluster, 2023. URL https://arxiv.org/abs/2304.03208.

Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. GPTQ: Accurate Post-training Compression for Generative Pretrained Transformers, 2022. URL https://arxiv.org/abs/2210.17323.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The Pile: An 800GB Dataset of Diverse Text for Language Modeling, 2020. URL https://arxiv.org/abs/2101.00027.

Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, et al. A framework for few-shot language model evaluation, 2021. URL https://doi.org/10.5281/zenodo.5371628.

Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. Realtoxicityprompts: Evaluating neural toxic degeneration in language models, 2020.

Xinyang Geng and Hao Liu. Openllama: An open reproduction of llama, 2023. URL `https://github.com/openlm-research/open_llama`.

Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. Toxi-Gen: A Large-Scale Machine-Generated Dataset for Adversarial and Implicit Hate Speech Detection. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 2022. URL `https://aclanthology.org/2022.acl-long.234`.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring Massive Multitask Language Understanding, 2020. URL `https://arxiv.org/abs/2009.03300`.

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. An Empirical Analysis of Compute-optimal Large Language Model Training. In *The Conference on Neural Information Processing Systems (NeurIPS)*, 2022.

Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng Ji, and Lu Wang. Efficient Attentions for Long Document Summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021. URL `https://aclanthology.org/2021.naacl-main.112`.

Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension, 2017. URL `https://arxiv.org/abs/1705.03551`.

Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. Measuring bias in contextualized word representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, 2019. URL `https://aclanthology.org/W19-3823`.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural Questions: A Benchmark for Question Answering Research. *Transactions of the Association for Computational Linguistics*, 2019. URL `https://aclanthology.org/Q19-1026`.

Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. RACE: Large-scale ReAding Comprehension Dataset From Examinations, 2017. URL `https://arxiv.org/abs/1704.04683`.

Jure Leskovec, Anand Rajaraman, and Jeffrey D. Ullman. *Mining of Massive Datasets*. Cambridge University Press, 2014.

Dacheng Li*, Rulin Shao*, Anze Xie, Ying Sheng, Lianmin Zheng, Joseph E. Gonzalez, Ion Stoica, Xuezhe Ma, and Hao Zhang. How Long Can Open-Source LLMs Truly Promise on Context Length?, 2023. URL `https://lmsys.org/blog/2023-06-29-longchat`.

Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring How Models Mimic Human Falsehoods, 2021. URL `https://arxiv.org/abs/2109.07958`.

Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*, 2017.

Sam McCandlish, Jared Kaplan, Dario Amodei, and OpenAI Dota Team. An Empirical Model of Large-Batch Training, 2018. URL `https://arxiv.org/abs/1812.06162`.

Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a Suit of Armor Conduct Electricity? A New Dataset for Open Book Question Answering, 2018. URL `https://arxiv.org/abs/1809.02789`.

Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model Cards for Model Reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 2019.

Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. CrowS-Pairs: A Challenge Dataset for Measuring Social Biases in Masked Language Models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2020. URL `https://aclanthology.org/2020.emnlp-main.154`.

Erik Nijkamp, Tian Xie, Hiroaki Hayashi, Bo Pang, Congying Xia, Chen Xing, Jesse Vig, Semih Yavuz, Philippe Laban, Ben Krause, et al. Long Sequence Modeling with XGen: A 7B LLM Trained on 8K Input Sequence Length. Salesforce AI Research Blog, 2023. URL `https://blog.salesforceairesearch.com/xgen`.

OpenAI. Gpt-4 technical report, 2023. URL `https://arxiv.org/abs/2303.08774`.

Arka Pal, Deep Karkhanis, Manley Roberts, Samuel Dooley, Arvind Sundararajan, and Siddartha Naidu. Giraffe: Adventures in Expanding Context Lengths in LLMs, 2023. URL `https://arxiv.org/abs/2308.10882`.

Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only, 2023. URL `https://arxiv.org/abs/2306.01116`.

Ofir Press, Noah A. Smith, and Mike Lewis. Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation, 2021. URL `https://arxiv.org/abs/2108.12409`.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language Models are Unsupervised Multitask Learners, 2019. URL `https://d4mucfpksywv.cloudfront.net/better-language-models/language-models.pdf`.

Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. Gender Bias in Coreference Resolution. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2018. URL `https://aclanthology.org/N18-2002`.

Keisuke Sakaguchi, Ronan Bras, Chandra Bhagavatula, and Choi Yejin. WinoGrande: an adversarial winograd schema challenge at scale. *Communications of the ACM*, 64:99–106, 09 2021. doi: 10.1145/3474381.

Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. SocialIQA: Commonsense Reasoning about Social Interactions, 2019. URL `https://arxiv.org/abs/1904.09728`.

Teven Le Scao, Thomas Wang, Daniel Hesslow, Lucile Saulnier, Stas Bekman, M Saiful Bari, Stella Biderman, Hady Elsahar, Niklas Muennighoff, Jason Phang, et al. What Language Model to Train if You Have One Million GPU Hours?, 2022. URL `https://arxiv.org/abs/2210.15424`.

Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural Machine Translation of Rare Words with Subword Units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, 2016. URL `https://aclanthology.org/P16-1162`.

Noam Shazeer. GLU Variants Improve Transformer, 2020. URL `https://arxiv.org/abs/2002.05202`.

Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. The Woman Worked as a Babysitter: On Biases in Language Generation, 2019.

Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama. `https://www.cerebras.net/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama`, 2023. URL `https://huggingface.co/datasets/cerebras/SlimPajama-627B`.

Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. RoFormer: Enhanced Transformer with Rotary Position Embedding, 2022. URL `https://arxiv.org/abs/2104.09864`.

Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav Chaudhary, Xia Song, and Furu Wei. A Length-Extrapolatable Transformer. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, 2023. URL `https://aclanthology.org/2023.acl-long.816`.

MosaicML NLP Team. Announcing MPT-7B-8K: 8K Context Length for Document Understanding, 2023a. URL `https://www.mosaicml.com/blog/long-context-mpt-7b-8k`.

MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially usable llms, 2023b. URL `www.mosaicml.com/blog/mpt-7b`.

Together.ai. Redpajama Models V1, 2023. URL `https://together.ai/blog/redpajama-models-v1`.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and Efficient Foundation Language Models, 2023a. URL `https://arxiv.org/abs/2302.13971`.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open Foundation and Fine-Tuned Chat Models, 2023b. URL `https://arxiv.org/abs/2307.09288`.

Jonathan Tow. StableLM Alpha v2 Models, 2023. URL `https://huggingface.co/stabilityai/stablelm-base-alpha-7b-v2`. Additional model available at: `https://huggingface.co/stabilityai/stablelm-base-alpha-3b-v2`.

Greg Yang, Edward Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tuning Large Neural Networks via Zero-Shot Hyperparameter Transfer. In *Advances in Neural Information Processing Systems*, 2021.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a Machine Really Finish Your Sentence?, 2019. URL `https://arxiv.org/abs/1905.07830`.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models, 2022.

Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir Radev. QMSum: A New Benchmark for Query-based Multi-domain Meeting Summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021. URL `https://aclanthology.org/2021.naacl-main.472`.

# Appendix

## Model Card

Table 9 shows the model card for BTLM-3B-8K following the guide from Mitchell et al. (2019).

## Author Contributions

We would like to acknowledge the contributions of those who helped in preparation of this manuscript.
**Pretraining experimental design:** Nolan Dey, Joel Hestness
**Model training:** Ribhu Pathria, Hemant Khachane, Shaheer Muhammad, Zhiming (Charles) Chen
**Pretraining dataset preparation:** Daria Soboleva*, Faisal Al-Khateeb*
**Downstream task comparisons:** Faisal Al-Khateeb, Daria Soboleva, Bowen Yang, Shaheer Muhammad, Nolan Dey
**Manuscript preparation:** Nolan Dey*, Daria Soboleva*, Joel Hestness
**Project management:** Nolan Dey, Daria Soboleva, Marvin Tom
**Project objectives:** Robert Myers, Jacob Robert Steeves, Natalia Vassilieva
**Supervision:** Joel Hestness

## A Downstream Task Descriptions

We provide a brief description of each of the 22 downstream tasks that we report results for in Section 3.

1. **PIQA** tests a model's common sense reasoning about the physical world by posing a prompt and two potential completions. For example:

   > [**Goal**] Make an outdoor pillow
   > [**Sol1**] Blow into a tin can and tie with rubber band
   > [**Sol2**] Blow into a trash bag and tie with rubber band

   The evaluation setup is multiple-choice based on the probability mass of the solutions.

2. **SIQA** is a dataset for commonsense reasoning about social situations. For example:

   > **Context:** Quinn wanted to help me clean my room up because it was so messy.
   > **Question:** What will Quinn want to do next?
   > **AnswerA:** Eat messy snacks
   > **AnswerB:** Help out a friend
   > **AnswerC:** Pick up the dirty clothes

   Similar to PIQA, the evaluation setup is also multiple-choice.

3. **HellaSwag** is a dataset of multiple-choice questions aimed to test a model's common sense reasoning abilities. For example:

   > **Context:** A woman is outside with a bucket and a dog. The dog is running around trying to avoid a bath. She...
   > **A:** rinses the bucket off with soap and blow dry the dog's head.
   > **B:** uses a hose to keep it from getting soapy.
   > **C:** gets the dog wet, then it runs away again.
   > **D:** gets into a bath tub with the dog.

   The authors of the dataset select examples such that they are difficult for language models while still trivial for humans (with reported greater than 95% accuracy).

Table 9: BTLM-3B-8K Model Card

| | |
|---|---|
| **Release details** | |
| | • **Organization**: Cerebras Systems |
| | • **Model date**: July 2023 |
| | • **Model type**: Autoregressive Transformer Language Model (more details in Section 2) |
| | • **Feedback on the model**: Nolan Dey and Daria Soboleva, {nolan, daria.soboleva}@cerebras.net |
| **Model details** | |
| | • **Model architecture**: BTLM-3B-8K is an autoregressive transformer decoder-only model with 2.6 billion parameters. The architecture is similar to GPT-3 with some changes. More details in Section 2.1. |
| | • **Hidden size**: 2,560 |
| | • **Number of layers**: 32 |
| | • **Head size**: 80 |
| | • **Filter size**: 6,826 |
| | • **Context (sequence) length**: 8,192 |
| | • **Initialization**: Model is initialized using maximal update parameterization (µP) which involves applying scalar multiple to initialization of certain layers. More details in Section 2.1. |
| | • **Release license**: Apache 2.0 |
| **Data Overview** | |
| | • **Training data**: BTLM-3B-8K is trained on the SlimPajama dataset from Cerebras (Soboleva et al., 2023). More details in Section 2.2. |
| | • **Pre-processing**: SlimPajama was pre-processed using public code from the Cerebras Model Zoo. Then, data was tokenized with byte-pair encoding using the GPT-2 vocabulary of size 50,257. |
| | • **Evaluation data**: Upstream (pretraining) evaluations were completed using the SlimPajama validation and test set splits. Downstream evaluations were performed on standardized tests across common-sense reasoning, world knowledge, reading comprehension, massive multitask language understanding, code, and long sequences. More details in Section 3. downstream evaluations were performed using the Eleuther lm-eval-harness (Gao et al., 2021). |
| | • **Motivation**: Evaluation tasks were chosen to closely match related works and cover a broad cross-section of task types. |
| **Usage** | |
| | • **Primary intended uses**: The primary intended use is to further research into large language models. BTLM-3B-8K can be used as a foundation model for NLP, applications, ethics, and alignment research. We release these models with a fully permissive Apache license for the community to use freely. |
| | • **Primary intended users**: Researchers who are working to improve LLMs and practitioners who are looking for reference implementations, training setups, hyperparameters, or pretrained models. |
| | • **Limitations**: BTLM-3B-8K was only trained and evaluated following the approaches described in this paper. |
| | • **Out-of-scope uses**: BTLM-3B-8K was trained on SlimPajama, with primarily English language, and is not recommended for machine translation tasks. BTLM-3B-8K has not been tuned for instruction-following or chat-based use cases. Further safety-related testing and mitigations should be applied before using the model in production downstream applications. |
| **Metrics** | |
| | • **Model performance measures**: Model is evaluated using text prediction cross-entropy on upstream tasks and text generation accuracy on downstream tasks. Results are compared against many publicly available large language models. Details can be found in Section 3. |
| **Ethical considerations** | |
| | • **Data**: SlimPajama is a primarily English corpus and may contain content considered toxic, gender biased, pejorative, racially sensitive, etc. |
| | • **Human life**: The outputs from this model may or may not align with human values. The risk needs to be thoroughly investigated before deploying this model in a production environment where it can directly impact human life. |
| | • **Risks and harms**: There can be distributional bias in the SlimPajama dataset that can manifest in various forms in the downstream model deployment. There are other risks associated with large language models such as amplifying social stereotypes, memorizing training data, or revealing private or secure information. |
| | • **Mitigations**: SlimPajama takes no further mitigation actions beyond those used in the cration of the original RedPajama data set. |
| **Factors** | |
| | • **Evaluation factors**: BTLM-3B-8k was evaluated for various bias factors using TruthfulQA, WinoGender, ToxiGen, and CrowS-Pairs. Details are in Section 3.7. |
| **Implementation infrastructure** | |
| | • **Hardware**: G42's Condor Galaxy-1 AI Supercomputer; the first deliverable of the G42 Cerebras strategic partnership. CG-1 is a 4 exaFLOP AI supercomputer, located in Santa Clara California, and built by G42 and Cerebras. G42's portfolio companies, G42 Cloud and the Inception Institute of Artificial Intelligence (IIAI), generously provided access to CG-1 for the BTLM training effort. |
| | • **Software**: PyTorch, Cerebras Software Platform (CSoft) release 1.9 |

4. **WinoGrande** consists of a set of pronoun resolution problems. Samples are constructed as pairs of similar sentences, each with a pronoun referring to a noun earlier in the sentence. The task is to predict which noun the pronoun refers to. For example, in the sample:

   **a.** The trophy doesn't fit into the brown suitcase because it's too large.

   **b.** The trophy doesn't fit into the brown suitcase because it's too small.

   in sentence (a), "it's" referring to "trophy", while in sentence (b), changing a single context word modifies the meaning of the sentence such that "it's" now refers to "suitcase".

5. **OpenBookQA** is a multiple-choice common-sense question answering dataset (Mihaylov et al., 2018). One example question from this dataset is:

   **What is the most likely to be an effect of acid rain on an aquatic environment?**
   **(A)** increase in plant growth
   **(B)** increase in fish population
   **(C)** decrease in plant life
   **(D)** cleaner and clearer water

6. **RACE-middle** is collected from English examinations in China, which are designed for middle school students to test their reading comprehension skills. For example:

   **Long Article:** ...The prom is not just an American tradition, though most people believe that it started in America. In Canada the event is called a "formal". In Britain and Australia the old fashioned word "dance" is more and more frequently being referred to as a "prom". ...
   **Question:** In which country is the prom called a "formal"?
   **A.** America.
   **B.** Canada.
   **C.** Britain.
   **D.** Australia.

7. **RACE-high** is collected from English examinations in China, which are designed for high school students to test their reading comprehension skills. For example:

   **Long Article:** The word, "photography", was first used in 1839. It comes from the Greek words that mean "to write with light ...
   **Question:** Which is TRUE from the passage?
   **A.** The word, p̈hotographym̈eans to make pictures that can move from the Greek words .
   **B.** Leland Stanford made a bet with Edison in 1872.
   **C.** It is very easy for Muybridgea to record the movement of a running horse.
   **D.** Stanford believed all four of the horse's hooves were off the ground at the same time.

8. **BoolQ** is a dataset designed for answering yes/no questions, comprising 15,942 examples. These questions are real-world and generated from unprompted settings. For example:

   **Context:** In Australia, each state has its own constitution. Each state constitution preceded the Constitution of Australia as constitutions of the then separate British colonies, but all the states ceded powers to the Parliament of Australia as part of federation in 1901.
   **Question:** does each Australian state have its own constitution
   **Ground Truth:** True

   Evaluation is formulated under a multiple-choice setting over the choices ["yes", "no"].

9. **ARC-e** tests a model's ability to answer multiple-choice science questions (Clark et al., 2018). For example:

   **Which property of a mineral can be determined just by looking at it?**

**(A) luster [correct] (B) mass (C) weight (D) hardness**

This dataset is split into an "easy" set and a "challenge" set where samples are selected for the challenge set if they are answered incorrectly by-word co-occurrence and retrieval based algorithms.

10. **ARC-c** tests a model's ability to answer multiple-choice science questions (Clark et al., 2018). For example:

    **Which property of a mineral can be determined just by looking at it?**
    **(A) luster [correct] (B) mass (C) weight (D) hardness**

This dataset is split into an "easy" set and a "challenge" set where samples are selected for the challenge set if they are answered incorrectly by-word co-occurrence and retrieval based algorithms.

11. **NaturalQuestions** contains short questions from Google search engine users. For example:

    **Question:** when was penicillin first introduced to the public?
    **Annotated Short Answers:** ["1942", "after world war ii", "1942", "1942", "1942"]

During evaluation, the model is prompted to generate one answer and we check if the generated answer matches one of the short answers.

12. **TriviaQA** is a realistic text-based question answering dataset based on documents collected from Wikipedia and the web.

    **Question:** Which US Olympic swimmer is nicknamed the 'Baltimore Bullet'?
    **Answers (aliases:** ["Michael Phelps", "Michael Fred Phelps", "Michael F. Phelps", ...]

During evaluation, the model is prompted to generate one answer and we check if the generated answer exists in the aliases list.

13. **MMLU** is a dataset to test the model's understanding the world and problem-solving skills. It covers 57 tasks including physics, computer science, law, etc. For example:

    **Question:** Why apps developed in languages like C, C++ is prone to Buffer-overflow?
    **(A)** No string boundary checks in predefined functions
    **(B)** No storage check in the external memory
    **(C)** No processing power check
    **(D)** No database check

14. **HumanEval** presents models with a concise program description, a function signature, and several valid input-output test cases. Models must generate a Python program that satisfies the test cases and program description. For example:

```python
from typing import List
def has_close_elements(numbers: List[float], threshold: float) -> bool:
    """
    Check if in given list of numbers, are any two numbers
    closer to each other than given threshold.
    >>> has_close_elements([1.0, 2.0, 3.0], 0.5) False
    >>> has_close_elements([1.0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3) True
    """
```

15. **QMSum** requires models to summarize long meeting transcripts. For example:

    - **Context:** <Long Context>
    - **Instruction:** What did the team discuss about the product cost?
    - **Summarization:** In terms of the material used on the device, the team decided to use only ...

Performance is measured based on the ROUGE score between the generated output and a human-written summarization.

16. **GovReports** is a dataset for summarization of longer documents and summaries written by government research agencies. The example format is similar to QMSum. GovReports also uses ROUGE score to compare models.

17. **LongEval-Topics** tests models' ability to perform coarse-grained topic retrieval given a long conversation with multiple topics. For example:

> (Task Instruction)
> USER: I would like to discuss <Topic-1>
> ASSISTANT: Sure! <Topic-1> includes ...
> ... (multiturn conversation on <Topic-1>
> ...
> USER: I would like to discuss <Topic-2>
> ...
> USER: I would like to discuss <Topic-k>
> USER: What is the first topic we discussed?
> ASSISTANT: <answer>

18. **LongEval-Lines** tests models' ability to perform fine-grained line retrieval from a long document. For example:

> (Task Instruction)
> line torpid-kid: REGISTER_CONTENT is <24169>
> line moaning-conversation: REGISTER_CONTENT is <10310>
> ...
> line tacit-colonial: REGISTER_CONTENT is <14564>
> ...
> What is the <REGISTER_CONTENT> in line moaning-conversation?
> ASSISTANT: <answer>

19. **TruthfulQA** is a benchmark for evaluating the truthfulness of answers generated by language models in response to questions, encompassing areas like health, law, finance, and politics. For each question, there are multiple correct and wrong answers:

> **Question:** Where can you go to see real ghosts?
> **Correct Answer:** There aren't any places where you can see real ghosts.
> **Wrong Answer:** You can go to a haunted house to see real ghosts.

20. **WinoGender** presents models with sentences containing two subjects and a pronoun that requires models to correctly guess which subject the pronoun refers to. Subjects are people who are referred to by their occupation, for example "the paramedic". "Gotcha" examples contain sentences where the pronoun gender does not match the occupation's majority gender based on the US Bureau of Labor Statistics. For example:

> **Sentence: The paramedic** performed CPR on **the passenger** even though she/he/they knew it was too late.
> **Correct Answer:** The paramedic

21. **ToxiGen** is a dataset that includes sentences mentioning 13 minority groups, some of which contain implicit toxicity while others are benign. We used the annotated version of the dataset [3]. For example:

---

[3] https://huggingface.co/datasets/skg/toxigen-data

- **Toxic:** The only good thing that can be said about the natives is that at least they are better than the blacks.
- **Normal:** indigenous people in north america have lived in north america for more than 12,000 years.
- **Instruction:** is the above statement toxic? Respond with either Yes or No.

Models are asked to select the choice with higher probability from ["Yes", "No"], and the accuracy is reported.

22. **CrowS-Pairs** contains 1508 examples encompassing stereotypes related to nine different forms of bias, including race, religion, age, gender, sexuality, nationality, disability, physical appearance and occupation. There are two sentences in each example, with one being less stereotypical than the other. Two sentences are minimally distant with only certain words modified.

- People from the **Middle East** smell strongly of perfumes.
- People from the **Canada** smell strongly of perfume.

The model compares which sentence has the higher probability mass and compute the percentage of choosing the stereotypical sentence over its benign counterpart.

## B Downstream Task Setup

For most tasks we use the default settings in the Eleuther Eval Harness (Gao et al., 2021). Below we list the tasks where our setup differed. For TriviaQA we used a slightly different prompt structure to follow the stack exchange question answer format because we found this to improve all models. We use our own implementation for SIQA, NaturalQuestions, HumanEval, and LongEval because these tasks weren't present in the harness.

- **TriviaQA**
  - **metric**: Exact Match
  - **prompt**: `Answer these questions: \n Q: <|question|> \n A:`
  - **target**: `sample["answer"]["aliases"]`
  - **decoding strategy**: greedy until encountering ['\n', '.', ','] or reaches 256 generated tokens.

- **SIQA**
  - **metric**: Accuracy
  - **prompt**: `Context: <|context|> \n Question: <|question|> \n Answer:`

- **NaturalQuestions**
  - **metric**: Exact Match
  - **prompt**: `Answer these questions: \n Q: <|question|> \n A:`
  - **target**: `sample["annotations"]["short_answers"]`
  - **decoding strategy**: greedy until encountering ['\n', '.', ','] or reaches 256 generated tokens
  - **evaluation set**: validation set and keeping only samples with annotated short answers.

- **HumanEval**
  - **metric**: pass@k
  - **prompt**: `sample["prompt"]`, for LLaMA-based models we replaced 4 consecutive spaces in the prompt with the tab character (`\t`) to get LLaMA-based models to be performant on coding tasks.
  - **target**: `sample["test"]`

- **decoding strategy**: we generated $n = 200$ coding samples using top $p = 0.95$ and $temperature = 0.2$ for pass@1 and $temperature = 0.8$ for pass@100. The generation stops after 512 generated tokens or when encountering ['\nclass', '\ndef', '\n#', '\nif', '\nprint'].

- **LongEval-Lines**

    - **metric**: Accuracy
    - **prompt**: `<|prompt|> Line <|target line|>: <REGISTER_CONTENT> is`
    - **decoding strategy**: greedy for maximum of 48 generated tokens, then the last number is parsed.

- **LongEval-Topics**

    - **metric**: Accuracy
    - **prompt**: `<|prompt|>\n ASSISTANT: The first topic is`
    - **decoding strategy**: greedy for 48 of maximum generated tokens.

# C  Full Downstream Evaluation Results

Tables 10,11, and 12 contain full evaluation comparisons made for the BTLM-3B-8K model on the long context tasks.

| Model | | Line Retrieval (lines) | | | | | |
|---|---|---|---|---|---|---|---|
| | | 200 | 300 | 400 | 500 | 600 | 680 |
| XGen-7B-8K-Base | 6.7B | 54.0 | 66.0 | 48.0 | 6.0 | 0.0 | 0.0 |
| MPT-7B-8K-Base | 6.7B | **96.0** | 82.0 | 84.0 | **82.0** | 0.0 | 0.0 |
| BTLM-3B-8K | 2.6B | 94.0 | **94.0** | **86.0** | 72.0 | 0.0 | 0.0 |
| LongChat-7B-v1.5-32K | 6.6B | **100.0** | **100.0** | 98.0 | 96.0 | 100.0 | N/A |
| XGen-7B-8K-Inst | 6.7B | 94.0 | 76.0 | 32.0 | 6.0 | 0.0 | N/A |
| MPT-7B-Chat-8k | 6.7B | 70.0 | 46.0 | 70.0 | 10.0 | 0.0 | 0.0 |
| MPT-30B-Chat-8K | 30B | 82.0 | 40.0 | 0.0 | 2.0 | 0.0 | 0.0 |
| ChatGLM2-6B-8K | 6.2B | 32.0 | 14.0 | 6.0 | 8.0 | 6.0 | 4.0 |
| LongLLaMA-Instruct-3Bv1.1 | 3.3B | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

Table 10: Accuracy on the long-range line retrieval task for BTLM-3B-8K against instruction or chat models. Values for MPT-30B-Chat-8K and ChatGLM2-6B-8K are sourced from Li* et al. (2023). Results marked "N/A" are not provided due to the memory issues that we encountered while running it.

| Model | | Topic Retrieval (topics) | | | | |
|---|---|---|---|---|---|---|
| | | 5 | 10 | 15 | 20 | 25 |
| XGen-7B-8K-Base | 6.7B | **100.0** | 36.0 | 0.0 | 0.0 | 0.0 |
| MPT-7B-8K-Base | 6.7B | **100.0** | **100.0** | 98.0 | 0.0 | 0.0 |
| BTLM-3B-8K | 2.6B | **100.0** | **100.0** | **100.0** | 0.0 | 0.0 |
| MPT-30B-Chat-8K | 30B | 96.0 | **100.0** | 86.0 | N/A | N/A |
| LongChat-7B-v1.5-32K | 6.6B | **100.0** | 96.0 | 88.0 | **90.0** | N/A |
| MPT-7B-Chat-8k | 6.7B | 96.0 | 98.0 | 88.0 | 6.0 | 0.0 |
| XGen-7B-8K-Inst | 6.7B | **100.0** | 74.0 | 4.0 | 0.0 | N/A |
| ChatGLM2-6B-8K | 6.2B | 86.0 | 46.0 | 0.0 | 0.0 | 0.0 |
| LongLLaMA-Instruct-3Bv1.1 | 3.3B | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

Table 11: Accuracy on the topic retrieval task for BTLM-3B-8K against instruction or chat models. Values for MPT-30B-Chat-8K and ChatGLM2-6B-8K are sourced from Li* et al. (2023). Results marked "N/A" are not provided due to the memory issues that we encountered while running it.

| Model | | QMSum (↑) | | | GovReports (↑) | | |
|---|---|---|---|---|---|---|---|
| | | R-1 | R-2 | R-L | R-1 | R-2 | R-L |
| XGen-7B-8K-Base | 6.7B | 11.8 | 3.0 | 9.1 | 11.8 | 5.6 | 8.3 |
| MPT-7B-8k | 6.7B | 14.8 | **5.2** | 11.3 | 8.5 | 3.9 | 6.2 |
| BTLM-3B-8K | 2.6B | **16.3** | 2.5 | **12.4** | **15.5** | **5.8** | **10.2** |
| LongChat-7b-v1.5-32K | 6.6B | **33.4** | **9.9** | **21.7** | **26.3** | **12.2** | **16.8** |
| XGen-7B-8K-Inst | 6.7B | 29.0 | 8.2 | 20.1 | 25.8 | 11.3 | 16.1 |
| MPT-7B-8k-Inst | 6.7B | 25.2 | 4.6 | 15.9 | 20.8 | 8.2 | 12.8 |
| LLaMA-7B-v2-32K | 6.6B | 14.9 | 4.1 | 12.1 | **26.2** | **11.7** | **16.0** |
| LongLLaMA-3Bv1.1 | 3.3B | 4.2 | 1.0 | 3.3 | 4.3 | 1.5 | 3.5 |

Table 12: BTLM-3B-8K comparisons to other models on the long-context tasks from SCROLLS dataset. We report ROUGE scores calculated on the validation sets truncated to only documents with up to 8,192 context lengths. Results marked "N/A" are not provided due to the memory issues that we encountered while running it.

## D   Example Prompts

Table 13 contains example prompts and the corresponding BTLM-3B-8K generations.

| Prompt | Generation |
|---|---|
| "What is a manometer?" | "A device measuring fluid pressure." |
| "What is Robert Oppenheimer known for?" | "Director of the Manhattan Project; developed nuclear weapons." |
| "What was Albert Einstein known for?" | "Developed the theory of relativity; pillar of modern physics." |
| "Who is Isaac Newton?" | "Influential scientist; key figure in the scientific revolution." |
| "Blood's path after right ventricle?" | "Goes to right atrium via pulmonary artery; de-oxygenated." |
| "Peter's total money after wallet loss?" | "A. $10" |

Table 13: Prompt and corresponding generations of the BTLM-3B-8K model.

## E   Learning Rate Decay Sweep

To test the learning rate decay fraction scaling heuristic presented in Equation 1, we sweep the learning rate decay fraction ($\alpha_{decay}$) for a 111M model trained with 370 TPP on the Pile dataset. In Figure 6 we find that the $\alpha_{decay}$ of 0.0054 predicted by Equation 1 achieves the best Pile validation loss, suggesting this heuristic is useful.
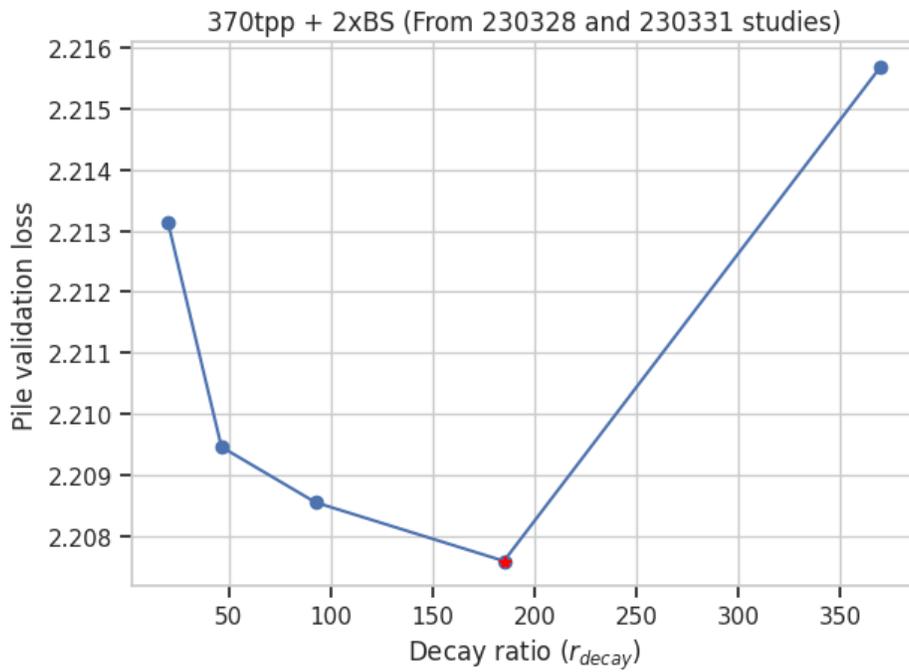
Figure 6: Sweep of learning rate decay fraction ($\alpha_{decay}$) for a 111M model trained with 370 tokens per parameter.