

Skywork: A More Open Bilingual Foundation Model

Tianwen Wei, Liang Zhao, Lichang Zhang, Bo Zhu, Lijie Wang, Haihua Yang, Biye Li, Cheng Cheng, Weiwei Lü, Rui Hu
Chenxia Li, Liu Yang, Xilin Luo, Xuejie Wu, Lunan Liu, Wenjun Cheng, Peng Cheng, Jianhao Zhang, Xiaoyu Zhang
Lei Lin, Xiaokun Wang, Yutuan Ma, Chuanhai Dong, Yanqi Sun, Yifu Chen, Yongyi Peng, Xiaojuan Liang
Shuicheng Yan, Han Fang, Yahui Zhou*

Skywork Team, Kunlun Inc.

Abstract

In this technical report, we present Skywork-13B, a family of large language models (LLMs) trained on a corpus of over 3.2 trillion tokens drawn from both English and Chinese texts. This bilingual foundation model is the most extensively trained and openly published LLMs of comparable size to date. We introduce a two-stage training methodology using a segmented corpus, targeting general purpose training and then domain-specific enhancement training, respectively. We show that our model not only excels on popular benchmarks, but also achieves *state of the art* performance in Chinese language modeling on diverse domains. Furthermore, we propose a novel leakage detection method, demonstrating that data contamination is a pressing issue warranting further investigation by the LLM community. To spur future research, we release Skywork-13B along with checkpoints obtained during intermediate stages of the training process. We are also releasing part of our SkyPile corpus, a collection of over 150 billion tokens of web text, which is the largest high quality open Chinese pre-training corpus to date. We hope Skywork-13B and our open corpus will serve as a valuable open-source resource to democratize access to high-quality LLMs.

1 Introduction

Natural Language Processing (NLP), a vital branch of artificial intelligence, has experienced a transformative surge in recent years. Pivotal to this revolution has been the advent and advancement of large language models (LLMs) (Ouyang et al., 2022; OpenAI, 2023; Bubeck et al., 2023; Chowdhery et al., 2022; Anil et al., 2023; Touvron et al., 2023a,b). These complex computational structures, composed of billions of parameters, are capable of understanding,

generating, and translating human language with an unprecedented degree of accuracy and sophistication. However, the proliferation of these models has also been accompanied by a growing trend towards commercialization and a lack of transparency, a phenomenon that is increasingly influencing the dynamics of the open-source community.

Historically, the open-source community has thrived on the principles of collaboration, transparency, and unrestricted sharing of ideas. However, as the commercial potential of LLMs has been recognized, this openness has begun to diminish. The reality is that many organizations only make model checkpoints publicly accessible, while withholding vital information on model reproduction. This practice significantly hampers the progress of the field.

In an effort to revive the spirit of the open-source community and contribute to the ongoing dialogue about transparency in AI, we present Skywork-13B: a family of bilingual large language models with 13 billion parameters, trained on a colossal corpus of more than 3.2 trillion tokens drawn from both English and Chinese texts. To our knowledge, our Skywork-13B is the most thoroughly trained family of open LLMs of comparable size to date.

In this technical report, we offer a comprehensive disclosure of the Skywork-13B developmental journey. We detail the composition of our training data, provide insights into the evolutionary trajectory of the model’s abilities during training, and share methodologies that could be employed to enhance model ability in specific domains. We believe that such an open approach not only aids in the reproducibility of our work but also provides a valuable resource for other researchers seeking to explore and expand the capabilities of large language models. This technical report is also a call to

* Email: {forename}.{surname}@kunlun-inc.com

action for renewed transparency in the field of NLP. Through it, we hope to inspire a return to a more collaborative, open-source community, where progress is not hampered by commercial considerations but propelled by collective intelligence and shared wisdom.

Our contributions are the following:

- We release Skywork-13B¹, a family of LLMs that is the most extensively trained and openly published LLMs of comparable size to date. Our Skywork-13B family includes 1) Skywork-13B-Base, a strong foundation model with *state of the art* Chinese language modeling capability, and 2) Skywork-13B-Chat, a fined-tuned version optimized for conversation².
- We disclose detailed information on the training process and data composition. We also release intermediate checkpoints, which provide a valuable resource for understanding how the model’s capabilities develop over the course of training. It enables other researchers to leverage these checkpoints for their specific use-cases.
- We release a portion of our high quality training corpus, totaling more than 150 billion tokens. To our knowledge, this is the largest open Chinese corpus for language model pre-training to date.
- We develop a novel method that detects the level of in-domain data usage during the training stage. To facilitate reproduction of the experiments presented in this report, we have released the relevant data.

2 Methodology

2.1 Two Pre-training Stages

In order to train Skywork-13B, we constructed SkyPile (see Section 3.1), a massive training corpus primarily constituted by publicly accessible web pages. We identified a small subset of SkyPile, encompassing exercises and solutions that span a broad spectrum of subjects from primary to graduate school. This includes

¹Github repository: <https://github.com/SkyworkAI/Skywork>.

²In this technical report we focus on the development of the base model. Details on Skywork-13B-Chat can be found in our Github repository.

coding problems, national exam questions, textbook exercises, and others. Given the majority of these exercises are STEM-related, we henceforth refer to this subset and its complement as SkyPile-STEM and SkyPile-Main, respectively.

Rather than training the Skywork-13B foundation model directly on SkyPile as a whole, we adopted a two-stage training approach. The first stage, which constitutes the primary pre-training phase, involves training the model from scratch on SkyPile-Main. In the second stage, our Skywork-13B is enriched with STEM-related domain knowledge and problem-solving skills through continual pre-training on SkyPile-STEM. To circumvent the potential issue of catastrophic forgetting, this continual pre-training is performed on a mix of SkyPile-STEM and SkyPile-Main, rather than exclusively on SkyPile-STEM.

The decision to segregate Stage-1 and Stage-2 pre-training serves a dual purpose. Firstly, we acknowledge that a significant proportion of the samples from SkyPile-STEM are, by their nature, supervised data. Those data are closely related to popular benchmarks such as CEVAL (Huang et al., 2023), MMLU (Hendrycks et al., 2021) and GSM8K (Cobbe et al., 2021), and can be utilized in a supervised fine-tuning (SFT) process to directly enhance model performance on related downstream tasks. In this context, the separation between Stage-1 and Stage-2 training enables us to more effectively assess the impacts of general-purpose pre-training (on web texts) and targeted pre-training (on in-domain/supervised data). Such insights could inform future data collection and compilation strategies for foundational model training.

Secondly, by restricting first stage pre-training to general-purpose data, we are able to produce a version of foundation model as an alternative to the one with targeted enhancement. While the latter demonstrates superior performance on certain downstream tasks, it is less capable in language modeling of natural texts. We posit that this alternative is a valuable contribution to the community, given its potential to excel in applications that do not require STEM-related competencies.

2.2 Training Progress Monitoring

It is of vital importance to monitor and assess progress made during pre-training in real-time.

Existing methods such as monitoring training loss and benchmark results on intermediate checkpoints, however, have their limitations.

The main issue of monitoring training loss lies in that its effectiveness comes into question when considering the potential of overfitting. The training loss is equivalent to validation loss only if the training data is utilized exactly once (i.e., in one epoch). Yet, in practical scenarios of training LLMs, high-quality data often go through the training process multiple times (Taylor et al., 2022; Touvron et al., 2023a; Rozière et al., 2023; Gunasekar et al., 2023; Li et al., 2023b). Besides, even after explicit de-duplication, there may still exist significant amount of duplicated data in the training set (Soboleva et al., 2023; Abbas et al., 2023). In either cases, solely relying on training loss can lead to overlooking the issue of overfitting, thereby producing overly optimistic estimates of model performance. The top left subplot in Figure 3 illustrates the trajectory of the pre-training loss for our Skywork-13B model. Consistent with findings reported in (Touvron et al., 2023a,b; Baichuan Inc., 2023), the loss demonstrates a steady decline throughout the training process. However, an observation not disclosed in these cited works is the behavior of the validation loss on held-out sets. From the figure it can be clearly seen that the validation losses seem to level off as training approaches its final stages.

Benchmarking based on intermediate checkpoints is another common monitoring approach (Touvron et al., 2023a; Baichuan Inc., 2023). Nevertheless, it presents several challenges. Firstly, there is a high variance in benchmark results, which can lead to unstable and unreliable assessments of training progress. Secondly, benchmark results are not sensitive to minor progress in training. This insensitivity makes it difficult to accurately track gradual improvements during the training process. Besides, weaker models do not follow instructions well. Hence benchmark results may not accurately reflect their true learning progress or potential. Finally, an inconvenience posed by most benchmarks is the necessity for model generation. This process is notably resource-intensive, demanding substantial computational power.

During the pre-training of Skywork-13B, we

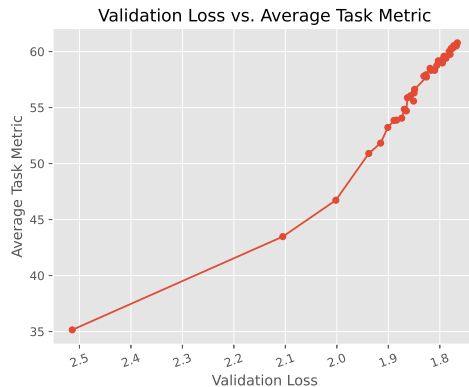


Figure 1: Validation loss on English web texts vs. average task metric during the pre-training of Skywork-13B. The tasks include BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2019), Winogrande (Sakaguchi et al., 2021), TriviaQA (Joshi et al., 2017) and RACE (Lai et al., 2017).

embrace the method of monitoring the language modeling loss across numerous reserved validation sets, each reflecting a distinct data distribution. More specifically, we have created separate validation sets for code, academic publications, social media posts, web texts in Chinese and English, among others. Conventional monitoring metrics are also utilized, but they serve merely as supplementary tools. In Figure 1 we plot the curve of language model validation loss on English web texts against the average metric of several English downstream tasks. It is apparent that there is a very high correlation between the two quantities, showing that validation loss can serve as a valid proxy metric for downstream task performance. In the context of LLM pre-training, this approach also yields several other benefits:

- **Ease of construction:** Crafting multiple validation sets is a relatively effortless task. This enables the evaluation of a model’s language modeling performance across varied domains.
- **Simplicity in computation:** Calculation of validation loss is straightforward, significantly reducing the computational and logistical overhead associated with tracking model training.
- **High sensitivity to training progress:** Validation loss is finely attuned to the progression of training, thereby offering a more detailed

perspective on how models evolve and improve over time.

- **Model-agnosticism:** Validation loss is indifferent to the composition of the training corpus or the model architecture. It allows for comparison not only between different checkpoints produced within a single training session, but also across varied models from the community. This ensures a consistent and equitable basis for model comparison.

Note that monitoring the validation loss on a held-out set sharing the same distribution as the training set is a ubiquitous practice in machine learning. However, the observation of validation loss across multiple held-out sets, each with deliberate, unique distributions, is not common. We also note that the perspective asserting the primacy of language modeling loss as the paramount performance metric for models is not a recent revelation. This principle has been either explicitly or implicitly adopted in a number of research studies, as exemplified in (Kaplan et al., 2020; Hoffmann et al., 2022; Anil et al., 2023; Xia et al., 2023; Delétang et al., 2023).

3 Pre-training

3.1 SkyPile Corpus

In order to train Skywork-13B, we build SkyPile, a vast, high quality corpus comprising more than 6 trillion tokens. A segment of the corpus, comprising over 150 billion tokens of web text, has been open sourced to facilitate research and training on Chinese LLMs³.

Our SkyPile is an amalgamation of several sources, the overwhelming majority of which is gleaned from publicly accessible channels. Numerous prior research works, exemplified by initiatives such as LLaMA (Touvron et al., 2023a) and RefinedWeb (Penedo et al., 2023), have substantiated the notion that publicly accessible web data can yield exceptionally high-quality LLMs. In alignment with this empirical evidence, we subscribe to the premise of leveraging publicly accessible webpages as our primary source for training data.

³huggingface.co/datasets/Skywork/SkyPile-150B

The construction of SkyPile is characterized by a dedicated emphasis on two primary dimensions: text quality and information distribution. Our data processing pipeline, inspired by (Wenzek et al., 2020; Touvron et al., 2023a; Penedo et al., 2023), incorporates the following stages:

- **Structural Extraction:** Due to the predominant source of our dataset being publicly accessible web pages, the objective of the first stage is the extraction of pertinent content while concurrently expunging extraneous textual elements that are deemed non-contributory to the training of our language model, e.g. these superfluous components include navigational bars, site-specific contact information, disjunctive title texts devoid of substantive content, etc. Subsequent to this culling process, the retained information predominantly consists of contiguous, medium to long-form textual passages.
- **Distribution Filtering:** In the pursuit of cultivating a profoundly adept LLM, the model’s exposure must encompass a diverse array of content spanning an extensive spectrum of domains. Prior endeavors within the field have entailed the task of assigning categorical labels to each individual document or webpage, thereby manually dictating the composition of the training corpus. However, we posit that the corpus employed for LLM training has burgeoned to such an extent that the knowledge it encapsulates can not be compartmentalized discretely. Consequently, eschewing a label-centric approach, our methodology centers on benchmarking the semantic affinities existing between textual segments, thereby identifying and omitting those text blocks characterized by an exceedingly high recurrence rate.
- **Deduplication:** Deduplication has demonstrated its remarkable efficacy in enhancing the overall quality of a training corpus, and it has found extensive application in virtually all prominent datasets (Hernandez et al., 2022; Kandpal et al., 2022; Abbas et al., 2023; Lee et al., 2022). Within the framework of SkyPile, we regard deduplication as an integral component of the Distribution Filtering process. When considering the broader perspective, it becomes evident

that duplication constitutes a paramount factor influencing the semantic distribution of a corpus. Consequently, the techniques and strategies we employed during the distribution filtering phase autonomously eliminated a substantial portion of duplicated content.

- **Quality Filtering:** In this phase, we deploy the CCNet (Wenzek et al., 2020) pipeline to perform two critical filtration tasks: the elimination of content of inferior quality and the exclusion of pages that are neither in English nor Chinese. We trained a binary classifier that predicts the likelihood that a given webpage is suitable for inclusion as a reference within the Wikipedia corpus. The outcome of this stage is organized into distinct quality-based categories, and we retain exclusively the high quality groups, opting to discard the remaining groups in its entirety.

Above we described our pre-processing pipeline for natural text. As for Github content, we employ an approach that is similar to (Together Computer, 2023). We have devised a collection of straightforward yet efficacious heuristics, encompassing criteria such as line length filtration and alphanumeric thresholds, designed to discern and exclude content of low quality. Our criteria are specifically oriented toward enhancing content quality, as opposed to merely curbing its volume. Notably, in contrast to prevailing practices that involve the wholesale removal of a significant portion of json, xml, yaml, and html content, we have made a deliberate choice to retain a judiciously proportionate representation of these data formats.

Note that in pursuit of harmonizing the model’s proficiency in both English and Chinese, we include in SkyPile a curated high-quality parallel corpora. This data is meticulously structured to pair a complete English paragraph with its corresponding Chinese counterpart, ensuring a seamless alignment of linguistic capabilities between the two languages.

3.2 Training Data Composition

Our Skywork-13B is pre-trained for 3.2 trillion tokens, sampled from SkyPile. Texts from certain sources are deemed as of high quality, e.g.

	Category	Percentage
English	Webpages	39.8%
	Books	3.6%
	Academic Papers	3.0%
	Encyclopedia	0.5%
	Miscellany	2.9%
Chinese	Webpages	30.4%
	Social Media	5.5%
	Encyclopedia	0.8%
	Miscellany	3.1%
Other Lang.	Encyclopedia	2.4%
Code	Github	8.0%

Table 1: Breakdown of training data in Stage-1 pre-training of Skywork-13B.

Wikipedia, hence have undergone upsampling. However, we generally stick to the rule that the number of repetition does not exceed five, as is recommended by recent studies (Taylor et al., 2022; Muennighoff et al., 2023).

We report in Table 1 a breakdown of the constituent components of the training tokens during Stage-1 pre-training. The training tokens are primarily composed of English and Chinese texts, constituting 49.8% and 39.6% of the data, respectively. Code contributes 8.0% to the total, with texts in other languages accounting for the remaining 2.4%. The category labeled as “miscellany” encompasses a diverse range of texts, including but not limited to, legal articles, court documents, company annual reports, and classical literature.

3.3 Tokenizer

We tokenize the data using byte-pair encoding (BPE) as implemented in SentencePiece (Kudo and Richardson, 2018), following the approach of LLaMA (Touvron et al., 2023a). Since our model is intended to be English-Chinese bilingual, we extend the original vocabulary of LLaMA, which primarily consists of latin-based words and subwords, with frequently used Chinese characters and words. Specifically, we add 8000 single-character tokens from BERT’s vocabulary (Devlin et al., 2019) to LLaMA’s vocabulary. We further expand the vocabulary with 25k frequent Chinese multi-character words. This results in a total vocabulary size of 65,536 tokens, of which 17 are reserved as

special symbols.

As in LLaMA, we split all numbers into individual digits, and fall back to bytes to decompose unknown UTF-8 characters.

	Category	Size
	Latin based words & subwords	32000
	Chinese characters & Unicode symbols	8000
	Chinese words	25519
	Reserved symbols	17
	Total	65536

Table 2: Breakdown of the vocabulary used in Skywork-13B.

3.4 Architecture

Our Skywork-13B is based on the transformer architecture (Vaswani et al., 2017), consisting of stacks of transformer-decoder layers. In contrast to the original transformer model, we have incorporated several modifications, inspired by LLaMA (Touvron et al., 2023a,b). Our preliminary experiments, as illustrated in Figure 2, validate these changes, demonstrating the improved performance they confer. Details on this experiment can be found in Appendix A.

While our network architecture takes after the LLaMA model to a great extent, there exists a notable difference in our preference for a deeper, yet narrower, network. A comparative exploration of the Skywork-13B and LLaMA2-13B network configurations is presented in Table 3.

The specific modifications made are described in detail below.

- **Positional Embedding:** We use Rotary Positional Embedding (RoPE) (Su et al., 2022), that was motivated by its extensive adoption in various prominent large language models, such as LLaMA and PaLM, as well as its demonstrated effectiveness in extending the length of context windows, as evidenced by recent studies (Chen et al., 2023; Rozière et al., 2023; Xiong et al., 2023).
- **Layer Normalization:** We replaced the conventional layer normalization with RMSNorm (Zhang and Sennrich, 2019). Additionally, we adopted pre-normalization in each layer instead of post-normalization, which has been shown to enhance the training stability of transformer models.

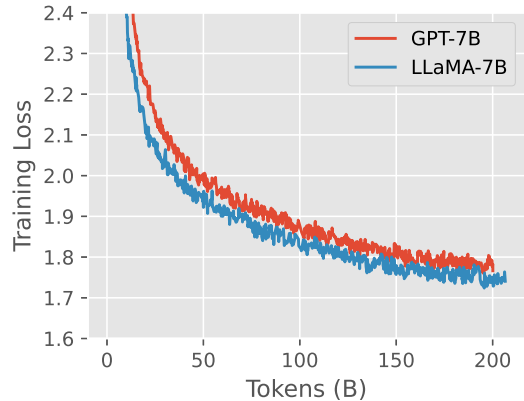


Figure 2: Preliminary Experiments: Comparison of conventional GPT architecture and more recent LLaMA architecture. For each of the two transformer variants, a model with 7 billion parameters is trained from Scratch on 200 Billion Tokens. The plot clearly shows that the LLaMA architecture achieves a lower training loss than GPT, demonstrating the former’s superiority.

- **Activation:** We employed the SwiGLU activation function (Shazeer, 2020). In line with established conventions in prior studies, we reduced the dimension of the feed-forward network (FFN) from four times the hidden size to eight-thirds of the hidden size. This adjustment was made to maintain parity between the total parameters in a layer and those in the vanilla transformer layer.

	LLaMA2-13B	Skywork-13B
Vocab. Size	32,000	65,536
Hidden Dim.	5,120	4,608
FFN Dim.	13,696	12,288
Head Dim.	128	128
Num. Heads	40	36
Num. Layers	40	52
Seq. Len.	4,096	4,096
#Tokens per Batch	4M	16M
Peak LR	3e-4	6e-4
Minimum LR	3e-5	6e-5

Table 3: Comparisons in architecture and important hyper-parameters of Skywork-13B and LLaMA2-13B.

3.5 Infrastructure

Our Skywork-13B is trained on a cluster of 64 NVIDIA-HGX-A800 nodes, a total of 512 A800-80G SXM GPUs. Each node in the cluster is outfitted with high-speed 400GB/s NVLinks

for intra-node communication and an 800Gb/s RoCE network for inter-node connectivity. Our training framework is based on Megatron-LM (Shoeybi et al., 2020) library, designed to support the stable, prolonged training of large-scale models, accommodating thousands of GPUs and model sizes in the order of hundreds of billions parameters.

Considering the relatively moderate size of our Skywork-13B model, we have avoided the use of GPU memory optimization techniques and parallel schemes that could impede speed. These include Tensor Model Parallelism (Shoeybi et al., 2020), Sequence Parallelism (Korthikanti et al., 2022), ZeRO-Stage2 (Rajbhandari et al., 2020), and Checkpointing (Chen et al., 2016). Instead, we have leveraged Data Parallelism (DP) with ZeRO-1 (Rajbhandari et al., 2020) and Pipeline Parallelism (PP) (Narayanan et al., 2021) as the primary parallelization strategies for training Skywork-13B. ZeRO-1 substantially diminishes the GPU memory footprint of the Adam optimizer state without increasing the burden on intercommunication. Pipeline Parallelism offers memory optimization at a minimal communication overhead, which decreases as the gradient accumulation step increases, thereby mitigating the slowdown of all-reduce as DP Size increases. Regarding operator optimization, we adopted Flash Attention V2 (Dao et al., 2022; Dao, 2023), a strategy that both optimizes GPU memory and expedites the training process.

Upon extensive preliminary experiments, we have decided to adopt the combination of DP256, PP2, and ZeRO-1 as our distributed training strategy for Skywork-13B. With this configuration, we achieved a token throughput of 1873 per GPU per second and a model flops utilization (MFU) of 56.5%. An overview of these experiments is provided in Appendix B. The training process of Skywork-13B spanned a total of 39 days.

3.6 Training Details

As outlined in Section 2.1, the pre-training of Skywork-13B is executed in two stages:

- **Stage-1:** General purpose pre-training on SkyPile-Main.
- **Stage-2:** STEM-oriented continual pre-training on SkyPile-STEM.

In both stages, the model is trained using the standard auto-regressive language modeling objective, with context lengths fixed at 4096 tokens. The AdamW optimizer (Loshchilov and Hutter, 2019), applied for the training process, uses β_1 and β_2 values of 0.9 and 0.95, respectively. Throughout the pre-training, we applied a weight decay of 0.1 and gradient clipping of 1.0. Our model was trained with bfloat16 mixed precision.

3.6.1 Stage-1 Pre-training

In the first stage, our Skywork-13B model is trained from scratch on SkyPile-Main for over three trillion tokens. This stage consists of two sequential training sessions, covering the first $0 \sim 2T$ tokens and the subsequent $2 \sim 3T$ tokens, respectively.

Our initial plan was to train Skywork-13B for two trillion tokens. We launched a training session accordingly, with a cosine learning rate schedule that gradually decays from a peak learning rate of $6e-4$ to a final learning rate of $6e-5$. In Figure. 3, we report in red curves the evolution of language modeling losses and several benchmark results of our Skywork-13B during this session. It is evident that by the end of this session, the model had not reached saturation. We hypothesized that the model could further benefit from additional pre-training, prompting us to launch a secondary training session targeting an additional one trillion tokens.

The second training session utilized a slightly different composition of training data compared to the initial $0 \sim 2T$ session, as data from certain sources had been depleted and fresh sources were introduced. Owing to the shift in the training distribution, we meticulously tuned the learning rate parameter, eventually deciding on a constant learning rate of $6e-5$ for the $2 \sim 3T$ session. In Figure. 4, we illustrate the model losses under varying learning rate conditions. Results indicate that a higher learning rate leads to escalations in training loss which we deem too costly to reverse. The impact of the second training session is depicted in blue curves of Fig. 3. The enhancement in the model’s performance continues, albeit at a decelerating pace. Interestingly, although our Skywork-13B trails in the realm of English language modeling, it significantly surpasses all

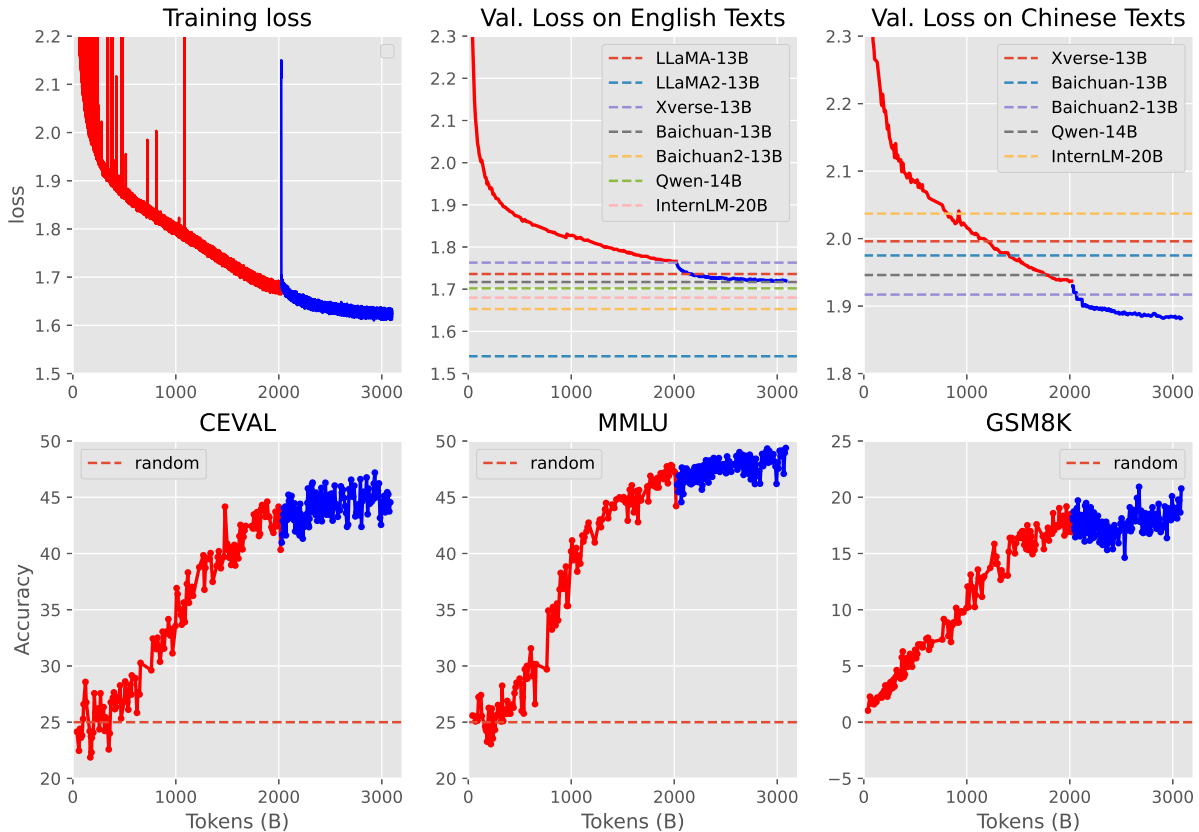


Figure 3: Trajectory of important monitoring metrics during Stage-1 pre-training. Top Left: Training loss. Top Middle and Right: Validation loss on English and Chinese held-out sets of web texts. The horizontal dashed lines in the middle and right plots correspond to the evaluated language modeling loss for several similar-sized open LLMs. Bottom: Benchmark results on CEVAL, MMLU and GSM8K respectively. Stage-1 pre-training consists of two sequential training sessions, represented by different colors in the loss curves (red for session 0 \sim 2T and blue for session 2 \sim 3T).

other comparable open LLMs in Chinese language modeling. In Section 4.3, we will confirm that the superiority of our Skywork-13B in Chinese language modeling is not only true on our validation set, it also holds true on a number of test sets sourced from diverse domains.

More results can be found in Appendix (see Figure 6).

3.6.2 Stage-2 Pre-training

The primary aim of Stage-2 pre-training is to augment the model with capabilities pertinent to STEM disciplines. The data utilized in this stage comprises an approximate 20% from SkyPile-STEM and 80% from SkyPile-Main, amassing a total of roughly 130 billion tokens. A constant learning rate of $6e-5$ is adopted, maintaining parity with the terminal learning rate used in Stage-1 pre-training

Consequent to the data distribution shift from Stage-1 to Stage-2, it becomes crucial

to meticulously calibrate the sampling ratio between the different data sources. Initial experiments revealed that a gradual increment in the SkyPile-STEM ratio yielded the most effective results. Therefore, for the actual Stage-2 pre-training phase, we implemented a sampling plan that commenced with 10% of SkyPile-STEM initially, gradually escalating to a peak of 40% towards the conclusion of the training.

This training strategy proved successful in maintaining the stability of the model’s language modeling validation loss while enabling an optimum transfer of STEM knowledge. The extended training period ensures a comprehensive assimilation of STEM-related knowledge into the model without causing significant disturbance to the pre-existing learned information.

The impact of Stage-2 pre-training is illustrated in Figure 5, which presents the progres-

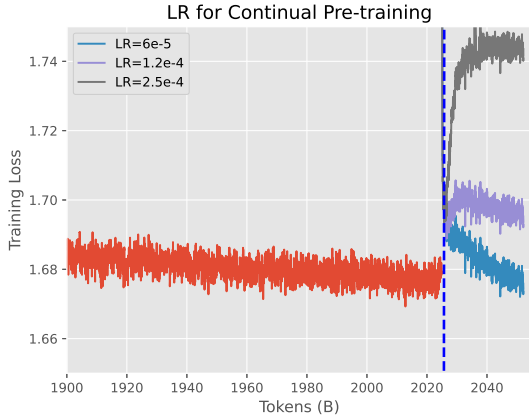


Figure 4: Test runs for tuning the learning rate of the 2 ~ 3T training session. It can be seen that 6e-5, which is the terminal learning rate from 0 ~ 2T training session, yields the best result.

sion of the CEVAL benchmark score. The evolution of scores on other STEM-related benchmarks, such as GSM8K, mirrors a similar trend. Improvements in individual subjects of the CEVAL can be found in Table 12 (see appendix).

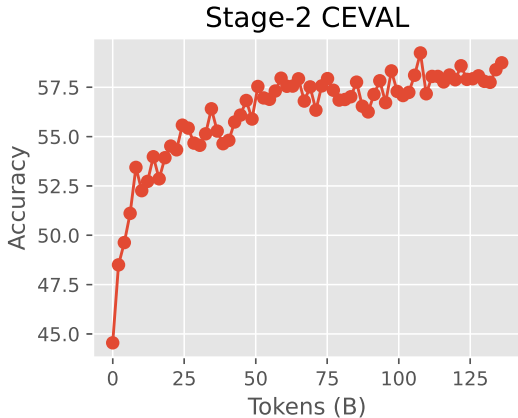


Figure 5: Evolution of CEVAL score during Stage-2 pre-training.

4 Evaluation

4.1 Baselines

We compare the performance of our Skywork-13B with open models that are similar in size, including LLaMA-13B (Touvron et al., 2023a), LLaMA2-13B (Touvron et al., 2023b), Baichuan-13B, Baichuan2-13B (Baichuan Inc., 2023), Xverse-13B (Xverse-AI, 2023), InternLM-20B (InternLM Team, 2023). A summary of these models can be found in Table 4.

Model	#Tokens	Language
OpenLLaMA-13B	1.0T	English
LLaMA-13B	1.0T	English
LLaMA2-13B	2.0T	English
Baichuan-13B	1.4T	English & Chinese
Baichuan2-13B	2.6T	English & Chinese
Xverse-13B	1.4T	English & Chinese
InternLM-20B	2.3T	English & Chinese
Skywork-13B	<u>3.2T</u>	English & Chinese

Table 4: Details of various models. The column labeled "#Tokens" indicates the quantity of training tokens used by each model, whereas the "Language" column specifies the primary languages supported by each model.

4.2 Benchmark Evaluation

We focus on the following popular benchmarks:

- MMLU (Hendrycks et al., 2021): MMLU is a benchmark designed to measure knowledge acquired during pre-training. The benchmark covers 57 subjects across STEM, the humanities, the social sciences, and more, ranging in difficulty from an elementary level to an advanced professional level. It tests both world knowledge and problem solving ability.
- CEVAL (Huang et al., 2023) and CMMLU (Li et al., 2023a): Those are Chinese benchmarks that mimic MMLU. CEVAL consists of 13948 multi-choice questions spanning 52 diverse disciplines and four difficulty levels. CMMLU covers 67 disciplines that span from elementary to advanced professional levels.
- GSM8K (Cobbe et al., 2021): This dataset consists of 8500 high-quality grade school math word problems created by human writers. These multi-step problems require between 2 and 8 steps to solve. GSM8K is usually used in benchmarking multi-step mathematical reasoning ability of LLMs.

In Table 5 we present a comparison of performance results from different models on these benchmarks. The metrics for CEVAL, CMMLU and MMLU are 5-shot accuracy, while for GSM8K it is 8-shot accuracy. Higher numbers indicate better performance. It can be seen that our Skywork-13B achieves the highest score on both the CEVAL and MMLU and

GSM8K benchmarks, with scores of 60.6 and 62.1 and 55.8 respectively. On the CMMLU benchmark, Baichuan2-13B achieves the highest performance with a score of 62.0.

In summary, our Skywork model has demonstrated exceptional performance across a diverse range of comprehensive benchmark tests.

Results of individual subjects of the CEVAL can be found in Table 12. Results of other benchmarks can be found in Appendix C.

4.3 Language Modeling Results

4.3.1 LM as a solution to benchmark overfitting

Conventional benchmarks for evaluating LLMs often rely on static datasets of human-annotated examples. A core issue with this approach is that updating the test samples regularly is difficult and costly. Over time, the static test sets tend to be overfitted, producing misleading benchmark results.

We propose language modeling evaluations as a compelling alternative. Perplexity in language modeling acts as a proxy metric strongly linked to performance on diverse downstream tasks (see Figure 1). Since language modeling solely requires unlabeled natural text, it eliminates the need for expensive human annotation. Constructing and revising language modeling test sets is low-cost, as new data can be readily sampled from newly published content. Additionally, if a test set becomes compromised, fresh test data can quickly be sampled as a replacement.

4.3.2 Construction of diverse LM testsets

We compare the language modeling capabilities of various language models with our Skywork-13B, focusing on Chinese language.

To conduct a robust evaluation of language modeling capability, we have separately collected a diverse corpus of texts from a myriad of websites, each labeled according to its respective domain. The domains we cover span a wide spectrum, encompassing areas such as technology, movies, finance, to name a few. These domain-specific evaluation datasets have also been open-sourced for public access⁴.

⁴Github: https://github.com/SkyworkAI/Skywork/tree/main/data/eval_loss

We ensure that every test sample consists of documents or user posts published *after* September 1, 2023. This cut-off date guarantees that no test sample was inadvertently included during the pre-training of any evaluated language model. Specifically, SkyPile’s cut-off date is June 30, 2023, and the majority of models under evaluation were released prior to August 31.

Note that while the held-out validation set used to monitor the training progress (as shown in Figure 3) of our model can also serve this purpose, it has the same distribution (web texts) as the bulk of the training corpus, thus may lead to overly optimistic estimate of the actual language modeling capability of the model. More details on the sources of the test samples and the underlying data collection pipeline can be found in Appendix D.

4.3.3 Results

The results of our language modeling evaluation are presented in Table 6, where results from ChatGLM3-6B (THUDM, 2023), MOSS-7B (Sun and Qiu, 2023), Baichuan2-7B (Baichuan Inc., 2023), Qwen-7B (Qwen Team, 2023), InternLM-7B (InternLM Team, 2023) and Aquilla2-34B are also included.

It can be seen that our Skywork-13B model shows the best performance overall, obtaining the lowest average perplexity score of 9.42. It also exhibits the best performance across individual domains, achieving the lowest perplexity scores in tech (11.58), movie (21.84), government (4.76), and finance (4.92) domains. It excels not only in surpassing the performance of models of a similar size, but also in outperforming significantly larger models such as InternLM-20B and Aquila2-34B.

We attribute the excellent language modeling performance of our Skywork-13B to the quality of our training corpus. Details on rigorous data filtering pipeline are described in Section 3.1.

5 Discussion

In this section, we delve into the benefits and associated risks of pre-training on the in-domain data⁵ of benchmark tasks.

⁵The term “in-domain data” is a vague one that refers to any data with distribution closely resembling to that of the task data. For instance, the training data of a task is trivially in-domain data for that task.

Model	CEVAL	CMMLU	MMLU	GSM8K
OpenLLaMA-13B	27.1	26.7	42.7	12.4
LLaMA-13B	35.5	31.2	46.9	17.8
LLaMA-2-13B	36.5	36.6	54.8	28.7
Baichuan-13B	52.4	55.3	51.6	26.6
Baichuan2-13B	58.1	<u>62.0</u>	59.2	52.8
XVERSE-13B	54.7	-	55.1	-
InternLM-20B	58.8	-	62.0	52.6
Skywork-13B	<u>60.6</u>	61.8	<u>62.1</u>	<u>55.8</u>

Table 5: Comparison of results on popular benchmarks. Best result in each column is underlined. It can be seen that our Skywork-13B consistently perform well across the different benchmarks, indicating its overall robustness.

	Tech	Movie	Gov.	Game	Finance	General	Average
ChatGLM3-6B	12.48	23.48	5.07	18.45	5.67	7.47	10.25
MOSS-7B	20.83	39.66	11.08	31.24	10.59	13.25	18.50
InternLM-7B	13.43	24.9	5.88	19.78	6.17	8.10	11.17
Qwen-7B	13.39	25.16	5.55	19.26	5.76	7.78	10.83
Baichuan2-7B	12.89	23.26	5.34	18.36	5.68	7.62	10.41
LLaMA2-13B	23.26	50.66	18.09	32.52	14.85	16.55	23.54
Xverse-13B	12.55	23.49	5.20	17.69	5.54	7.46	10.19
Baichuan-13B	12.38	22.46	5.21	17.59	5.42	7.37	10.03
Baichuan2-13B	12.14	21.85	5.05	17.15	5.35	7.24	9.81
Qwen-14B	11.90	22.43	4.89	<u>16.94</u>	5.24	7.03	9.67
InternLM-20B	12.34	22.06	5.75	17.45	5.73	7.78	10.34
Aquila2-34B	14.62	29.09	5.72	21.78	5.83	8.45	11.73
Skywork-13B	<u>11.58</u>	<u>21.84</u>	<u>4.76</u>	17.28	<u>4.92</u>	<u>6.82</u>	<u>9.42</u>

Table 6: Comparative analysis of language modeling capabilities across diverse domains. Performance is measured using perplexity (lower values is better). Underlined figures correspond to the best result in each column.

5.1 Effect of pre-training on in-domain data

Pre-trained language models, or foundation models, are intended to be used in transfer learning as a general purpose backbone. As a foundation model in itself has little usage other than sentence completion, the quality of a foundation model is typically evaluated in terms of its performance in those tasks. Apparently, when it comes to improve a foundation model’s quality as measured by its task performance, it is always far more efficient to train the model on in-domain data of that task (Hernandez et al., 2021; Chung et al., 2022), as

GPT-4 generated data with few-shot task examples can also be considered as in-domain data for that task.

compared to general-purpose data (web texts).

We have shown that Stage-2 pre-training significantly amplifies our Skywork-13B’s STEM related capabilities, leading to a substantial improvement in performance on STEM-related tasks. Now we show that it is even possible to enhance a much weaker base model, i.e., an intermediate checkpoint, using only a fraction of the data and compute used in Stage-2 pre-training.

Table 7 presents the CEVAL and GSM8K scores before and after pre-training on in-domain data, utilizing a relatively weak model checkpoint that has only undergone 0.5T pre-training. The results indicate that after pre-training with merely 1B tokens of in-domain

	CEVAL	GSM8K	En Loss	Zh Loss
Before	28.3	6.9	1.86	2.08
After	50.8	40.7	2.09	2.21
Δ	+22.5	+33.8	+0.23	+0.13

Table 7: The impact of pre-training on a 0.5T checkpoint of Skywork-13B using only 1B tokens. The training data is sourced from a subset of our SkyPile-STEM corpus. The columns “En Loss” and “Zh Loss” show the model’s validation loss on held-out sets of English and Chinese web texts, respectively.

data, a weak model, initially performing only slightly better than random at CEVAL and GSM8K, can surpass the performance of our strongest Skywork-13B (3T) backbone without in-domain pre-training. However, this comes at the cost of significant degradation in language modeling performance, as evidenced by the higher loss on both tasks, shown in the two rightmost columns of the table.

5.2 Pre-training on in-domain data: a common practice?

It is of interest to explore whether popular foundational models are pre-trained on in-domain data. In pursuit of this, we delve into the GSM8K datasets, equipped with official train/test splits and comprehensive solutions. We evaluate an LLM’s language modeling loss on three datasets drawn from the same distribution: 1) The official GSM8K training set, 2) The official GSM8K test set, 3) A set composed of GSM8K-like samples generated by GPT-4. The corresponding losses are denoted as L_{train} , L_{test} , and L_{ref} , respectively. Theoretically, if a language model has not been exposed to any of the three datasets during pre-training, the three losses L_{train} , L_{test} , and L_{ref} should be approximately equivalent. However, if the model has been pre-trained on the training set or if the test data has been inadvertently exposed during the pre-training process, we would anticipate a notable discrepancy between L_{train} , L_{test} , and L_{ref} .

Our results are outlined in Table 8, which also reports the differences in losses $\Delta_1 = L_{test} - L_{ref}$ and $\Delta_2 = L_{test} - L_{train}$. Notably, the Δ_2 column reveals that for most models, the language modeling loss on the GSM8K training and test splits are almost iden-

tical. However, models such as ChatGLM3-6B, Baichuan2-13B, Qwen-7B/14B, and Aquila2-34B display markedly lower loss on the training split than on the test split. Consequently, we postulate that these models may have been considerably pre-trained on GSM8K training split or similar data.

Moreover, we notice one particular anomaly in the Δ_1 column, indicating the significantly lower L_{test} loss compared to L_{ref} , which is interesting to further study for better understanding.

5.3 Pre-Training or Supervised Fine-Tuning?

In the era preceding the advent of LLMs such as GPT-4 (Bubeck et al., 2023; OpenAI, 2023) and Claude (Bai et al., 2022), supervised data for NLP tasks was generally scarce. This was because the process of data collection and annotation was both time-consuming and costly. Due to the scarcity of supervised data, NLP researchers rely on unsupervised pre-training techniques (Mikolov et al., 2013; Peters et al., 2018; Radford et al., 2018; Devlin et al., 2019) to improve downstream task performance via transfer learning, where supervised data is to be used only in the fine-tuning stage. In this context, pre-training on in-domain (supervised) data was pointless, as it would defeat the purpose of pre-training itself (transfer learning).

This reality has significantly shifted, however, with the emergence of powerful LLMs. This is because procuring large amounts of high quality supervised/in-domain data is now as simple as making a few API requests to these LLMs, and it is comparatively low-cost (Wang et al., 2023; Taori et al., 2023). This new reality blurs the boundary between pre-training and supervised fine-tuning, making it feasible to incorporate substantial amounts of supervised data into the pre-training phase (Gunasekar et al., 2023; Li et al., 2023b). After all, curated in-domain data, whether written by human annotators or generated by LLM, are all form of human knowledge, and there is good reason for this knowledge to be absorbed into a foundation model.

That said, we believe that there is valid risk on the practice of targeted pre-training, in that it compromise fairness in benchmarking. While through pre-training on in-domain data a model

	L_{test}	L_{train}	L_{ref}	Δ_1	Δ_2
ChatGLM3-6B	0.99	0.78	0.99	0.0	0.21
MOSS-7B	1.51	1.52	1.49	0.02	-0.01
InternLM-7B	1.21	1.12	1.27	-0.06	0.09
Qwen-7B	1.07	0.64	1.10	-0.03	0.43
Baichuan2-7B	1.41	1.42	1.36	0.05	-0.01
LLaMA-13B	1.41	1.42	1.36	0.05	-0.01
LLaMA2-13B	1.36	1.38	1.33	0.03	-0.01
Xverse-13B	1.42	1.43	1.39	0.03	-0.01
Baichuan-13B	1.41	1.42	1.37	0.04	-0.01
Baichuan2-13B	1.09	0.72	1.12	-0.03	0.37
Qwen-14B	1.03	0.42	1.14	-0.11	0.61
InternLM-20B	1.20	1.09	1.19	0.01	0.11
Aquila2-34B	0.78	0.39	1.29	-0.51	0.39
Skywork-13B	1.01	0.97	1.00	0.01	0.04

Table 8: We evaluate the language modeling (LM) loss on samples (a sample is a concatenation of question and answer) from GSM8K dataset for several foundation models. For each LLM, we compare LM loss on the training split (L_{train}), the test split (L_{test}), and a specially curated reference set (L_{ref}), generated by GPT-4, designed to mimic the GSM8K dataset. We also reports two key metrics: $\Delta_1 = L_{test} - L_{ref}$, serving as an indicator of potential test data leakage during the training of the LLM, i.e., a lower value suggests possible leakage; and $\Delta_2 = L_{test} - L_{train}$, which measures the degree of overfitting on the training split of the dataset. A higher value of Δ_2 implies excessive overfitting. Outliers for both Δ_1 and Δ_2 are highlighted in gray.

may excel at specific tasks, it remains uncertain how well it would perform on unseen tasks. Its capabilities may be overestimated based on the benchmark alone, which can lead to unfair comparisons between models and mislead users or stakeholders about the true capabilities of the model.

6 Limitation

Our pre-training approach for Skywork-13B involved a two-stage process: general purpose pre-training followed by domain-specific enhancement pre-training. However, it remains unclear whether this methodology can produce a model on par with, or superior to, a model trained in one stage on a mixed corpus. Further investigation is needed to determine the comparative effectiveness of these pre-training approaches.

Additionally, we have proposed using language modeling loss or perplexity as proxy metrics for monitoring and evaluating large language models. A limitation is that language modeling evaluation relies on the specific distribution used to sample test data, of which there are infinite possibilities. While language mod-

eling perplexity over a given data distribution may predict performance on some tasks, it may not translate to other tasks. The correlation between language modeling and downstream performance could vary across different distributions and tasks.

7 Conclusion

Our work on Skywork-13B represents a significant leap forward in the development of open large language models. We believe that our comprehensive and transparent approach to the model’s development will be a valuable resource for researchers in the field, fostering collaboration and open-source principles. Our two-stage training methodology, leveraging a segmented corpus, offers a novel approach for enhancing model capability in specific domain, while our method of monitoring the training progress provides a practical solution to the challenges of tracking the improvement of these models over time.

However, our work is more than just the creation of a new LLM. It is a call to action for the broader NLP community, urging a return to

the principles of fairness, transparency, and the sharing of ideas that have historically fueled progress in the field. We hope that Skywork-13B will not only serve as a powerful tool for a wide range of applications but also inspire a renewed commitment to openness and cooperation in the development of future models.

References

- Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S. Morcos. 2023. [Semdedup: Data-efficient learning at web-scale through semantic deduplication](#).
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. 2022. [Training a helpful and harmless assistant with reinforcement learning from human feedback](#).
- Baichuan Inc. 2023. [Baichuan 2: Open large-scale language models](#). https://github.com/baichuan-inc/Baichuan2/blob/main/README_EN.md.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2019. [Piqa: Reasoning about physical commonsense in natural language](#).
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. [Sparks of artificial general intelligence: Early experiments with gpt-4](#).
- Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. 2023. [Extending context window of large language models via positional interpolation](#).
- Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. 2016. [Training deep nets with sublinear memory cost](#).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. [Palm: Scaling language modeling with pathways](#). *arXiv preprint arXiv:2204.02311*.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. [Scaling instruction-finetuned language models](#).
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. [Boolq: Exploring the surprising difficulty of natural yes/no questions](#).
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? try arc, the ai2 reasoning challenge](#). *arXiv preprint arXiv:1803.05457*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#).
- Tri Dao. 2023. [Flashattention-2: Faster attention with better parallelism and work partitioning](#).
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. [Flashattention: Fast and memory-efficient exact attention with io-awareness](#).
- Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christopher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, Marcus Hutter, and Joel Veness. 2023. [Language modeling is compression](#).
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno,

- Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. 2023. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#).
- Danny Hernandez, Tom Brown, Tom Conerly, Nova DasSarma, Dawn Drain, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Tom Henighan, Tristan Hume, Scott Johnston, Ben Mann, Chris Olah, Catherine Olsson, Dario Amodei, Nicholas Joseph, Jared Kaplan, and Sam McCandlish. 2022. [Scaling laws and interpretability of learning from repeated data](#).
- Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam McCandlish. 2021. [Scaling laws for transfer](#).
- 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](#).
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. 2023. [C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models](#). *arXiv preprint arXiv:2305.08322*.
- InternLM Team. 2023. [Internlm: A multilingual language model with progressively enhanced capabilities](#). <https://github.com/InternLM/InternLM>.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. [Deduplicating training data mitigates privacy risks in language models](#).
- 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](#).
- Vijay Korthikanti, Jared Casper, Sangkug Lym, Lawrence McAfee, Michael Andersch, Mohammad Shoeybi, and Bryan Catanzaro. 2022. [Reducing activation recomputation in large transformer models](#).
- Taku Kudo and John Richardson. 2018. [SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. [Race: Large-scale reading comprehension dataset from examinations](#). *arXiv preprint arXiv:1704.04683*.
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2022. [Deduplicating training data makes language models better](#).
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023a. [Cmmlu: Measuring massive multitask language understanding in chinese](#).
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023b. [Textbooks are all you need ii: phi-1.5 technical report](#). *arXiv preprint arXiv:2309.05463*.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#).
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. [Efficient estimation of word representations in vector space](#).
- Niklas Muennighoff, Alexander M. Rush, Boaz Barak, Teven Le Scao, Aleksandra Piktus, Noumane Tazi, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. 2023. [Scaling data-constrained language models](#).
- Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vijay Anand Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, Amar Phanishayee, and Matei Zaharia. 2021. [Efficient large-scale language model training on gpu clusters using megatron-lm](#).
- OpenAI. 2023. [GPT-4 technical report](#).
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#).

- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refined-web dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. **Deep contextualized word representations**. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Qwen Team. 2023. **QWEN technical report**. <https://github.com/QwenLM/Qwen>.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training.
- Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. 2020. **Zero: Memory optimizations toward training trillion parameter models**.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. **Code llama: Open foundation models for code**.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- Noam Shazeer. 2020. **Glu variants improve transformer**.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2020. **Megatron-lm: Training multi-billion parameter language models using model parallelism**.
- Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. 2023. **SlimPajama: A 627B token cleaned and deduplicated version of RedPajama**.
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. 2022. **Roformer: Enhanced transformer with rotary position embedding**.
- Tianxiang Sun and Xipeng Qiu. 2023. **MOSS**. https://github.com/OpenLMLab/MOSS/blob/main/README_en.md.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.
- Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. **Galactica: A large language model for science**.
- THUDM. 2023. **ChatGLM3-6B**. <https://github.com/THUDM/ChatGLM3> Webpage in Chinese.
- Together Computer. 2023. **Redpajama: An open source recipe to reproduce llama training dataset**.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. **Llama: Open and efficient foundation language models**.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. **Llama 2: Open foundation and fine-tuned chat models**. *arXiv preprint arXiv:2307.09288*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 6000–6010, Red Hook, NY, USA. Curran Associates Inc.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. **Self-instruct: Aligning language models with self-generated instructions**.
- Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. 2020. **CCNet: Extracting high quality monolingual datasets from web crawl data**. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4003–4012, Marseille, France. European Language Resources Association.
- Mengzhou Xia, Mikel Artetxe, Chunting Zhou, Xi Victoria Lin, Ramakanth Pasunuru, Danqi Chen, Luke Zettlemoyer, and Ves Stoyanov. 2023. **Training trajectories of language models across scales**.

Wenhan Xiong, Jingyu Liu, Igor Molybog, Hejia Zhang, Prajjwal Bhargava, Rui Hou, Louis Martin, Rashi Rungta, Karthik Abinav Sankararaman, Barlas Oğuz, Madian Khabza, Han Fang, Yashar Mehdad, Sharan Narang, Kshitiz Malik, Angela Fan, Shruti Bhosale, Sergey Edunov, Mike Lewis, Sinong Wang, and Hao Ma. 2023. [Effective long-context scaling of foundation models](#).

Xverse-AI. 2023. [Xverse-13B](#). <https://github.com/xverse-ai/XVERSE-13B> Webpage in Chinese.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.

Biao Zhang and Rico Sennrich. 2019. [Root Mean Square Layer Normalization](#). In *Advances in Neural Information Processing Systems 32*, Vancouver, Canada.

A Details on GPT-7B vs. LLaMA-7B Experiment

In a preliminary experiment, we compared the language modeling performance between GPT and LLaMA architecture in a controlled environment. We trained a 7B model with GPT architecture and a comparable 7B model with LLaMA architecture for 200B tokens sampled from the same corpus and with the same training parameters. Details are given in Table 9.

B Preliminary Experiments on Distributed Training

In Table 10 we report preliminary results obtained with various distributed training configurations on LLaMA-13B and Skywork-13B model architecture. In both cases, the best throughput is achieved with DP256 and PP2 with ZERO-1 setting.

C More Benchmark Results

We also provide results of the following benchmarks in Table 11:

- **TriviaQA** (Joshi et al., 2017): TriviaQA is a realistic text-based question answering dataset which includes 950K question-answer pairs from 662K documents collected from Wikipedia and the web.

- **HellaSwag** (Zellers et al., 2019): HellaSWAG is a dataset that focuses on grounded commonsense inference.
- **Winogrande** (Sakaguchi et al., 2021): Winogrande is a dataset that focuses on commonsense reasoning.
- **BoolQ** (Clark et al., 2019) BoolQ is a question answering dataset for yes/no questions.
- **PIQA** (Bisk et al., 2019): PIQA is a dataset for commonsense reasoning, and was created to investigate the physical knowledge of existing models in NLP.
- **ARC** (Clark et al., 2018): ARC is a dataset consisting of multiple-choice question-answering tasks that focus on commonsense reasoning.
- **RACE** (Lai et al., 2017) RACE is a dataset that focuses on reading comprehension.

D Details on LM Test Sets

We established a daily crawl of published articles and user posts from a selection of widely used Chinese websites. This data collection process is distinct from the pipeline utilized to construct SkyPile. The purpose of gathering this data is to create independent language modeling test sets, categorized by their domain, for the evaluation of current open Language Learning Models (LLMs).

Below we describe the sources of these domain testsets:

- **Technology**: AI related articles from ([36kr.com](#)). This website provides timely and comprehensive news articles about startups, technology, and business trends, primarily in the Chinese market.
- **Movie**: User written movie reviews from Douban ([douban.com](#)). Douban is a popular social networking service in China that offers a platform for users to share their opinions and create content related to movies, books, and music. It is one of the most influential web 2.0 websites in China and has a strong focus on user-generated content.
- **Government**: News from website of People’s Daily ([www.people.com.cn](#)), which is the

	GPT-7B	LLaMA-7B
Positional Embedding	Absolute	Rotary
Max Position Embeddings	4096	4096
Normalization	LayerNorm	RMSNorm
Activation	Gelu	SwiGlu
Attention	MHA	MHA
Num. Layers	32	32
Hidden Size	4096	4096
Num. Heads	32	32
FFN Size	16384	11008
Context Size	4096	4096
Global Batch Size	1024	1024
Adam β_1	0.95	0.95
Adam β_2	0.9	0.9
Adam ϵ	1.00e-8	1.00-8
Precision	bf16	bf16
Peak Learning Rate	3e-4	3e-4
Min Learning Rate	3e-5	3e-5
Learning Rate Decay Steps	43945	43945
Learning Rate Decay Style	Cosine	Cosine
Warm-up Steps	2000 steps	2000 steps
Weight Decay	0.1	0.1
Dropout Probability	0.1	0
Gradient Clip	1	1
Total Steps	51200	51200

Table 9: Comparison of GPT-7B and LLaMA-7B. All variables are controlled in our experiment except for the differences in architecture.

Model	Strategy	Throughput	MFU	TFlops	Memory
LLaMA2	DP512	-	-	-	OOM
LLaMA2	DP256+PP2	<u>2045</u>	<u>58.5</u>	<u>182.6</u>	<u>70.7</u>
LLaMA2	DP256+TP2	1928	55.2	172.2	65.5
LLaMA2	DP128+TP2+PP2	1936	55.4	172.9	39.4
LLaMA2	DP128+PP4	1964	56.2	175.4	53.4
LLaMA2	DP128+TP4	1744	44.4	138.5	35.4
Skywork	DP512	-	-	-	OOM
Skywork	DP256+PP2	<u>1873</u>	<u>56.5</u>	<u>176.2</u>	<u>77.1</u>
Skywork	DP256+TP2	1775	53.5	167.0	67.9
Skywork	DP128+TP2+PP2	1776	53.5	167.0	42.5
Skywork	DP128+PP4	1828	55.1	171.9	58.7
Skywork	DP128+TP4	1417	43.1	134.6	36.6

Table 10: Compute efficiency achieved with different distributed training configurations. We tested both LLaMA2-13B and Skywork-13B. Throughout the experiments, we use a global batch size of 4096 and a micro batch size of 1. When Tensor Parallelism is enabled, Sequence Parallelism is enabled as well. Throughput is measured in tokens processed per GPU per second, while Model Flops Utilization (MFU) is expressed as a percentage (%). Memory usage is reported in Gigabytes (GB).

Models	BoolQ	PIQA	Winogrande	TriviaQA	RACE	Hellaswag	ARC-E	ARC-C
OpenLLaMA-13B	77.6	79.5	72.0	60.2	42.4	76.0	78.9	48.6
LLaMA-13B	80.7	81.0	<u>76.2</u>	65.0	43.4	80.1	82.1	54.7
LLaMA2-13B	<u>83.3</u>	<u>81.7</u>	75.8	<u>68.2</u>	43.9	<u>81.5</u>	<u>83.7</u>	<u>57.0</u>
Baichuan-13B	78.8	77.2	70.4	51.6	35.8	74.2	77.2	48.4
Baichuan2-13B	80.3	79.3	72.1	58.0	25.2	76.4	81.1	53.2
Xverse-13B	79.8	80.0	71.1	53.3	43.2	77.2	78.5	49.1
Skywork-13B	82.9	79.9	72.2	54.0	<u>45.2</u>	77.4	78.5	50.2

Table 11: More English benchmarks results. As all of these models are more or less sensitive to the prompt template or number of shots, the reported results, which are reproduced by us, may be different to those from other sources.

most influential and authoritative newspapers in China. The language used in the news is typically formal Standard Mandarin and carries an authoritative tone.

- **Game:** Articles from Gcores (www.gcores.com). This is a Chinese digital media platform dedicated to video games, tech trends, and geek culture. The platform features a wide range of original content, including news articles, podcast episodes, videos, and independent games.
- **Finance:** News from finance section of Sina (finance.sina.com.cn). It is one of China’s leading online media companies, offers a comprehensive suite of financial information and services. It covers a broad range of topics including stock markets, forex, commodities, real estate, and personal finance.
- **General:** News from Jiemian News (www.jiemian.com). Jiemian is a prominent Chinese digital media platform known for its in-depth and high-quality journalism. It covers a wide range of topics, including politics, economy, culture, technology, finance, and lifestyle.

Subject	Stage-1	Stage-2	Boost
Accountant	40.8	49.0	8.2
Advanced Mathematics	26.3	42.1	15.8
Art Studies	60.6	72.7	12.1
Basic Medicine	42.1	57.9	15.8
Business Administration	42.4	48.5	6.1
Chinese Language and Literature	47.8	56.5	8.7
Civil Servant	40.4	66.0	25.5
Clinical Medicine	36.4	40.9	4.5
College Chemistry	37.5	50.0	12.5
College Economics	52.7	47.3	-5.5
College Physics	15.8	36.8	21.1
College Programming	51.4	51.4	0.0
Computer Architecture	33.3	52.4	19.0
Computer Network	21.1	26.3	5.3
Discrete Mathematics	50.0	18.8	-31.3
Education Science	44.8	75.9	31.0
Electrical Engineer	35.1	35.1	0.0
Environmental Impact Assessment Engineer	45.2	51.6	6.5
Fire Engineer	45.2	51.6	6.5
High School Biology	42.1	78.9	36.8
High School Chemistry	36.8	63.2	26.3
High School Chinese	26.3	42.1	15.8
High School Geography	36.8	78.9	42.1
High School History	80.0	80.0	0.0
High School Mathematics	27.8	16.7	-11.1
High School Physics	42.1	57.9	15.8
High School Politics	47.4	84.2	36.8
Ideological and Moral Cultivation	84.2	100.0	15.8
Law	33.3	45.8	12.5
Legal Professional	39.1	52.2	13.0
Logic	50.0	45.5	-4.5
Mao Zedong Thought	70.8	83.3	12.5
Marxism	57.9	63.2	5.3
Metrology Engineer	37.5	58.3	20.8
Middle School Biology	76.2	95.2	19.0
Middle School Chemistry	30.0	95.0	65.0
Middle School Geography	41.7	83.3	41.7
Middle School History	59.1	81.8	22.7
Middle School Mathematics	15.8	36.8	21.1
Middle School Physics	42.1	73.7	31.6
Middle School Politics	52.4	90.5	38.1
Modern Chinese History	47.8	73.9	26.1
Operating System	52.6	47.4	-5.3
Physician	46.9	57.1	10.2
Plant Protection	63.6	63.6	0.0
Probability and Statistics	27.8	33.3	5.6
Professional Tour Guide	69.0	65.5	-3.4
Sports Science	42.1	52.6	10.5
Tax Accountant	30.6	49.0	18.4
Teacher Qualification	61.4	84.1	22.7
Urban and Rural Planner	50	67.4	17.4
Veterinary Medicine	26.1	60.9	34.8

Table 12: Details on CEVAL benchmark results.

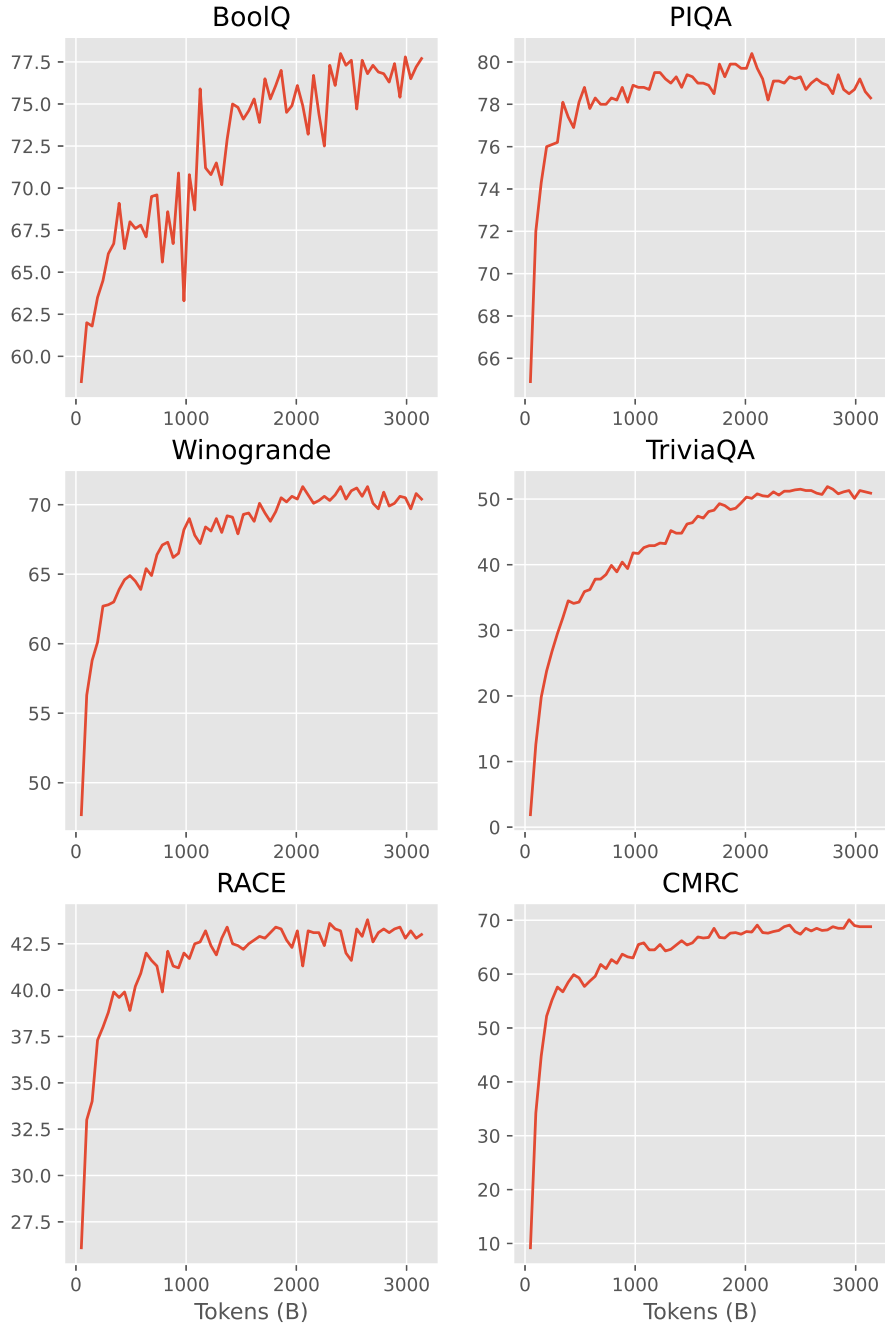


Figure 6: Performance of the Skywork-13B on various benchmarks during Stage-1 pre-training. Benchmarks include BoolQ, PIQA, Winogrande, TriviaQA, RACE, and CMRC.