

# DeepSeek LLM

## Scaling Open-Source Language Models with Longtermism

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, Huazuo Gao, Kaige Gao, Wenjun Gao, Ruiqi Ge, Kang Guan, Daya Guo, Jianzhong Guo, Guangbo Hao, Zhewen Hao, Ying He, Wenjie Hu, Panpan Huang, Erhang Li, Guowei Li, Jiashi Li, Yao Li, Y.K. Li, Wenfeng Liang, Fangyun Lin, A.X. Liu, Bo Liu, Wen Liu, Xiaodong Liu, Xin Liu, Yiyuan Liu, Haoyu Lu, Shanghao Lu, Fuli Luo, Shirong Ma, Xiaotao Nie, Tian Pei, Yishi Piao, Junjie Qiu, Hui Qu, Tongzheng Ren, Zehui Ren, Chong Ruan, Zhangli Sha, Zhihong Shao, Junxiao Song, Xuecheng Su, Jingxiang Sun, Yaofeng Sun, Minghui Tang, Bingxuan Wang, Peiyi Wang, Shiyu Wang, Yaohui Wang, Yongji Wang, Tong Wu, Y. Wu, Xin Xie, Zhenda Xie, Ziwei Xie, Yiliang Xiong, Hanwei Xu, R.X. Xu, Yanhong Xu, Dejian Yang, Yuxiang You, Shuiping Yu, Xingkai Yu, B. Zhang, Haowei Zhang, Lecong Zhang, Liyue Zhang, Mingchuan Zhang, Minghua Zhang, Wentao Zhang, Yichao Zhang, Chenggang Zhao, Yao Zhao, Shangyan Zhou, Shunfeng Zhou, Qihao Zhu, Yuheng Zou \*

\*DeepSeek-AI

### Abstract

The rapid development of open-source large language models (LLMs) has been truly remarkable. However, the scaling laws described in previous literature presents varying conclusions, which casts a dark cloud over scaling LLMs. We delve into the study of scaling laws and present our distinctive findings that facilitate the scaling of large scale models in two prevalent used open-source configurations, 7B and 67B. Guided by the scaling laws, we introduce DeepSeek LLM, a project dedicated to advancing open-source language models with a long-term perspective. To support the pre-training phase, we have developed a dataset that currently consists of 2 trillion tokens and is continuously expanding. We further conduct supervised fine-tuning (SFT) and direct preference optimization (DPO) on DeepSeek LLM Base models, resulting in the creation of DeepSeek Chat models. Our evaluation results demonstrate that DeepSeek LLM 67B surpasses LLaMA-2 70B across a range of benchmarks, especially in the domains of code, mathematics, and reasoning. Furthermore, open-ended evaluations reveal that our DeepSeek LLM 67B Chat exhibits superior performance compared to GPT-3.5.

---

\*Authors are ordered alphabetically by the last name.

**Contents**

- 1 Introduction** **3**
  
- 2 Pre-Training** **4**
  - 2.1 Data . . . . . 4
  - 2.2 Architecture . . . . . 5
  - 2.3 Hyperparameters . . . . . 5
  - 2.4 Infrastructures . . . . . 6
  
- 3 Scaling Laws** **7**
  - 3.1 Scaling Laws for Hyperparameters . . . . . 8
  - 3.2 Estimating Optimal Model and Data Scaling . . . . . 9
  - 3.3 Scaling Laws with Different Data . . . . . 12
  
- 4 Alignment** **12**
  
- 5 Evaluation** **13**
  - 5.1 Public Benchmark Evaluation . . . . . 13
    - 5.1.1 Base Model . . . . . 14
    - 5.1.2 Chat Model . . . . . 14
  - 5.2 Open-Ended Evaluation . . . . . 17
    - 5.2.1 Chinese Open-Ended Evaluation . . . . . 17
    - 5.2.2 English Open-Ended Evaluation . . . . . 18
  - 5.3 Held-Out Evaluation . . . . . 18
  - 5.4 Safety Evaluation . . . . . 19
  - 5.5 Discussion . . . . . 20
  
- 6 Conclusion, Limitation, and Future Work** **23**
  
- A Appendix** **30**
  - A.1 Acknowledgments . . . . . 30
  - A.2 Different Model Scale Representations . . . . . 30
  - A.3 Benchmark Metrics Curves . . . . . 31
  - A.4 Comparison with Code or Math Specific Models . . . . . 32
  - A.5 Benchmark Results w/ DPO Stage . . . . . 32
  - A.6 Evaluation Formats . . . . . 32

# 1. Introduction

Over the past few years, Large Language Models (LLMs) based on decoder-only Transformers (Vaswani et al., 2017) have increasingly become the cornerstone and pathway to achieving Artificial General Intelligence (AGI). By predicting the next word in continuous text, LLMs undergo self-supervised pre-training on massive datasets, enabling them to achieve various purposes and possess many abilities, such as novel creation, text summarization, code completion, and more. Subsequent developments like supervised fine-tuning and reward modeling have enabled Large Language Models (LLMs) to better follow user intentions and instructions. This has endowed them with more versatile conversational capabilities and rapidly expanded their influence.

This wave is sparked with closed products, such as ChatGPT (OpenAI, 2022), Claude (Anthropic, 2023), and Bard (Google, 2023), which are developed with extensive computational resources and substantial annotation costs. These products have significantly raised the community’s expectations for the capabilities of open-source LLMs, consequently inspiring a series of work (Bai et al., 2023; Du et al., 2022; Jiang et al., 2023; Touvron et al., 2023a,b; Yang et al., 2023). Among these, the LLaMA series models (Touvron et al., 2023a,b) stand out. It consolidates a range of works to create an efficient and stable architecture, building well-performing models ranging from 7B to 70B parameters. Consequently, the LLaMA series has become the de facto benchmark for architecture and performance among open-source models.

Following LLaMA, the open-source community has primarily focused on training fixed-size (7B, 13B, 34B, and 70B), high-quality models, often neglecting research exploration into LLM scaling laws (Hoffmann et al., 2022; Kaplan et al., 2020). Nonetheless, research on scaling laws is of utmost importance, considering that the current open-source models are merely at the initial stage of Artificial General Intelligence (AGI) development. In addition, early works (Hoffmann et al., 2022; Kaplan et al., 2020) reached varying conclusions on the scaling of model and data with increased compute budgets and inadequately addressed hyperparameter discussions. In this paper, we extensively investigate the scaling behavior of language models and apply our findings in two widely used large-scale model configurations, namely 7B and 67B. Our study aims to lay the groundwork for future scaling of open-source LLMs, paving the way for further advancements in this domain. Specifically, we first examined the scaling laws of batch size and learning rate, and found their trends with model size. Building on this, we conducted a comprehensive study of the scaling laws of the data and model scale, successfully revealing the optimal model/data scaling-up allocation strategy and predicting the expected performance of our large-scale models. Additionally, during development, we discovered that the scaling laws derived from different datasets show significant differences. This suggests that choice of dataset remarkably affects the scaling behavior, indicating that caution should be exercised when generalizing scaling laws across datasets.

Under the guidance of our scaling laws, we build from scratch open-source large language models, and release as much information as possible for community reference. We collect 2 trillion tokens for pre-training, primarily in Chinese and English. At the model level, we generally followed the architecture of LLaMA, but replaced the cosine learning rate scheduler with a multi-step learning rate scheduler, maintaining performance while facilitating continual training. We collected over 1 million instances for supervised fine-tuning (SFT) (Ouyang et al., 2022) from diverse sources. This paper shares our experiences with different SFT strategies and findings in data ablation techniques. Additionally, we have utilized direct preference optimization (DPO) (Rafailov et al., 2023) to improve the conversational performance of the model.

We conduct extensive evaluations using our base and chat models. The evaluation results demonstrate that DeepSeek LLM surpasses LLaMA-2 70B across various benchmarks, particularly in the fields of code, mathematics, and reasoning. Following SFT and DPO, the DeepSeek 67B chat model outperforms GPT-3.5 in both Chinese and English open-ended evaluations. This highlights the superior performance of DeepSeek 67B in generating high-quality responses and engaging in meaningful conversations in both languages. Furthermore, the safety evaluation indicates that DeepSeek 67B Chat can provide harmless responses in practice.

In the rest of this paper, we first introduce our pre-training basic concepts of DeepSeek LLM in Section 2, including the composition of data, model architecture, infrastructure, and hyperparameters. In Section 3, we provide a detailed explanation of the scaling laws we have discovered and its implications. Additionally, we discuss the rationale behind our selection of pre-training hyperparameters, taking into account the insights gained from the scaling laws analysis. In Section 4, we discuss our fine-tuning methodology, encompassing the composition of fine-tuning data and specific methods during the SFT and DPO stages. We then present the detailed evaluation results of DeepSeek LLM in Section 5, covering both the base and chat models, as well as their performance in open-ended evaluations and safety evaluations. Finally, we discuss the current limitations and future directions of DeepSeek LLM in Section 6.

## 2. Pre-Training

### 2.1. Data

Our main objective is to comprehensively enhance the richness and diversity of the dataset. We have gained valuable insights from reputable sources such as (Computer, 2023; Gao et al., 2020; Penedo et al., 2023; Touvron et al., 2023a). To achieve these goals, we have organized our approach into three essential stages: deduplication, filtering, and remixing. The deduplication and remixing stages ensure a diverse representation of the data by sampling unique instances. The filtering stage enhances the density of information, thereby enabling more efficient and effective model training.

We adopted an aggressive deduplication strategy, expanding the deduplication scope. Our analysis revealed that deduplicating the entire Common Crawl corpus results in higher removal of duplicate instances compared to deduplicating within a single dump. Table 1 illustrates that deduplicating across 91 dumps eliminates four times more documents than a single dump method.

Dumps Used	1	2	6	12	16	22	41	91
Deduplication Rate (%)	22.2	46.7	55.7	69.9	75.7	76.3	81.6	89.8

Table 1 | Deduplication ratios for various Common Crawl dumps.

In the filtering stage, we focus on developing robust criteria for document quality assessment. This involves a detailed analysis incorporating both linguistic and semantic evaluations, providing a view of data quality from individual and global perspectives. In the remixing phase, we adjust our approach to address data imbalances, focusing on increasing the presence of underrepresented domains. This adjustment aims to achieve a more balanced and inclusive dataset, ensuring that diverse perspectives and information are adequately represented.

For our tokenizer, we implemented the Byte-level Byte-Pair Encoding (BBPE) algorithm based on the tokenizers library (Huggingface Team, 2019). Pre-tokenization was employed to

prevent the merging of tokens from different character categories such as new lines, punctuation, and Chinese-Japanese-Korean (CJK) symbols, similar to GPT-2 (Radford et al., 2019). We also chose to split numbers into individual digits following the approach used in (Touvron et al., 2023a,b). Based on our prior experience, we set the number of conventional tokens in the vocabulary at 100000. The tokenizer was trained on a multilingual corpus of approximately 24 GB, and we augmented the final vocabulary with 15 special tokens, bringing the total size to 100015. To ensure computational efficiency during training and to reserve space for any additional special tokens that might be needed in the future, we configured the model’s vocabulary size to 102400 for training.

## 2.2. Architecture

Params	$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{heads}}$	$n_{\text{kv\_heads}}$	Context Length	Sequence Batch Size	Learning Rate	Tokens
7B	30	4096	32	32	4096	2304	4.2e-4	2.0T
67B	95	8192	64	8	4096	4608	3.2e-4	2.0T

Table 2 | Detailed specs of DeepSeek LLM family of models. We choose the hyper-parameters based on our findings in Section 3

The micro design of DeepSeek LLM largely follows the design of LLaMA (Touvron et al., 2023a,b), adopting a Pre-Norm structure with RMSNorm (Zhang and Sennrich, 2019) function and using SwiGLU (Shazeer, 2020) as the activation function for the Feed-Forward Network (FFN), with an intermediate layer dimension of  $\frac{8}{3}d_{\text{model}}$ . It also incorporates Rotary Embedding (Su et al., 2024) for positional encoding. To optimize inference cost, the 67B model uses Grouped-Query Attention (GQA) (Ainslie et al., 2023) instead of the traditional Multi-Head Attention (MHA).

However, in terms of macro design, DeepSeek LLM differs slightly. Specifically, DeepSeek LLM 7B is a 30-layer network, while DeepSeek LLM 67B has 95 layers. These layer adjustments, while maintaining parameter consistency with other open-source models, also facilitate model pipeline partitioning to optimize training and inference.

Unlike most works using Grouped-Query Attention (GQA), we expanded the 67B model’s parameters in network depth rather than the common practice of widening the intermediate width of FFN layers, aiming for better performance. Detailed network specifications can be found in Table 2.

## 2.3. Hyperparameters

DeepSeek LLM is initialized with a standard deviation of 0.006 and trained using the AdamW optimizer (Loshchilov and Hutter, 2017), with the following hyperparameters:  $\beta_1 = 0.9$ ,  $\beta_2 = 0.95$ , and  $\text{weight\_decay} = 0.1$ .

A multi-step learning rate scheduler is employed during pre-training instead of the typical cosine scheduler. Specifically, the learning rate of the model reaches its maximum value after 2000 warmup steps, and then decreases to 31.6% of the maximum value after processing 80% of the training tokens. It further reduces to 10% of the maximum value after 90% of the tokens. The gradient clipping during the training phase is set to 1.0.

Based on our empirical findings, we observed that despite differences in the loss reduction

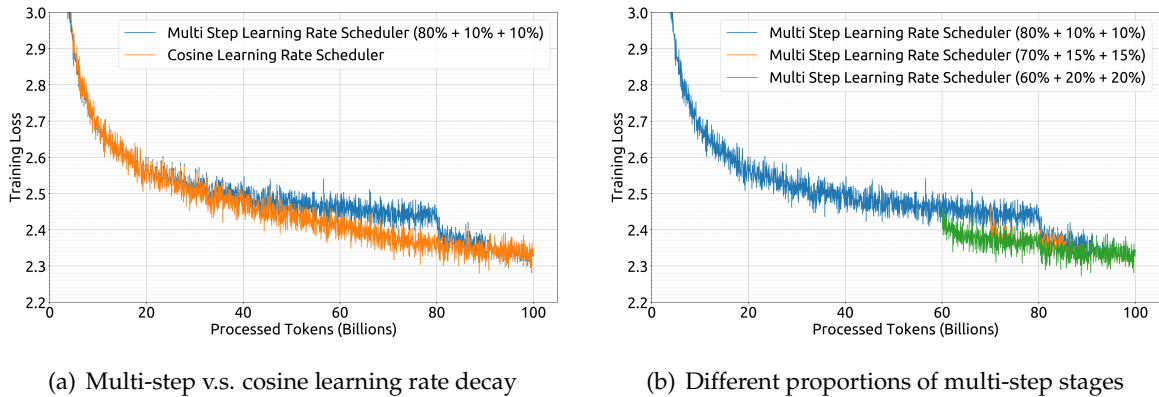


Figure 1 | Training loss curves with different learning rate schedulers or different parameters for schedulers. The model size is 1.6 billion parameters, trained on a dataset of 100 billion tokens.

trend during training, the final performance using a multi-step learning rate scheduler is essentially consistent with that of a cosine scheduler, as shown in Figure 1(a). When adjusting the training scale while keeping the model size fixed, the multi-step learning rate scheduler allows for the reuse of training from the first phase, offering a unique convenience for continual training. Therefore, we chose the multi-step learning rate scheduler as our default setting. We also demonstrate in Figure 1(b) that adjusting the proportions of different stages in the multi-step learning rate scheduler can yield slightly better performance. However, for the sake of balancing reuse ratios in continual training and model performance, we opted for the aforementioned distribution of 80%, 10%, and 10% for the three stages respectively.

The batch size and learning rate vary with the model size. Specific parameters for the pre-training phases of the 7B and 67B models can be found in Table 2.

## 2.4. Infrastructures

We use an efficient and light-weight training framework named HAI-LLM (High-flyer, 2023) to train and evaluate large language models. Data parallelism, tensor parallelism, sequence parallelism, and 1F1B pipeline parallelism are integrated into this framework as done in Megatron (Korthikanti et al., 2023; Narayanan et al., 2021; Shoeybi et al., 2019). We also leverage the flash attention (Dao, 2023; Dao et al., 2022) technique to improve hardware utilization. ZeRO-1 (Rajbhandari et al., 2020) is exploited to partition optimizer states over data parallel ranks. Efforts are also made to overlap computation and communication to minimize additional waiting overhead, including the backward procedure of the last micro-batch and reduce-scatter operation in ZeRO-1, and GEMM computation and all-gather/reduce-scatter in sequence parallel. Some layers/operators are fused to speed up training, including LayerNorm, GEMM whenever possible, and Adam updates. To improve model training stability, we train the model in bf16 precision but accumulate gradients in fp32 precision. In-place cross-entropy is performed to reduce GPU memory consumption, i.e.: we convert bf16 logits to fp32 precision on the fly in the cross-entropy CUDA kernel (instead of converting it beforehand in HBM), calculate the corresponding bf16 gradient, and overwrite logits with its gradient.

Model weights and optimizer states are saved every 5 minutes asynchronously, which means we will lose no more than 5 minutes of training in the worst case of occasional hardware or network failures. These temporary model checkpoints are cleared up regularly to avoid

consuming too much storage space. We also support resuming training from a different 3D parallel configuration to cope with dynamic changes in computing cluster load.

As for evaluation, we employ vLLM (Kwon et al., 2023) in generative tasks, and continuous batching in non-generative tasks to avoid manual batch size tuning and reduce token padding.

### 3. Scaling Laws

Research on scaling laws (Hestness et al., 2017) predates the emergence of large language models. Scaling laws (Henighan et al., 2020; Hoffmann et al., 2022; Kaplan et al., 2020) suggest that model performance can be predictably improved with increases in compute budget  $C$ , model scale  $N$ , and data scale  $D$ . When model scale  $N$  is represented by model parameters and data scale  $D$  by the number of tokens,  $C$  can be approximated as  $C = 6ND$ . Therefore, how to optimize the allocation between model and data scales when increasing the compute budget is also a crucial research objective in scaling laws.

The development of LLMs (Dai et al., 2019; Radford et al., 2019), with larger models achieving unexpected and significant performance improvements, has brought scaling laws research to a new peak. Results in scaling laws demonstrate that expanding the compute budget continues to yield significant benefits, which further encourages the increase in model scales (Brown et al., 2020; Smith et al., 2022).

However, as shown in Table 4, early works (Hoffmann et al., 2022; Kaplan et al., 2020) on the optimal model/data scaling-up allocation strategy have shown varying conclusions, raising doubts about the general applicability of scaling laws. Moreover, these studies often lacked a complete description of hyperparameter settings, leaving it uncertain whether models under different compute budgets reached optimal performance. Therefore, we revisit scaling laws in this section to address these uncertainties and ensure we are on the right path to efficiently scale-up compute, which reflects the long-term perspective and is key to developing continuously improving models.

To ensure that models under different compute budgets can achieve optimal performance, we first studied the scaling laws of hyperparameters. Empirically, it has been observed that the optimal values of most parameters during training do not change when varying compute budgets. Therefore, these parameters are consistent with those outlined in Section 2.3 and remain unchanged across different compute budgets. However, the hyperparameters that have the most significant impact on performance, namely batch size and learning rate, were re-examined.

Early works (Goyal et al., 2017; McCandlish et al., 2018; Shallue et al., 2019; Smith et al., 2017; Zhang et al., 2019) provided some empirical observations for setting batch size and learning rate, but we found these observations have limited applicability in our preliminary experiments. Through extensive experiments, we modeled the power law relationship between the compute budget  $C$  and the optimal batch size and learning rate. This relationship, which we refer to as the scaling laws of hyperparameters, provides an empirical framework for determining the optimal hyperparameters. This methodology ensures that models across different compute budgets can reach their near-optimal performance.

We then study the scaling laws of the model and data scales. To reduce experimental costs and fitting difficulties, we adopted the IsoFLOP profile approach from Chinchilla (Hoffmann et al., 2022) to fit the scaling curve. To represent the model scale more accurately, we utilized a new model scale representation, non-embedding FLOPs/token  $M$ , replacing the earlier-used model parameters  $N$ , and substituted the approximate compute budget formula  $C = 6ND$

with the more precise  $C = MD$ . The experimental results provided insights into the optimal model/data scaling-up allocation strategy and performance predictions, and also accurately forecasted the expected performance of DeepSeek LLM 7B and 67B models.

Additionally, in the process of exploring scaling laws, the data we used underwent multiple iterations, continually improving in quality. We attempted to fit the scaling curve on various datasets and found that the data quality significantly influences the optimal model/data scaling-up allocation strategy. The higher the data quality, the more the increased compute budget should be allocated to model scaling. This implies that high-quality data can drive the training of larger models given the same data scale. The differences in the optimal model/data scaling-up allocation strategy may also serve as an indirect approach to assess the quality of data. We will continue to pay close attention to the changes in data quality and its impact on scaling laws, and provide more analysis in future works.

In summary, our contributions and findings in scaling laws can be summarized as follows:

- We established the scaling laws for hyperparameters, providing an empirical framework for determining the optimal hyperparameters.
- Instead of model parameters  $N$ , we adopt non-embedding FLOPs/token  $M$  to represent the model scale, leading to a more accurate optimal model/data scaling-up allocation strategy and a better prediction of generalization loss for large-scale models.
- The quality of pre-training data impacts the optimal model/data scaling-up allocation strategy. The higher the data quality, the more the increased compute budget should be allocated to model scaling.

### 3.1. Scaling Laws for Hyperparameters

We initially conducted a grid search for batch size and learning rate on small-scale experiments with a compute budget of  $1e17$ , and the results of a specific model size (177M FLOPs/token) are illustrated in Figure 2(a). The results demonstrate that the generalization error remains stable across a wide range of choices of batch sizes and learning rates. This indicates that near-optimal performance can be achieved within a relatively wide parameter space.

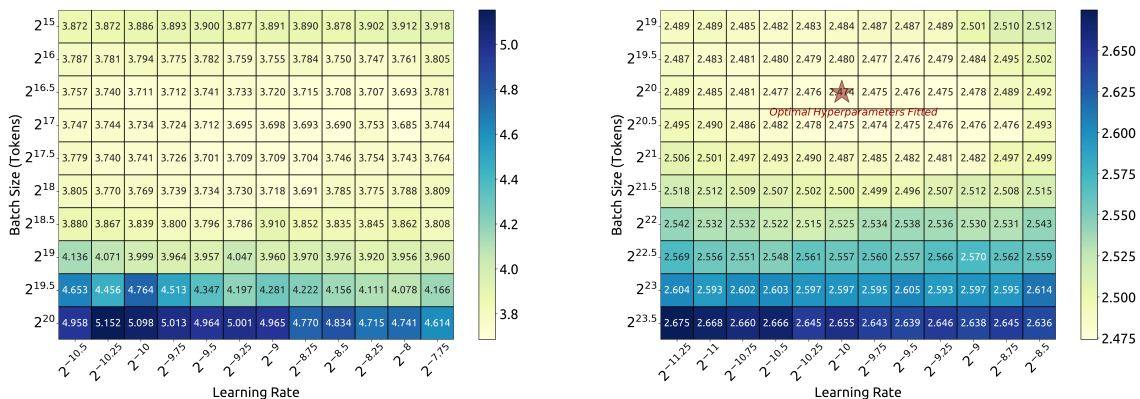


Figure 2 | Training loss w.r.t. batch size and learning rate with  $1e17$  and  $1e20$  FLOPs.

Then, we utilized the aforementioned multi-step learning rate scheduler to effectively train multiple models with different batch sizes, learning rates, and compute budgets ranging from



1e17 to 2e19 by reusing the first stage. Considering the redundancy in the parameter space, we regarded the parameters used by models whose generalization error exceeded the minimum by no more than 0.25% as near-optimal hyperparameters. We then fitted the batch size  $B$  and learning rate  $\eta$  with respect to the compute budget  $C$ . The fitting results, as shown in Figure 3, reveal that the optimal batch size  $B$  gradually increases with the increase in compute budget  $C$ , while the optimal learning rate  $\eta$  gradually decreases. This is in line with the intuitive empirical settings for batch size and learning rate when scaling up models. Moreover, all near-optimal hyperparameters fall within a broad band range, indicating that it is relatively easy to choose near-optimal parameters within this interval. The final formulae we fitted for batch size and learning rate are as follows:

$$\begin{aligned} \eta_{\text{opt}} &= 0.3118 \cdot C^{-0.1250} \\ B_{\text{opt}} &= 0.2920 \cdot C^{0.3271} \end{aligned} \quad (1)$$

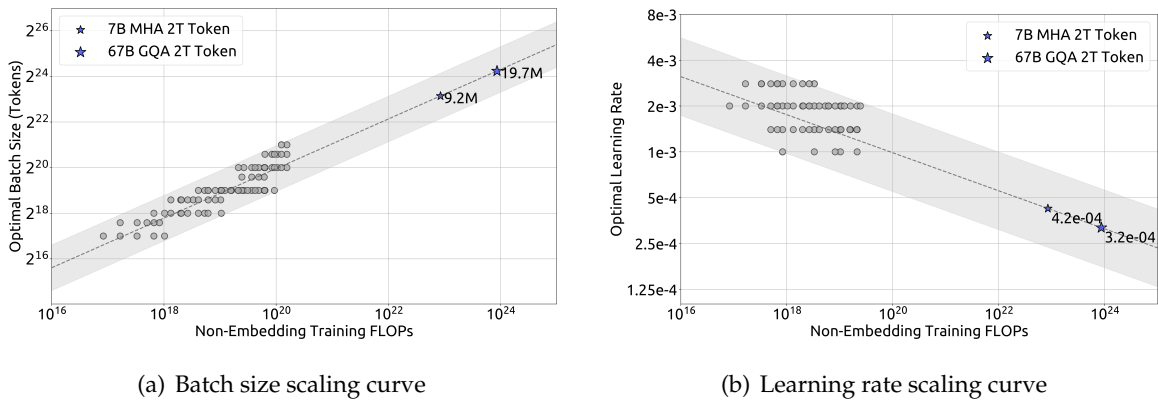


Figure 3 | Scaling curves of batch size and learning rate. The grey circles represent models whose generalization error exceeded the minimum by no more than 0.25%. The dotted line represents the power law fitting the smaller model. The blue stars represent DeepSeek LLM 7B and 67B.

We validated our formulae on a series of models with a 1e20 compute budget, and the results of a specific model size (2.94B FLOPs per token) are shown in Figure 2(b). The results indicate that the fitted parameters are centered in the optimal parameter space. Subsequent sections also show that the parameters we fitted for DeepSeek LLM 7B and 67B models similarly achieved good performance.

However, it’s important to note that we have not yet considered the impact of factors beyond the compute budget  $C$  on the optimal hyperparameters. This is inconsistent with some earlier works (Kaplan et al., 2020; McCandlish et al., 2018) which suggested that the optimal batch size can be modeled as being solely related to the generalization error  $L$ . Furthermore, we observed that in models with the same compute budget but different model/data allocations, the optimal parameter space varies slightly. This suggests that further research is needed to understand the selection of hyperparameters and training dynamics. We will explore these aspects in future works.

### 3.2. Estimating Optimal Model and Data Scaling

After deriving the formulae for fitting near-optimal hyperparameters, we started fitting the scaling curve and analyzing the optimal model/data scaling-up allocation strategy. This strategy involves finding model scaling exponent  $a$  and data scaling exponent  $b$  that satisfy  $N_{\text{opt}} \propto C^a$

and  $D_{\text{opt}} \propto C^b$ , respectively. The data scale  $D$  can be consistently represented by the number of tokens in the dataset. In previous works, the model scale was typically represented by model parameters, with non-embedding parameters  $N_1$  (Kaplan et al., 2020) and complete parameters  $N_2$  (Hoffmann et al., 2022). The relationship between compute budget  $C$  and model/data scale could be approximately described as  $C = 6ND$ , meaning we could use  $6N_1$  or  $6N_2$  to approximate the model scale. However, since both  $6N_1$  and  $6N_2$  do not account for the computational overhead of attention operation, and  $6N_2$  also includes the vocabulary computation, which contributes less to the model’s capacity, they both have significant approximation errors under certain settings.

To mitigate these errors, we introduced a new model scale representation: non-embedding FLOPs/token  $M$ .  $M$  includes the computational overhead of attention operation but does not take into account the vocabulary computation. With the model scale represented by  $M$ , the compute budget  $C$  can be simply expressed as  $C = MD$ . The specific differences between  $6N_1$ ,  $6N_2$ , and  $M$  are as shown in the following formulae:

$$\begin{aligned} 6N_1 &= 72 n_{\text{layer}} d_{\text{model}}^2 \\ 6N_2 &= 72 n_{\text{layer}} d_{\text{model}}^2 + 6 n_{\text{vocab}} d_{\text{model}} \\ M &= 72 n_{\text{layer}} d_{\text{model}}^2 + 12 n_{\text{layer}} d_{\text{model}} l_{\text{seq}} \end{aligned} \quad (2)$$

where  $n_{\text{layer}}$  represents the number of layers,  $d_{\text{model}}$  represents the model width,  $n_{\text{vocab}}$  is the vocabulary size, and  $l_{\text{seq}}$  is the sequence length. We assessed the differences between these three representations across models of varying scales, as shown in Table 3. The results indicate that both  $6N_1$  and  $6N_2$  either overestimate or underestimate the computational cost in models of different scales. This discrepancy is particularly pronounced in small-scale models, with differences reaching up to 50%. Such inaccuracies can introduce substantial statistical errors when fitting the scaling curve. Please refer to Appendix A.2 for further analysis regarding different representations of model scale.

$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{vocab}}$	$l_{\text{seq}}$	$N_1$	$N_2$	$M$	$\frac{6N_1}{M}$	$\frac{6N_2}{M}$
8	512			25.2M	77.6M	352M	0.43	1.32
12	768			84.9M	164M	963M	0.53	1.02
24	1024			302M	407M	3.02B	0.60	0.81
24	2048	102400	4096	1.21B	1.42B	9.66B	0.75	0.88
32	4096			6.44B	6.86B	45.1B	0.85	0.91
40	5120			12.6B	13.1B	85.6B	0.88	0.92
80	8192			64.4B	65.3B	419B	0.92	0.94

Table 3 | Difference in model scale representations and disparities of non-embedding parameters  $N_1$  and complete parameters  $N_2$  relative to non-embedding FLOPs/token  $M$ .

After adopting  $M$  to represent the model scale, our objective could be described more clearly as: Given a computing budget  $C = MD$ , find the optimal model scale  $M_{\text{opt}}$  and data scale  $D_{\text{opt}}$  that minimize the generalization error of the model. This target could be formalized as:

$$M_{\text{opt}}(C), D_{\text{opt}}(C) = \underset{M, D \text{ s.t. } C=MD}{\text{argmin}} L(N, D) \quad (3)$$

To reduce experimental costs and fitting difficulties, the IsoFLOP profile approach from Chinchilla (Hoffmann et al., 2022) was used to fit the scaling curve. We selected 8 different

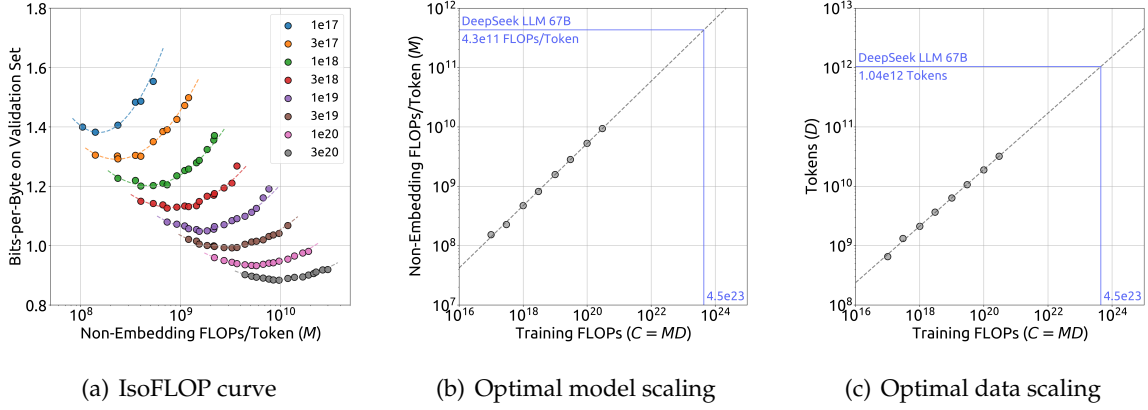


Figure 4 | IsoFLOP curve and optimal model/data allocation. The metric in IsoFLOP curve is bits-per-byte on the validation set. The dotted lines in optimal model/data scaling curves represent the power law fitting the smaller model (grey circles).

compute budgets ranging from  $1e17$  to  $3e20$ , and designed around 10 different model/data scale allocations for each budget. The hyperparameters for each budget were determined by Formula(1), and the generalization error was calculated on an independent validation set, distributed similarly to the training set and containing 100M tokens.

Figure 4 demonstrates the IsoFLOP curve and model/data scaling curves, which are fitted by using the optimal model/data allocation for each compute budget. The specific formulae for the optimal non-embedding FLOPs/token  $M_{opt}$  and optimal tokens  $D_{opt}$  are as follows:

$$\begin{aligned} M_{opt} &= M_{base} \cdot C^a, & M_{base} &= 0.1715, & a &= 0.5243 \\ D_{opt} &= D_{base} \cdot C^b, & D_{base} &= 5.8316, & b &= 0.4757 \end{aligned} \quad (4)$$

Additionally, we fitted the loss scaling curve according to compute budget  $C$  and optimal generalization error, and predicted the generalization error for DeepSeek LLM 7B and 67B, as shown in Figure 5. The results indicate that using small-scale experiments can accurately predict

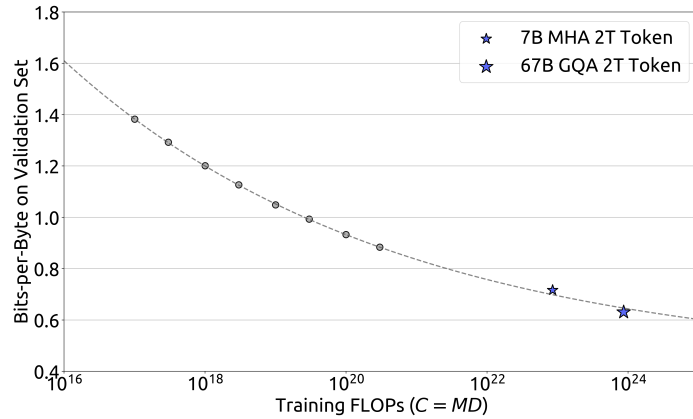


Figure 5 | Performance scaling curve. The metric is the bits-per-byte on the validation set. The dotted line represents the power law fitting the smaller model (grey circles). The blue stars represent DeepSeek LLM 7B and 67B. Their performance is well-predicted by the scaling curve.

the performance of models with 1000× compute budget. This provides both confidence and guidance for training models on a larger scale.

### 3.3. Scaling Laws with Different Data

In the development process of DeepSeek LLM, the dataset was iteratively refined multiple times, with adjustments in the proportions of different data sources while enhancing the overall quality. This allowed us to further analyze the impact of different datasets on scaling laws.

We studied the scaling laws using three different datasets: early in-house data, current in-house data, and OpenWebText2, which was utilized in the previous study of scaling laws (Kaplan et al., 2020). Our internal data assessment revealed that current in-house data has higher data quality than early in-house data. Furthermore, the quality of OpenWebText2 even surpasses the current in-house data, due to its smaller scale which allows for more meticulous processing.

Approach	Coeff. $a$ where $N_{\text{opt}}(M_{\text{opt}}) \propto C^a$	Coeff. $b$ where $D_{\text{opt}} \propto C^b$
OpenAI (OpenWebText2)	0.73	0.27
Chinchilla (MassiveText)	0.49	0.51
Ours (Early Data)	0.450	0.550
Ours (Current Data)	0.524	0.476
Ours (OpenWebText2)	0.578	0.422

Table 4 | Coefficients of model scaling and data scaling vary with training data distribution.

An interesting observation from the analysis is that the optimal model/data scaling-up allocation strategy across these three datasets showed consistency with data quality. As illustrated in Table 4, as data quality improves, the model scaling exponent  $a$  gradually increases, while the data scaling exponent  $b$  decreases, which suggests that the increased compute budget should be allocated more to the model instead of the data. This finding might also explain the significant differences in optimal model/data scaling-up allocation observed in earlier studies of scaling laws.

An intuitive speculation for this finding is that high-quality data usually implies logical clarity and less predictive difficulty after sufficient training. Therefore, it’s more advantageous to scale up the model size when increasing compute budget. We will continue to pay close attention to the changes in data quality and its impact on scaling laws, and provide more analysis in future works.

## 4. Alignment

We collect around 1.5 million instruction data instances in English and Chinese, covering a wide range of helpfulness and harmlessness topics. Our helpful data contains 1.2 million instances, with a distribution of 31.2% for general language tasks, 46.6% for mathematical problems, and 22.2% for coding exercises. The safety data consists of 300K instances, covering various sensitive topics.

Our alignment pipeline contains two stages.

**Supervised Fine-Tuning:** We fine-tuned our 7B model with 4 epochs, but only 2 epochs for the 67B model, since we observed the overfitting problem is serious on the 67B model. We

observed that GSM8K (Cobbe et al., 2021) and HumanEval (Chen et al., 2021) are improved consistently for the 7B model, while the 67B model hits the upper bound soon. The learning rate is  $1e-5$  and  $5e-6$  for 7B and 67B models, respectively. In addition to monitoring the benchmark accuracy, we also assess the repetition ratio of a chat model during the fine-tuning process. We gathered a total of 3868 Chinese and English prompts and determined the proportion of generated responses that fail to terminate and instead endlessly repeat a sequence of text. We observed that the repetition ratio tends to rise as the quantity of math SFT data increases. This can be attributed to the fact that math SFT data occasionally includes similar patterns in reasoning. Consequently, weaker models struggle to grasp such reasoning patterns, resulting in repetitive responses. To tackle the problem, we tried two-stage fine-tuning and DPO (Rafailov et al., 2023), both of which could almost keep the benchmark score and reduce the repetition significantly.

**DPO:** To further enhance the model’s ability, we used the direct preference optimization algorithm (Rafailov et al., 2023), which is proven to be a simple but effective method for LLM alignment. We constructed the preference data for DPO training in terms of helpfulness and harmlessness. For helpfulness data, we collected multilingual prompts, which cover categories including creative writing, question answering, instruction following, and so on. Then we generated responses using our DeepSeek Chat models as response candidates. Similar operations are applied to harmlessness preference data construction.

We trained an epoch for DPO, with a learning rate of  $5e-6$  and batch size of 512, and we used a learning rate warmup and cosine learning rate scheduler. We found out that DPO can strengthen the model’s open-ended generation skill, while engendering little difference in performance among standard benchmarks.

## 5. Evaluation

### 5.1. Public Benchmark Evaluation

We evaluate our models on a series of public benchmarks both in English and Chinese, based on the internal evaluation framework.

**Multi-subject multiple-choice** datasets including MMLU (Hendrycks et al., 2020), C-Eval (Huang et al., 2023) and CMMLU (Li et al., 2023).

**Language understanding and reasoning** datasets including HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2020), ARC (Clark et al., 2018), OpenBookQA (Mihaylov et al., 2018) and BigBench Hard (BBH) (Suzgun et al., 2022).

**Closed-book question answering** datasets including TriviaQA (Joshi et al., 2017) and NaturalQuestions (Kwiatkowski et al., 2019).

**Reading comprehension** datasets including RACE Lai et al. (2017) and DROP (Dua et al., 2019), C3 (Sun et al., 2019).

**Reference disambiguation** datasets including WinoGrande Sakaguchi et al. (2019) and CLUEWSC (Xu et al., 2020).

**Language modeling** datasets including Pile (Gao et al., 2020).

**Chinese understanding and culture** datasets including CHID (Zheng et al., 2019) and CCPM (Li et al., 2021).

**Math** datasets including GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021) and CMATH (Wei et al., 2023).

**Code** datasets including HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021).

**Standardized exams** including AGIEval (Zhong et al., 2023).

We apply perplexity-based evaluation to datasets that require answers to be chosen from several options. These datasets include HellaSwag, PIQA, WinoGrande, RACE-Middle, RACE-High, MMLU, ARC-Easy, ARC-Challenge, OpenBookQA, CHID, C-Eval, CMMLU, C3 and CCPM. The perplexity-based evaluation here refers to calculating the perplexity of each option and selecting the lowest one as the model prediction. For ARC and OpenBookQA, we calculate the perplexity with unconditional normalization (Brown et al., 2020), and for other datasets we use length normalization.

We apply generation-based evaluation for TriviaQA, NaturalQuestions, DROP, MATH, GSM8K, HumanEval, MBPP, BBH, AGIEval, CLUEWSC, and CMATH. The generation-based evaluation here refers to letting the model generate free texts and parsing results from generated texts. For generation-based evaluation, we use greedy decoding.

We apply language-modeling-based evaluation for Pile-test, which means calculating the bits-per-byte on the test corpus.

We use 2048 or 4096 as the maximum sequence length for different benchmarks. Details of evaluation formats can be found in Appendix A.6.

### 5.1.1. Base Model

Table 5 presents the main results on the evaluation benchmark. Despite DeepSeek models are pre-trained on 2T bilingual corpus, they show comparable performance on English language understanding benchmarks with LLaMA2 models, which also consume 2T tokens but focus on English. Furthermore, DeepSeek 67B achieves considerably better performance on MATH, GSM8K, HumanEval, MBPP, BBH, and Chinese benchmarks compared to LLaMA2 70B. We show the benchmark curve in the Appendix A.3. We can see some task performance is boosted as model scaling, such as GSM8K and BBH. Given that we train both 7B and 67B on the same dataset, the emergence of this improvement can be attributed to the powerful few-shot learning ability of large models. However, as the proportion of mathematical data increases, the disparity between small and large models may diminish.

An interesting observation is that the advantage of DeepSeek 67B over LLaMA2 70B is larger than that of DeepSeek 7B over LLaMA2 7B. This phenomenon highlights the greater influence of language conflict on smaller models. Additionally, LLaMA2 demonstrates impressive performance on certain Chinese tasks, such as CMATH, despite not being specifically trained on Chinese data. This suggests that certain fundamental abilities, such as mathematical reasoning, can be effectively transferred across languages. However, tasks like CHID, which involve evaluating the usage of Chinese idioms, require the model to consume a significant number of Chinese tokens during pre-training. In this case, LLaMA2 significantly underperforms compared to DeepSeek LLM.

### 5.1.2. Chat Model

Table 6 demonstrates the results of the DeepSeek Chat models, showcasing overall improvements in most tasks following tuning. However, there were a few instances where the performance of

Language	Benchmark	Test-shots	LLaMA2	DeepSeek	LLaMA2	DeepSeek
			7B	7B	70B	67B
English	HellaSwag	0-shot	75.6	75.4	<b>84.0</b>	<b>84.0</b>
	PIQA	0-shot	78.0	79.2	82.0	<b>83.6</b>
	WinoGrande	0-shot	69.6	70.5	<b>80.4</b>	79.8
	RACE-Middle	5-shot	60.7	63.2	<b>70.1</b>	69.9
	RACE-High	5-shot	45.8	46.5	<b>54.3</b>	50.7
	TriviaQA	5-shot	63.8	59.7	<b>79.5</b>	78.9
	NaturalQuestions	5-shot	25.5	22.2	36.1	<b>36.6</b>
	MMLU	5-shot	45.8	48.2	69.0	<b>71.3</b>
	ARC-Easy	0-shot	69.1	67.9	76.5	<b>76.9</b>
	ARC-Challenge	0-shot	49.0	48.1	<b>59.5</b>	59.0
	OpenBookQA	0-shot	57.4	55.8	<b>60.4</b>	60.2
	DROP	1-shot	39.8	41.0	<b>69.2</b>	67.9
	MATH	4-shot	2.5	6.0	13.5	<b>18.7</b>
	GSM8K	8-shot	15.5	17.4	58.4	<b>63.4</b>
	HumanEval	0-shot	14.6	26.2	28.7	<b>42.7</b>
	MBPP	3-shot	21.8	39.0	45.6	<b>57.4</b>
	BBH	3-shot	38.5	39.5	62.9	<b>68.7</b>
	AGIEval	0-shot	22.8	26.4	37.2	<b>41.3</b>
	Pile-test	-	0.741	0.725	0.649	<b>0.642</b>
Chinese	CLUEWSC	5-shot	64.0	73.1	76.5	<b>81.0</b>
	CHID	0-shot	37.9	89.3	55.5	<b>92.1</b>
	C-Eval	5-shot	33.9	45.0	51.4	<b>66.1</b>
	CMMLU	5-shot	32.6	47.2	53.1	<b>70.8</b>
	CMath	3-shot	25.1	34.5	53.9	<b>63.0</b>
	C3	0-shot	47.4	65.4	61.7	<b>75.3</b>
	CCPM	0-shot	60.7	76.9	66.2	<b>88.5</b>

Table 5 | Main results. The evaluation results we report are based on the internal evaluation framework. **Bold** numbers indicate the best results among the 4 models. For Pile-test we report bits-per-byte (BPB), for DROP we report F1 score and for other tasks we report accuracy. Note that the test-shots is the maximum value and fewer shots might be applied because of limited context length or limited few-shot examples available in the same passage for reading comprehension tasks such as RACE.

certain tasks declined.

**Knowledge:** We have observed fluctuations of base and chat models in knowledge-related tasks, such as TriviaQA, MMLU, and C-Eval. However, we do not believe that such minor fluctuations indicate the acquisition or loss of knowledge after SFT. The value of SFT lies in the ability to learn to achieve comparable scores to the base model’s few-shot setting in the chat model’s zero-shot setting, which is aligned with real scenarios. For example, 0-shot MMLU performance of a chat model is comparable with 5-shot MMLU performance of a base model.

**Reasoning:** As a significant proportion of the SFT instances are in the CoT format Wei et al. (2022), the chat models demonstrate slight improvements in reasoning tasks, such as BBH and NaturalQuestions. However, we believe that the SFT stage does not learn reasoning capabilities but rather the correct format for reasoning paths.

Language	Benchmark	DeepSeek 7B Base	DeepSeek 7B Chat	DeepSeek 67B Base	DeepSeek 67B Chat
English	HellaSwag	75.4	68.5	<b>84.0</b>	75.7
	PIQA	79.2	77.6	<b>83.6</b>	82.6
	WinoGrande	70.5	66.9	<b>79.8</b>	76.0
	RACE-Middle	63.2	65.2	69.9	<b>70.9</b>
	RACE-High	46.5	50.8	50.7	<b>56.0</b>
	TriviaQA	59.7	57.9	78.9	<b>81.5</b>
	NaturalQuestions	22.2	32.5	36.6	<b>47.0</b>
	MMLU	48.2	49.4	<b>71.3</b>	71.1
	ARC-Easy	67.9	71.0	76.9	<b>81.6</b>
	ARC-Challenge	48.1	49.4	59.0	<b>64.1</b>
	GSM8K	17.4	63.0	63.4	<b>84.1</b>
	MATH	6.0	15.8	18.7	<b>32.6</b>
	HumanEval	26.2	48.2	42.7	<b>73.8</b>
	MBPP	39.0	35.2	57.4	<b>61.4</b>
	DROP	41.0	49.1	67.9	<b>71.9</b>
	OpenBookQA	55.8	54.8	60.2	<b>63.2</b>
	BBH	39.5	42.3	68.7	<b>71.7</b>
AGIEval	26.4	19.3	41.3	<b>46.4</b>	
Chinese	CLUEWSC	73.1	71.9	<b>81.0</b>	60.0
	CHID	89.3	64.9	<b>92.1</b>	72.6
	C-Eval	45.0	47.0	<b>66.1</b>	65.2
	CMMLU	47.2	49.7	<b>70.8</b>	67.8
	CMath	34.5	68.4	63.0	<b>80.3</b>
	C3	65.4	66.4	75.3	<b>77.0</b>
	CCPM	76.9	76.5	<b>88.5</b>	84.9

Table 6 | The comparison between base and chat models. We evaluate chat models with 0-shot for MMLU, GSM8K, MATH, C-Eval, and CMMLU, while base model results are still obtained in the few-shot setting.

**Performance Drop Tasks:** The performance of a few tasks consistently declines after fine-tuning, regardless of the model size or pre-trained checkpoint selected. These particular tasks typically involve cloze tasks or sentence completion tasks, such as HellaSwag. It is reasonable to assume that pure language models are better equipped to handle such tasks.

**Math and Code:** Our model exhibits significant improvements in math and coding tasks after fine-tuning. For instance, HumanEval and GSM8K scores are improved by over 20 points. Our explanation for this is that the base model was initially underfitted for these tasks, and the SFT stage has learned additional knowledge in coding and mathematics through the extensive SFT data. However, it is important to note that the model’s capabilities may be primarily focused on code completion and algebraic questions. To develop a comprehensive understanding of mathematics and coding, it is crucial to incorporate a diverse range of data during the pre-training stage, which is left as future work. We conducted a detailed analysis of code and math tasks in Appendix A.4.

In the 7B model fine-tuning, we initially fine-tune the model using all data. Subsequently, a second stage is introduced, which excludes math and code data. The motivation behind this approach is that the stage-1 model exhibits a repetition ratio of 2.0%, which is reduced to 1.4%



Model 模型	Overall 总分	Reasoning 中文推理			Language 中文语言						
		Avg. 推理 总分	Math. 数学 计算	Logi. 逻辑 推理	Avg. 语言 总分	Fund. 基本 任务	Chi. 中文 理解	Open. 综合 问答	Writ. 文本 写作	Role. 角色 扮演	Pro. 专业 能力
gpt-4-1106-preview	<b>8.01</b>	<b>7.73</b>	<b>7.80</b>	<b>7.66</b>	<b>8.29</b>	<b>7.99</b>	7.33	<b>8.61</b>	<b>8.67</b>	<b>8.47</b>	<b>8.65</b>
gpt-4-0613	<b>7.53</b>	7.47	7.56	7.37	7.59	7.81	6.93	7.42	7.93	7.51	7.94
<b>DeepSeek-67B-Chat-DPO*</b>	<b>6.69</b>	5.77	6.13	5.41	7.60	7.29	7.47	7.82	7.51	7.83	7.71
<b>DeepSeek-67B-Chat*</b>	<b>6.43</b>	5.75	5.71	5.79	7.11	7.12	6.52	7.58	7.20	6.91	7.37
chatglm-turbo (智谱清言)	<b>6.24</b>	5.00	4.74	5.26	7.49	6.82	7.17	8.16	7.77	7.76	7.24
erniebot-3.5 (文心一言)	<b>6.14</b>	5.15	5.03	5.27	7.13	6.62	<b>7.60</b>	7.26	7.56	6.83	6.90
gpt-3.5-turbo-0613	<b>6.08</b>	5.35	5.68	5.02	6.82	6.71	5.81	7.29	7.03	7.28	6.77
chatglm-pro (智谱清言)	<b>5.83</b>	4.65	4.54	4.75	7.01	6.51	6.76	7.47	7.07	7.34	6.89
spark_desk_v2 (讯飞星火)	<b>5.74</b>	4.73	4.71	4.74	6.76	5.84	6.97	7.29	7.18	6.92	6.34
Qwen-14B-Chat	<b>5.72</b>	4.81	4.91	4.71	6.63	6.90	6.36	6.74	6.64	6.59	6.56
Baichuan2-13B-Chat	<b>5.25</b>	3.92	3.76	4.07	6.59	6.22	6.05	7.11	6.97	6.75	6.43
ChatGLM3-6B	<b>4.97</b>	3.85	3.55	4.14	6.10	5.75	5.29	6.71	6.83	6.28	5.73
Baichuan2-7B-Chat	<b>4.97</b>	3.66	3.56	3.75	6.28	5.81	5.50	7.13	6.84	6.53	5.84
InternLM-20B	<b>4.96</b>	3.66	3.39	3.92	6.26	5.96	5.50	7.18	6.19	6.49	6.22
Qwen-7B-Chat	<b>4.91</b>	3.73	3.62	3.83	6.09	6.40	5.74	6.26	6.31	6.19	5.66
ChatGLM2-6B	<b>4.48</b>	3.39	3.16	3.61	5.58	4.91	4.52	6.66	6.25	6.08	5.08
InternLM-Chat-7B	<b>3.65</b>	2.56	2.45	2.66	4.75	4.34	4.09	5.82	4.89	5.32	4.06
Chinese-LLaMA-2-7B-Chat	<b>3.57</b>	2.68	2.29	3.07	4.46	4.31	4.26	4.50	4.63	4.91	4.13
LLaMA-2-13B-Chinese-Chat	<b>3.35</b>	2.47	2.21	2.73	4.23	4.13	3.31	4.79	3.93	4.53	4.71

Table 7 | AlignBench leaderboard rated by gpt-4-0613. Models are ranked in descending order of total score. Results with \* are our evaluation results based on the official AlignBench repository, whereas all other results are derived from the AlignBench paper. We found that our Deepseek-67B-Chat model surpasses ChatGPT and other baseline models by a clear margin, which indicates the superior performance of our model in both basic Chinese language tasks and advanced Chinese reasoning tasks. Besides, we can find that the DPO process has brought improvements in almost all fields.

after stage-2 tuning, while maintaining the benchmark score. In the case of the 67B model, the repetition ratio is already below 1% following the first stage fine-tuning, and the second stage hurts the model score on the benchmark. Therefore, only one stage of SFT is done for the 67B model.

## 5.2. Open-Ended Evaluation

For chat models, in addition to observing metrics on standard benchmarks, the quality of results generated in open domains and open-ended questions directly affects the actual user experience. Hence, we separately tested the open-ended generation capabilities of our chat model in both Chinese and English tasks.

### 5.2.1. Chinese Open-Ended Evaluation

For Chinese open-ended evaluation, we tested the comprehensive of our chat model in different domains on a high-quality open-ended question testset AlignBench (Liu et al., 2023). AlignBench includes a total of 8 primary categories, 36 secondary categories, and encompasses 683 questions. For each question, in addition to the prompt, AlignBench also provides professional reference answers and rating templates for GPT-4 to judge the quality of the response.

We utilized the official AlignBench Github code repository to implement the evaluation of

our model. We strictly aligned the key temperature parameter with the original setting: for role-playing, writing ability, and open-ended questions, the generation temperature was set to 0.7; whereas for other tasks, the generation temperature was set to 0.1.

The AlignBench leaderboard is shown in Table 7. We can find that our DeepSeek 67B Chat model surpasses ChatGPT and other baseline models, and is only after the two versions of GPT-4. This demonstrates the excellent performance of our model across various Chinese tasks, compared to other open-source or proprietary Chinese Large Language Models. The DPO model has shown improvement across almost all metrics, which demonstrates the positive impact of the DPO training process on model alignment.

For the basic Chinese Language tasks, our model is in the first tier among all models, and the Chinese fundamental language ability of our DPO model is even higher than the newest version of GPT-4. For the advanced Chinese Reasoning tasks, our model’s scores are significantly higher than those of other Chinese LLMs with a clear margin, demonstrating the superior performance of our model in more complex Chinese logical reasoning and mathematical calculations.

### 5.2.2. English Open-Ended Evaluation

For English open-ended evaluation, we use the MT-Bench benchmark (Zheng et al., 2023), which contains 8 different categories of multi-turn questions. As illustrated in Table 8, our DeepSeek LLM 67B Chat outperforms other open-source models such as LLaMA-2-Chat Touvron et al. (2023b) 70B, Xwin 70b v0.1, and TULU 2+DPO 70B (Iverson et al., 2023), and achieves 8.35 score comparable with GPT-3.5-turbo. Besides, after the DPO stage, our DeepSeek LLM 67B Chat DPO further improves the average score to 8.76, which is only behind GPT-4 (OpenAI, 2023). These results illustrate the strong multi-turn open-ended generation ability of DeepSeek LLM.

Model	STEM	Humanities	Reasoning	Coding	Math	Extraction	Roleplay	Writing	Average
GPT-4-1106-preview*	9.90	9.95	8.10	9.05	7.95	9.90	9.50	9.70	9.26
GPT-3.5-turbo-0613*	9.55	9.95	6.20	7.05	7.05	9.00	8.65	9.65	8.39
LLAMA-2-Chat 7B*	8.65	8.75	4.25	3.00	2.40	6.50	7.70	8.90	6.27
LLAMA-2-Chat 13B*	8.63	9.75	5.10	3.00	3.45	6.93	7.50	8.85	6.65
LLAMA-2-Chat 70B*	8.93	9.63	5.80	3.15	3.30	7.25	7.50	9.30	6.86
Zephyr-Beta 7B*	9.03	9.63	5.60	5.10	4.45	7.45	8.20	9.35	7.35
Xwin 70b v0.1*	9.68	9.95	6.55	4.25	3.30	8.75	8.25	9.55	7.53
Xwin 13b v0.2*	9.55	9.88	5.20	3.60	2.85	7.70	8.60	8.68	7.01
TULU 2+DPO 70B*	9.00	9.90	7.00	4.70	4.65	9.35	9.25	9.25	7.89
<b>DeepSeek LLM 67B Chat</b>	9.60	9.70	8.00	7.35	6.25	8.40	8.20	9.30	8.35
<b>DeepSeek LLM 67B Chat DPO</b>	9.70	9.80	9.05	6.75	6.65	9.30	9.10	9.75	<b>8.76</b>

Table 8 | MT-Bench Evaluation. Results with \* are reported in Iverson et al. (2023)

### 5.3. Held-Out Evaluation

Data contamination and benchmark overfitting are two challenges in evaluating LLMs. One common practice is to utilize testsets published recently to evaluate the model as held-out testsets.

**LeetCode:** To assess the coding proficiency of the model, we have utilized problems from the LeetCode Weekly Contest (Weekly Contest 351-372, Bi-Weekly Contest 108-117, from July 2023 to Nov 2023). We have obtained these problems by crawling data from LeetCode, which consists of 126 problems with over 20 test cases for each. The evaluation metric employed is akin to that of HumanEval. In this regard, if a model’s outputs successfully pass all test cases, the model is considered to have effectively solved the problem. The model’s coding capabilities are

depicted in the Figure below, where the y-axis represents the pass@1 score on in-domain human evaluation testing, and the x-axis represents the pass@1 score on out-domain LeetCode Weekly Contest problems. The LeetCode test data will be released accompanied with the DeepSeek Coder technique report soon.

**Hungarian National High-School Exam:** In line with Grok-1, we have evaluated the model’s mathematical capabilities using the Hungarian National High School Exam. This exam comprises 33 problems, and the model’s scores are determined through human annotation. We follow the scoring metric in the solution.pdf to evaluate all models.

**Instruction Following Evaluation:** On Nov 15th, 2023, Google released an instruction following the evaluation dataset (Zhou et al., 2023). They identified 25 types of verifiable instructions and constructed around 500 prompts, with each prompt containing one or more verifiable instructions. We use the prompt-level loose metric to evaluate all models.

Model	LeetCode	Hungarian Exam	IFEval
GPT-4	48.4	68	79.3
ChatGLM3 6B	2.4	32	29.7
DeepSeek LLM 7B Chat	4.7	28.5	41.2
Baichuan2-Chat 13B	1.6	19.5	44.5
Yi-Chat 34B	7.9	39	48.4
Qwen 72B Chat	12.7	52	50.8
DeepSeek LLM 67B Chat	<b>17.5</b>	<b>58</b>	<b>55.5</b>

Table 9 | Held-out Dataset Evaluation.

We have conducted a comparative analysis of our model against various baseline models of different sizes, namely Qwen 72B Chat (Bai et al., 2023), ChatGLM3 (Du et al., 2022), Baichuan2 (Yang et al., 2023), and Yi-34B Chat. Our observations indicate that there exists a significant performance gap between large models and small models on these held-out datasets, even if certain small models achieve promising results on conventional benchmarks. For instance, ChatGLM3 achieves a score of 52.4 on MBPP, a code testset, which is close to DeepSeek 67B. However, when evaluated on new benchmarks, its performance falls considerably short compared to DeepSeek 67B. A similar trend is also observed in math datasets, where ChatGLM3 is very strong on GSM8K (72.3), but its performance in the Hungarian Exam score is inferior to large models. Furthermore, the capability of instruction following demonstrates that total computing plays a crucial role.

The DeepSeek 7B and 67B models utilize the same training pipeline, but there is a significant disparity in their performance. Through our subjective evaluation, we have observed a notable discrepancy in intelligence across various tasks when scaling model size to 67B. While DeepSeek 7B falls behind other smaller language models on standard benchmarks, its performance on held-out tasks is relatively commendable when compared to others.

#### 5.4. Safety Evaluation

We profoundly recognize the importance of safety for general artificial intelligence. The premise for establishing a truly helpful artificial intelligence model is that it possesses values consistent with those of humans and exhibits friendliness towards humanity. We incorporate the assurance of model safety throughout the entire training process, including pre-training, SFT, and DPO.

To validate the safety of our model, we established a 20-person expert team from various

Category	Subcategory	#Safety Answers / #Total Cases
歧视偏见问题 (Discrimination and Prejudice Questions)	民族种族 (Ethnic and Racial), 宗教信仰 (Religious Belief), 国别地域 (Nationality and Geography), 性别 (Gender), 年龄 (Age), 职业 (Occupation), 健康 (Health), 其他方面歧视 (Discrimination in Other Aspects)	486/500
侵犯他人合法权益 (Infringement of Others' Legal Rights)	身心健康 (Physical and Mental Health), 合法财产 (Legitimate Property), 肖像权 (Portrait Rights), 名誉权 (Reputation Rights), 荣誉权 (Honor Rights), 隐私权 (Privacy Rights), 信息权益 (Information Rights), 其他合法权益 (Other Legal Rights)	473/500
商业秘密与知识产权 (Trade Secrets and Intellectual Property Rights)	侵犯他人知识产权 (Infringing Others' Intellectual Property Rights), 垄断和不正当竞争行为 (Monopolistic and Unfair Competitive Actions), 其他商业违法违规行为 (Other Commercially Illegal and Non-compliant Behaviors), 违反商业道德 (Violating Business Ethics), 泄露他人商业机密 (Disclosing Others' Trade Secrets)	281/300
违法违规行为 (Illegal and Non-compliant Behavior)	邪教迷信 (Cults and Superstition), 色情 (Pornography), 赌博 (Gambling), 毒品和违禁品 (Drugs and Prohibited Items), 侮辱谩骂 (Insults and Abuse), 暴力行为 (Violent Behavior), 涉黑涉恶 (Involvement in Organized Crime), 其他违法违规行为 (Other Illegal and Non-compliant Behaviors)	290/300
其他安全问题 (Other Safety Issues)	幻觉和真实性问题 (Issues of Illusion and Reality), 时效性问题 (Time-sensitive Issues), 自我认知问题 (Self-recognition Problems), 其他敏感话题 (Other Sensitive Topics)	767/800

Table 10 | Our taxonomy for safety evaluation. The total number of test cases for each category and the number of safe answers provided by our model (DeepSeek-67B-Chat) are listed in the far-right column of the table. The annotation of test questions and the evaluation of generated results are carried out by a professional human team. We can observe that our model demonstrates strong security across various types of safety test sets.

disciplines and constructed a safety content classification system that aligns with human values (the safety evaluation taxonomy is shown in Table 10). Subsequently, the expert team constructed dozens of high-quality test cases for each safety subcategory manually. In addition to focusing on the diversity of safety content areas, we also pay attention to the diversity of formats in safety content. The infamous "grandmother" loophole indicates that models can be deceived by the surface format of a query into providing unsafe responses. Therefore, when devising questions, the expert team also pays attention to diversifying the ways of inquiry. They construct diverse safety issues through means such as inducement, role-playing, multi-turn dialogues, preset positions, and etc.. Ultimately, we obtained a safety test set comprising 2400 questions. In addition, the expert team has constructed a basic guideline constitution for safety reviews for each different content type and format type.

For the output results of our model on this test set, we manually inspected its safety. Our review team was well-trained and cross-verification was performed on the annotation results. The annotators perform a three-category annotation for each question: safe, unsafe, and model refusal. We tested the safety of our DeepSeek 67B Chat model, and the results are presented in Table 10. The number of test questions for each safety category and the number of safety tests passed by our model are listed in the table. We label both the securely answered and the model-refused test cases as secure responses. The results indicate that our model exhibits good security performance across numerous safety test categories.

Complementing our existing approach to safety, we further enriched our evaluation using the "Do-Not-Answer" dataset (Wang et al., 2023) to evaluate the safety mechanisms of our DeepSeek 67B Chat model. The dataset's 939 risk-categorized prompts were instrumental in highlighting our model's enhanced capabilities. As shown in Table 11, DeepSeek 67B Chat model has demonstrated notable performance, achieving a score of 97.8, which is higher than both ChatGPT and GPT-4. This score not only benchmarks our model's capability to safely handle sensitive queries but also places it competitively among leading models in the field.

## 5.5. Discussion

Throughout the development process, we have discovered some interesting findings in building LLMs.

Model	Do-Not-Answer
LLAMA-2-7B-Chat	99.4
Claude	98.3
DeepSeek-67B-Chat*	97.8
ChatGPT	97.7
GPT-4	96.5
Vicuna-7B	94.9
ChatGLM2	92.9

Table 11 | Do-Not-Answer Score (Wang et al., 2023), a higher score signifies greater model safety. Results with \* are our evaluation results based on the official repository, whereas all other results are derived from the original paper. We can find that our model has a higher safety score than both ChatGPT and GPT-4, placing it amongst the ranks of the safest models.

**Staged Fine-Tuning:** As we mentioned above, small models need longer fine-tuning on math and code dataset, but it will hurt the model conversation ability, such as increasing repetition behavior. To address this issue, we have implemented a staged fine-tuning process. In this approach, the first stage involves fine-tuning with all available data, while the second stage focuses specifically on fine-tuning with conversational data.

Model	HumanEval	GSM8K	Repetition	IFEval
DeepSeek LLM 7B Chat Stage1	48.2	63.9	0.020	38.0
DeepSeek LLM 7B Chat Stage2	48.2	63.0	0.014	41.2

Table 12 | Two-stage fine-tuning results. The repetition ratio is computed when the temperature is 0. The lower repetition ratio is better. The IFEval result is the prompt-level loose accuracy.

Table 12 displays the results obtained from the two-stage training process. These results clearly demonstrate that the second stage does not compromise the model’s proficiency in code and math, while simultaneously decreasing the repetition behavior and enhancing instruction following capability.

**Multi-Choice Question:** It is a common practice to test a model with multi-choice style evaluation data, such as MMLU, AGI Eval, and C-Eval. Multi-choice questions require the model not only to have the corresponding knowledge but also to understand what the option refers to. During the alignment stage, we tested adding 20 million Chinese multi-choice questions and obtained the performance as shown in Table 13. It is important to note that we conducted deduplication for the C-Eval validation set and CMMLU test set to prevent data contamination.

Model	MMLU	C-Eval	CMMLU	TriviaQA	ChineseQA
DeepSeek LLM 7B Chat	49.4	47.0	49.7	57.9	75.0
DeepSeek LLM 7B Chat + MC	60.9	71.3	73.8	57.9	74.4

Table 13 | The impact of adding multi-choice question data.

The inclusion of an additional 20M MC (multiple-choice) data has proven to be beneficial not only for Chinese multiple-choice benchmarks but also for improving English benchmarks. This indicates that the model’s capability to solve MC problems has been enhanced. However, we have observed that this improvement does not extend to the model’s performance on other evaluations that do not utilize the multiple-choice format, such as TriviaQA and our in-house

ChineseQA testsets, which are generative evaluation benchmarks. This suggests that users may not perceive the model as becoming more intelligent during conversational interactions, as these interactions involve generating responses rather than solving multiple-choice problems.

Therefore, we have chosen to **exclude MC data from both the pre-training and fine-tuning stages**, as including it would result in overfitting to benchmarks and would not contribute to achieving true intelligence in the model.

**Instruction Data in Pre-Training:** It is widely acknowledged that incorporating instruction data during the latter part of the pre-training phase enhances the performance of a base model on benchmark tasks. In our study, we integrated 5 million instruction data, primarily consisting of multi-choice questions, during the final 10% of the pre-training stage. We observed that the base model did exhibit improved performance on the benchmark. However, the final outcomes were nearly identical to those achieved by adding the same data during the SFT stage. We conclude that while this approach strengthens the base model’s performance on the benchmark, its overall potential is equivalent to not incorporating these instruction data. If the instruction data is substantial in size, it is acceptable to incorporate it into the pre-training process. Due to our preference for excluding multi-choice questions and the limited availability of non-multi-choice questions we have, we made the decision not to include instruction data in the pre-training process.

**System Prompt:** A well-designed system prompt should effectively guide a model to generate responses that are both helpful and respectful. We slightly changed the prompt introduced by LLaMA-2 as our system prompt.

*System prompt: You are DeepSeek Chat, a helpful, respectful and honest AI assistant developed by DeepSeek. The knowledge cut-off date for your training data is up to May 2023. Always answer as helpfully as possible, while being safe. Your answers should not include any harmful, unethical, racist, sexist, toxic, dangerous, or illegal content. Please ensure that your responses are socially unbiased and positive in nature. If a question does not make any sense, or is not factually coherent, explain why instead of answering something not correct. If you don’t know the answer to a question, please don’t share false information.*

We have observed an intriguing phenomenon wherein the performance of a 7B LLM experiences a slight degradation when a system prompt is introduced. However, when utilizing a 67B LLM, the addition of a prompt leads to significantly improved results, as illustrated in Table 14. Our explanation for this disparity is that larger models possess a better understanding of the intended meaning behind the system prompt, enabling them to follow instructions more effectively and generate superior responses. On the other hand, smaller models struggle to grasp the system prompt adequately, and the inconsistency between training and testing might negatively impact their performance.

Model	MT Bench
DeepSeek LLM 7B Chat	7.15
DeepSeek LLM 7B Chat + System Prompt	7.11
DeepSeek LLM 67B Chat	8.35
DeepSeek LLM 67B Chat + System Prompt	8.58

Table 14 | The impact of adding a system prompt.

## 6. Conclusion, Limitation, and Future Work

We introduce DeepSeek LLMs, a series of open-source models trained from scratch on a vast dataset of 2 trillion tokens in both English and Chinese. In this paper, we provide an in-depth explanation of hyper-parameters selection, scaling laws, as well as the various fine-tuning attempts we made. We calibrate the scaling laws in the previous work and propose a new optimal model/data scaling-up allocation strategy. In addition, we present a method to predict the near-optimal batch size and learning rate with given compute budget. We further conclude that the scaling laws is related to the data quality, which might be the root cause of varying scaling behavior in different works. Guided by the scaling laws, we conduct pre-training with the best hyper-parameter and provide a comprehensive evaluation. We avoid benchmark decoration and dark secrets in all training stages.

DeepSeek Chat shares the acknowledged limitations commonly found in other LLMs, which include the lack of ongoing knowledge updates after pre-training, the possibility of generating non-factual information such as unverified advice, and a tendency to produce hallucinations. Moreover, it is important to note that our initial version of Chinese data is not exhaustive, which may result in suboptimal performance on certain Chinese-specific topics. Since our data primarily consists of Chinese and English sources, the model’s proficiency in other languages remains delicate and should be approached with caution.

DeepSeek LLM is a long-term project committed to advancing open-source language models.

- Soon, we will release our technique reports in code intelligence and Mixture-of-Experts (MoE), respectively. They show how we create high-quality code data for pre-training, and design a sparse model to achieve dense model performance.
- At present, we are constructing a larger and improved dataset for the upcoming version of DeepSeek LLM. We hope the reasoning, Chinese knowledge, math, and code capabilities will be significantly improved in the next version.
- Our alignment team is dedicated to studying ways to deliver a model that is helpful, honest, and safe to the public. Our initial experiments prove that reinforcement learning could boost model complex reasoning capability.

## References

- J. Ainslie, J. Lee-Thorp, M. de Jong, Y. Zemlyanskiy, F. Lebrón, and S. Sanghai. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. [arXiv preprint arXiv:2305.13245](#), 2023.
- Anthropic. Introducing Claude, 2023. URL <https://www.anthropic.com/index/introducing-claude>.
- J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le, et al. Program synthesis with large language models. [arXiv preprint arXiv:2108.07732](#), 2021.
- J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang, et al. Qwen technical report. [arXiv preprint arXiv:2309.16609](#), 2023.
- Y. Bisk, R. Zellers, R. L. Bras, J. Gao, and Y. Choi. PIQA: reasoning about physical commonsense in natural language. In [The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI](#)

- 2020, *The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020*, New York, NY, USA, February 7-12, 2020, pages 7432–7439. AAAI Press, 2020. doi: 10.1609/aaai.v34i05.6239. URL <https://doi.org/10.1609/aaai.v34i05.6239>.
- T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners, 2020.
- M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin, B. Chan, S. Gray, N. Ryder, M. Pavlov, A. Power, L. Kaiser, M. Bavarian, C. Winter, P. Tillet, F. P. Such, D. Cummings, M. Plappert, F. Chantzis, E. Barnes, A. Herbert-Voss, W. H. Guss, A. Nichol, A. Paino, N. Tezak, J. Tang, I. Babuschkin, S. Balaji, S. Jain, W. Saunders, C. Hesse, A. N. Carr, J. Leike, J. Achiam, V. Misra, E. Morikawa, A. Radford, M. Knight, M. Brundage, M. Murati, K. Mayer, P. Welinder, B. McGrew, D. Amodei, S. McCandlish, I. Sutskever, and W. Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374, 2021. URL <https://arxiv.org/abs/2107.03374>.
- P. Clark, I. Cowhey, O. Etzioni, T. Khot, A. Sabharwal, C. Schoenick, and O. Tafjord. Think you have solved question answering? try arc, the AI2 reasoning challenge. *CoRR*, abs/1803.05457, 2018. URL <http://arxiv.org/abs/1803.05457>.
- K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- T. Computer. Redpajama: an open dataset for training large language models, 2023. URL <https://github.com/togethercomputer/RedPajama-Data>.
- Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- T. Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. 2023.
- T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems*, 2022.
- Z. Du, Y. Qian, X. Liu, M. Ding, J. Qiu, Z. Yang, and J. Tang. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335, 2022.
- D. Dua, Y. Wang, P. Dasigi, G. Stanovsky, S. Singh, and M. Gardner. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In J. Burstein, C. Doran, and T. Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 2368–2378. Association for Computational Linguistics, 2019. doi: 10.18653/V1/N19-1246. URL <https://doi.org/10.18653/v1/n19-1246>.



- L. Gao, S. Biderman, S. Black, L. Golding, T. Hoppe, C. Foster, J. Phang, H. He, A. Thite, N. Nabeshima, et al. The Pile: An 800GB dataset of diverse text for language modeling. arXiv preprint arXiv:2101.00027, 2020.
- Google. An important next step on our AI journey, 2023. URL <https://blog.google/technology/ai/bard-google-ai-search-updates/>.
- Z. Gou, Z. Shao, Y. Gong, Y. Shen, Y. Yang, M. Huang, N. Duan, and W. Chen. Tora: A tool-integrated reasoning agent for mathematical problem solving. *CoRR*, abs/2309.17452, 2023. doi: 10.48550/ARXIV.2309.17452. URL <https://doi.org/10.48550/arXiv.2309.17452>.
- P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He. Accurate, large minibatch sgd: Training imagenet in 1 hour. arXiv preprint arXiv:1706.02677, 2017.
- D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300, 2020.
- D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt. Measuring mathematical problem solving with the math dataset. arXiv preprint arXiv:2103.03874, 2021.
- T. Henighan, J. Kaplan, M. Katz, M. Chen, C. Hesse, J. Jackson, H. Jun, T. B. Brown, P. Dhariwal, S. Gray, et al. Scaling laws for autoregressive generative modeling. arXiv preprint arXiv:2010.14701, 2020.
- J. Hestness, S. Narang, N. Ardalani, G. Diamos, H. Jun, H. Kianinejad, M. M. A. Patwary, Y. Yang, and Y. Zhou. Deep learning scaling is predictable, empirically. arXiv preprint arXiv:1712.00409, 2017.
- High-flyer. Hai-llm: 高效且轻量的大模型训练工具, 2023. URL <https://www.high-flyer.cn/en/blog/hai-llm>.
- J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas, L. A. Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. van den Driessche, B. Damoc, A. Guy, S. Osindero, K. Simonyan, E. Elsen, J. W. Rae, O. Vinyals, and L. Sifre. Training compute-optimal large language models. *CoRR*, abs/2203.15556, 2022. doi: 10.48550/ARXIV.2203.15556. URL <https://doi.org/10.48550/arXiv.2203.15556>.
- Y. Huang, Y. Bai, Z. Zhu, J. Zhang, J. Zhang, T. Su, J. Liu, C. Lv, Y. Zhang, J. Lei, et al. C-Eval: A multi-level multi-discipline chinese evaluation suite for foundation models. arXiv preprint arXiv:2305.08322, 2023.
- Huggingface Team. Tokenizers: Fast state-of-the-art tokenizers optimized for research and production, 2019. URL <https://github.com/huggingface/tokenizers>.
- F. i, M. Suzgun, M. Freitag, X. Wang, S. Srivats, S. Vosoughi, H. W. Chung, Y. Tay, S. Ruder, D. Zhou, D. Das, and J. Wei. Language models are multilingual chain-of-thought reasoners. In The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023. OpenReview.net, 2023. URL <https://openreview.net/pdf?id=fR3wGck-IXp>.

- H. Ivison, Y. Wang, V. Pyatkin, N. Lambert, M. Peters, P. Dasigi, J. Jang, D. Wadden, N. A. Smith, I. Beltagy, and H. Hajishirzi. Camels in a changing climate: Enhancing lm adaptation with tulu 2. 2023.
- A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. d. I. Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, et al. Mistral 7b. arXiv preprint arXiv:2310.06825, 2023.
- M. Joshi, E. Choi, D. Weld, and L. Zettlemoyer. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In R. Barzilay and M.-Y. Kan, editors, Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1601–1611, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1147. URL <https://aclanthology.org/P17-1147>.
- J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei. Scaling laws for neural language models. CoRR, abs/2001.08361, 2020. URL <https://arxiv.org/abs/2001.08361>.
- V. A. Korthikanti, J. Casper, S. Lym, L. McAfee, M. Andersch, M. Shoeybi, and B. Catanzaro. Reducing activation recomputation in large transformer models. Proceedings of Machine Learning and Systems, 5, 2023.
- T. Kwiatkowski, J. Palomaki, O. Redfield, M. Collins, A. P. Parikh, C. Alberti, D. Epstein, I. Polosukhin, J. Devlin, K. Lee, K. Toutanova, L. Jones, M. Kelcey, M. Chang, A. M. Dai, J. Uszkoreit, Q. Le, and S. Petrov. Natural questions: a benchmark for question answering research. Trans. Assoc. Comput. Linguistics, 7:452–466, 2019. doi: 10.1162/tacl\_a\_00276. URL [https://doi.org/10.1162/tacl\\_a\\_00276](https://doi.org/10.1162/tacl_a_00276).
- W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. E. Gonzalez, H. Zhang, and I. Stoica. Efficient memory management for large language model serving with pagedattention. In Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles, 2023.
- G. Lai, Q. Xie, H. Liu, Y. Yang, and E. H. Hovy. RACE: large-scale reading comprehension dataset from examinations. In M. Palmer, R. Hwa, and S. Riedel, editors, Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017, pages 785–794. Association for Computational Linguistics, 2017. doi: 10.18653/V1/D17-1082. URL <https://doi.org/10.18653/v1/d17-1082>.
- H. Li, Y. Zhang, F. Koto, Y. Yang, H. Zhao, Y. Gong, N. Duan, and T. Baldwin. CMMLU: Measuring massive multitask language understanding in Chinese. arXiv preprint arXiv:2306.09212, 2023.
- W. Li, F. Qi, M. Sun, X. Yi, and J. Zhang. Ccpm: A chinese classical poetry matching dataset, 2021.
- X. Liu, X. Lei, S. Wang, Y. Huang, Z. Feng, B. Wen, J. Cheng, P. Ke, Y. Xu, W. L. Tam, X. Zhang, L. Sun, H. Wang, J. Zhang, M. Huang, Y. Dong, and J. Tang. Alignbench: Benchmarking chinese alignment of large language models. CoRR, abs/2311.18743, 2023. doi: 10.48550/ARXIV.2311.18743. URL <https://doi.org/10.48550/arXiv.2311.18743>.
- I. Loshchilov and F. Hutter. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.

- H. Luo, Q. Sun, C. Xu, P. Zhao, J. Lou, C. Tao, X. Geng, Q. Lin, S. Chen, and D. Zhang. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. [arXiv preprint arXiv:2308.09583](#), 2023.
- S. McCandlish, J. Kaplan, D. Amodei, and O. D. Team. An empirical model of large-batch training. [arXiv preprint arXiv:1812.06162](#), 2018.
- T. Mihaylov, P. Clark, T. Khot, and A. Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering, 2018.
- D. Narayanan, M. Shoeybi, J. Casper, P. LeGresley, M. Patwary, V. Korthikanti, D. Vainbrand, P. Kashinkunti, J. Bernauer, B. Catanzaro, et al. Efficient large-scale language model training on gpu clusters using megatron-lm. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–15, 2021.
- OpenAI. Introducing ChatGPT, 2022. URL <https://openai.com/blog/chatgpt>.
- OpenAI. GPT4 technical report. [arXiv preprint arXiv:2303.08774](#), 2023.
- L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- G. Penedo, Q. Malartic, D. Hesslow, R. Cojocaru, A. Cappelli, H. Alobeidli, B. Pannier, E. Almazrouei, and J. Launay. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. [arXiv preprint arXiv:2306.01116](#), 2023.
- A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, and C. Finn. Direct preference optimization: Your language model is secretly a reward model. 2023.
- S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–16. IEEE, 2020.
- K. Sakaguchi, R. L. Bras, C. Bhagavatula, and Y. Choi. Winogrande: An adversarial winograd schema challenge at scale, 2019.
- C. J. Shallue, J. Lee, J. Antognini, J. Sohl-Dickstein, R. Frostig, and G. E. Dahl. Measuring the effects of data parallelism on neural network training. *Journal of Machine Learning Research*, 20(112):1–49, 2019.
- N. Shazeer. Glu variants improve transformer. [arXiv preprint arXiv:2002.05202](#), 2020.
- M. Shoeybi, M. Patwary, R. Puri, P. LeGresley, J. Casper, and B. Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. [arXiv preprint arXiv:1909.08053](#), 2019.
- S. Smith, M. Patwary, B. Norick, P. LeGresley, S. Rajbhandari, J. Casper, Z. Liu, S. Prabhunoye, G. Zerveas, V. Korthikanti, et al. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. [arXiv preprint arXiv:2201.11990](#), 2022.
- S. L. Smith, P.-J. Kindermans, C. Ying, and Q. V. Le. Don’t decay the learning rate, increase the batch size. [arXiv preprint arXiv:1711.00489](#), 2017.

- J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- K. Sun, D. Yu, D. Yu, and C. Cardie. Investigating prior knowledge for challenging chinese machine reading comprehension, 2019.
- M. Suzgun, N. Scales, N. Schärli, S. Gehrmann, Y. Tay, H. W. Chung, A. Chowdhery, Q. V. Le, E. H. Chi, D. Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.
- H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C. Canton-Ferrer, M. Chen, G. Cucurull, D. Esiobu, J. Fernandes, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini, R. Hou, H. Inan, M. Kardas, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P. S. Koura, M. Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov, P. Mishra, I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten, R. Silva, E. M. Smith, R. Subramanian, X. E. Tan, B. Tang, R. Taylor, A. Williams, J. X. Kuan, P. Xu, Z. Yan, I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic, S. Edunov, and T. Scialom. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023b. doi: 10.48550/arXiv.2307.09288. URL <https://doi.org/10.48550/arXiv.2307.09288>.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Y. Wang, H. Li, X. Han, P. Nakov, and T. Baldwin. Do-not-answer: A dataset for evaluating safeguards in llms. *CoRR*, abs/2308.13387, 2023. doi: 10.48550/ARXIV.2308.13387. URL <https://doi.org/10.48550/arXiv.2308.13387>.
- J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html).
- T. Wei, J. Luan, W. Liu, S. Dong, and B. Wang. Cmath: Can your language model pass chinese elementary school math test?, 2023.
- L. Xu, H. Hu, X. Zhang, L. Li, C. Cao, Y. Li, Y. Xu, K. Sun, D. Yu, C. Yu, Y. Tian, Q. Dong, W. Liu, B. Shi, Y. Cui, J. Li, J. Zeng, R. Wang, W. Xie, Y. Li, Y. Patterson, Z. Tian, Y. Zhang, H. Zhou, S. Liu, Z. Zhao, Q. Zhao, C. Yue, X. Zhang, Z. Yang, K. Richardson, and Z. Lan. CLUE: A chinese language understanding evaluation benchmark. In D. Scott, N. Bel, and C. Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 4762–4772. International Committee on Computational Linguistics, 2020. doi: 10.18653/V1/2020.COLING-MAIN.419. URL <https://doi.org/10.18653/v1/2020.coling-main.419>.
- A. Yang, B. Xiao, B. Wang, B. Zhang, C. Yin, C. Lv, D. Pan, D. Wang, D. Yan, F. Yang, F. Deng, F. Wang, F. Liu, G. Ai, G. Dong, H. Zhao, H. Xu, H. Sun, H. Zhang, H. Liu, J. Ji, J. Xie, J. Dai,

- K. Fang, L. Su, L. Song, L. Liu, L. Ru, L. Ma, M. Wang, M. Liu, M. Lin, N. Nie, P. Guo, R. Sun, T. Zhang, T. Li, T. Li, W. Cheng, W. Chen, X. Zeng, X. Wang, X. Chen, X. Men, X. Yu, X. Pan, Y. Shen, Y. Wang, Y. Li, Y. Jiang, Y. Gao, Y. Zhang, Z. Zhou, and Z. Wu. Baichuan 2: Open large-scale language models. Technical report, Baichuan Inc., 2023. URL <https://cdn.baichuan-ai.com/paper/Baichuan2-technical-report.pdf>.
- L. Yu, W. Jiang, H. Shi, J. Yu, Z. Liu, Y. Zhang, J. T. Kwok, Z. Li, A. Weller, and W. Liu. Metamath: Bootstrap your own mathematical questions for large language models. *CoRR*, abs/2309.12284, 2023. doi: 10.48550/ARXIV.2309.12284. URL <https://doi.org/10.48550/arXiv.2309.12284>.
- R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi. HellaSwag: Can a machine really finish your sentence? In A. Korhonen, D. R. Traum, and L. Màrquez, editors, *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 4791–4800. Association for Computational Linguistics, 2019. doi: 10.18653/v1/p19-1472. URL <https://doi.org/10.18653/v1/p19-1472>.
- B. Zhang and R. Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.
- G. Zhang, L. Li, Z. Nado, J. Martens, S. Sachdeva, G. Dahl, C. Shallue, and R. B. Grosse. Which algorithmic choices matter at which batch sizes? insights from a noisy quadratic model. *Advances in neural information processing systems*, 32, 2019.
- C. Zheng, M. Huang, and A. Sun. Chid: A large-scale chinese idiom dataset for cloze test. In A. Korhonen, D. R. Traum, and L. Màrquez, editors, *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 778–787. Association for Computational Linguistics, 2019. doi: 10.18653/V1/P19-1075. URL <https://doi.org/10.18653/v1/p19-1075>.
- L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. P. Xing, H. Zhang, J. E. Gonzalez, and I. Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. 2023.
- W. Zhong, R. Cui, Y. Guo, Y. Liang, S. Lu, Y. Wang, A. Saied, W. Chen, and N. Duan. AGIEval: A human-centric benchmark for evaluating foundation models. *CoRR*, abs/2304.06364, 2023. doi: 10.48550/arXiv.2304.06364. URL <https://doi.org/10.48550/arXiv.2304.06364>.
- J. Zhou, T. Lu, S. Mishra, S. Brahma, S. Basu, Y. Luan, D. Zhou, and L. Hou. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.

## A. Appendix

### A.1. Acknowledgments

This project was realized thanks to the efforts of numerous contributors. We offer our extended thanks to the following individuals for their help<sup>1</sup>:

- Data Annotation Team: Jialu Cai, Ruijian Chen, Ruyi Chen, Bei Feng, Yanping Huang, Zhen Huang, Pin Jiang, Rongli Jin, Xiangyue Jin, Ziyun Ke, Hui Li, Meng Li, Sangsang Li, Xiaoqian Li, Yaohui Li, Yunxian Ma, Jiaqi Ni, Xiaojin Shen, Xinnan Song, Tianyu Sun, Xiaosha Chen, Haoyuan Tian, Xiaohan Wang, Xiaoxiang Wang, Yuhao Wang, Fanyi Xia, Lei Xu, Zeyuan Xu, Zhipeng Xu, Tian Yuan, Zhongyu Zhang, Yi Zheng, Shuang Zhou, Xinyi Zhou, Yuchen Zhu, Yuxuan Zhu.
- Compliance Team: Jin Chen, Ying Tang, Miaojun Wang, Xianzu Wang, Shaoqing Wu, Leyi Xia, W.L. Xiao.
- Business Team: Jian Liang, Mingming Li, T. Wang, Xianzu Wang, Zhiniu Wen, Shengfeng Ye, Peng Zhang, Zhen Zhang.
- Design Team: Wei An, Yukun Zha.

### A.2. Different Model Scale Representations

We refitted the scaling curve for different model scale representations, reusing the experiments from the IsoFLOP profile. We recalculated the compute FLOPs using  $6N_1$  and  $6N_2$  as model scale representations and refitted the performance scaling curves. As shown in Figure 6, the results indicate that the deviation of optimal model/data allocation among these three representations is not significant at higher compute budgets, but there are noticeable differences at lower budgets.

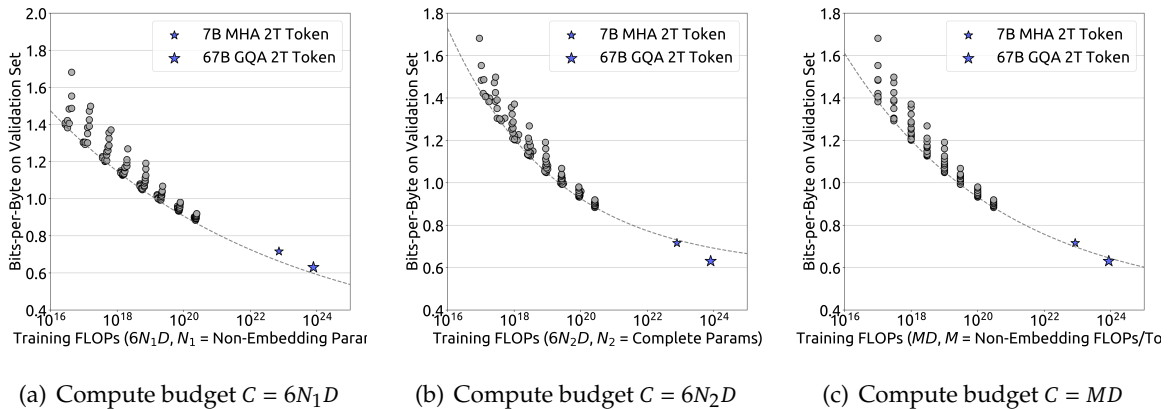


Figure 6 | Performance scaling curves using different model scale representations. The metric is the bits-per-byte on the validation set. The dotted line represents the power law fitting the smaller model (grey circles). The blue stars represent DeepSeek LLM 7B and 67B.  $N_1$ ,  $N_2$ , and  $M$  represent the non-embedding parameters, complete parameters, and non-embedding FLOPs/token of the model, respectively.

When using  $6N_1$  as the model scale representation, the fitted performance scaling curve tends to overestimate the performance of large-scale models. Conversely, when using  $6N_2$ , the

<sup>1</sup>Authors are ordered alphabetically by the last name.

curve tends to underestimate their performance. Using  $M$  as the model scale representation, however, achieves the most accurate predictions.

### A.3. Benchmark Metrics Curves

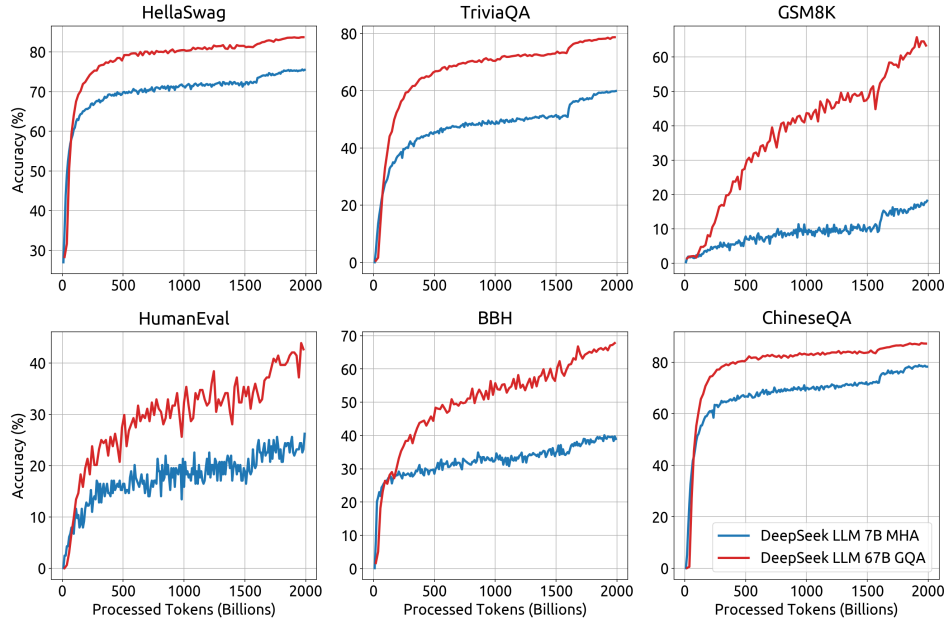


Figure 7 | Benchmark metrics curves of DeepSeek LLM Base. ChineseQA is our in-house test set, constructed in a manner akin to TriviaQA.

Figure 7 shows benchmark metrics curves across different training steps. We can see consistent improvement on these benchmarks from the start to the end of training. We believe the performance will further be improved if the training continues.

Model	Size	HumanEval		MBPP
		Python	Multilingual	
<b>Pre-Trained Models</b>				
Codex-001	-	33.5%	26.1%	45.9%
StarCoder	16B	36.0%	28.7%	46.8%
CodeGeeX2	6B	36.0%	24.5%	42.4%
CodeLlama	7B	31.7%	29.2%	41.6%
CodeLlama	13B	36.0%	35.4%	48.4%
CodeLlama	34B	<b>48.2%</b>	<b>41.0%</b>	55.2%
DeepSeek-LLM-Base	67B	42.7%	37.2%	<b>57.4%</b>
<b>Instruction-Tuned Models</b>				
Wizard-Coder	34B	73.2%	48.8%	61.2%
DeepSeek-LLM-Chat	67B	<b>73.8%</b>	<b>53.3%</b>	<b>61.4%</b>

Table 15 | Comparison with code-specific models.

#### A.4. Comparison with Code or Math Specific Models

We have conducted a comparison between our model and specific code and math language models (LLMs). Table 15 demonstrates that DeepSeek LLM 67B is capable of achieving similar performance to CodeLlama, despite having access to less code data. It is worth noting that DeepSeek LLM possesses greater capabilities in areas other than code.

Likewise, Table 16 presents the results obtained from various math-related benchmarks, such as GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), MGSM-zh (i et al., 2023), and CMath (Wei et al., 2023). DeepSeek 67B exhibits exceptional performance on math-related tasks across different languages, showcasing its superiority in this domain. In addition, DeepSeek LLM can utilize programs to solve math problems, which demonstrates better performance than chain-of-thoughts. It is significantly better than the previous SOTA model, ToRA (Gou et al., 2023), on the benchmarks.

	Inference	GSM8K	MATH	MGSM-zh	CMath
<b>Chain-of-Thoughts</b>					
MetaMath 70B (Yu et al., 2023)	CoT	82.3%	26.6%	66.4%	70.9%
WizardMath 70B (Luo et al., 2023)	CoT	81.6%	22.7%	64.8%	65.4%
DeepSeek LLM 67B Chat	CoT	<b>84.1%</b>	<b>32.6 %</b>	<b>74.0%</b>	<b>80.3%</b>
<b>Tool-Integrated Reasoning</b>					
ToRA-Code 34B (Gou et al., 2023)	Tool-Integrated	80.7%	50.8%	41.2%	53.4%
DeepSeek LLM 67B Chat	Tool-Integrated	<b>86.7%</b>	<b>51.1%</b>	<b>76.4%</b>	<b>85.4%</b>

Table 16 | Comparison with math-specific models.

#### A.5. Benchmark Results w/ DPO Stage

Table 17 presents the benchmark results obtained with the DPO stage. Based on these results, we can conclude that the DPO stage does not significantly impact the fundamental capability of an LLM.

	DeepSeek 67B Chat	DeepSeek 67B Chat DPO
HellaSwag	75.7	76.1
TriviaQA	81.5	82.9
NaturalQuestions	47.0	48.8
MMLU	71.1	70.9
GSM8K	84.1	85.2
MATH	32.6	30.2
HumanEval	73.8	71.3
BBH	71.7	70.8
AGIEval	46.4	46.1
CEval	65.2	64.3
CMMLU	67.8	68.2

Table 17 | The benchmark metrics before and after DPO stage.

#### A.6. Evaluation Formats

Table 18~Table 40 present examples of our evaluation formats on different benchmarks.



---

**PROMPT**

以下是一道中国高考生物选择题，请选择正确的答案。

问题：下列有关高尔基体、线粒体和叶绿体的叙述，正确的是选项：(A)三者都存在于蓝藻中(B)三者都含有DNA (C)三者都是ATP 合成的场所(D)三者的膜结构中都含有蛋白质

答案：从A到D, 我们应选择

---

Table 18 | An example of AGIEval.

---

**PROMPT**

Question: Use the information below to answer the question. Cotton is a plant product used to make fabric. Cotton is made of cellulose, a fiber not digestible by humans. Cellulose is composed of many sugar molecules bonded together into long chains. Each sugar molecule contains carbon, hydrogen, and oxygen atoms. When cotton fabric is washed, wrinkles often form. The clothing industry uses chemicals to manufacture some cotton fabrics that are wrinkle-free. Dyes are also added to color the cellulose fibers in cotton. How would a clothing manufacturer separate colors to determine the purity of the dyes?

Answer:

---

**OPTIONS**

- through filtration
  - by their boiling points
  - by their freezing points
  - through paper chromatography
- 

Table 19 | An example of ARC.

---

**PROMPT**

Evaluate the result of a random Boolean expression.

Q:  $\text{not} ( ( \text{not not True} ) )$  is

A: Let's think step by step.

Remember that (i) expressions inside brackets are always evaluated first and that (ii) the order of operations from highest priority to lowest priority is "not", "and", "or", respectively. We first simplify this expression "Z" as follows:  $Z = \text{not} ( ( \text{not not True} ) ) = \text{not} ( ( A ) )$  where  $A = \text{not not True}$ . Let's evaluate A:  $A = \text{not not True} = \text{not} ( \text{not True} ) = \text{not False} = \text{True}$ . Plugging in A, we get:  $Z = \text{not} ( ( A ) ) = \text{not} ( ( \text{True} ) ) = \text{not True} = \text{False}$ . So the answer is False.

Q: True and False and not True and True is

A: Let's think step by step.

Remember that (i) expressions inside brackets are always evaluated first and that (ii) the order of operations from highest priority to lowest priority is "not", "and", "or", respectively. We first simplify this expression "Z" as follows:  $Z = \text{True and False and not True and True} = A \text{ and B}$  where  $A = \text{True and False}$  and  $B = \text{not True and True}$ . Let's evaluate A:  $A = \text{True and False} = \text{False}$ . Let's evaluate B:  $B = \text{not True and True} = \text{not} ( \text{True and True} ) = \text{not} ( \text{True} ) = \text{False}$ . Plugging in A and B, we get:  $Z = A \text{ and B} = \text{False and False} = \text{False}$ . So the answer is False.

Q:  $\text{not not} ( \text{not} ( \text{False} ) )$  is

A: Let's think step by step.

Remember that (i) expressions inside brackets are always evaluated first and that (ii) the order of operations from highest priority to lowest priority is "not", "and", "or", respectively. We first simplify this expression "Z" as follows:  $Z = \text{not not} ( \text{not} ( \text{False} ) ) = \text{not not} ( A )$  where  $A = \text{not} ( \text{False} )$ . Let's evaluate A:  $A = \text{not} ( \text{False} ) = \text{not False} = \text{True}$ . Plugging in A, we get:  $Z = \text{not not} ( A ) = \text{not not} ( \text{True} ) = \text{not not False} = \text{True}$ . So the answer is True.

Q: False and False and False or not False is

A: Let's think step by step.

---

Table 20 | An example of BBH.

---

**PROMPT**

以下是中国关于教育学考试的单项选择题，请选出其中的正确答案。

根据我国心理学家冯忠良教授的学习分类，培养学生品德要通过\_\_\_\_\_。

- A. 知识的学习
- B. 技能的学习
- C. 行为规范的学习
- D. 态度的学习

答案: C

开设跨学科课程或建立跨学科专业体现了高等教育课程发展的\_\_\_\_\_。

- A. 综合化趋势
- B. 多样化趋势
- C. 人文化趋势
- D. 科学化趋势

答案: A

心智技能的特点有\_\_\_\_\_。

- A. 物质性、外显性、简缩性
- B. 观念性、内潜性、简缩性
- C. 物质性、外显性、展开性
- D. 观念性、内潜性、展开性

答案: B

下列关于大学生的情绪与理智关系的说法中正确的是\_\_\_\_\_。

- A. 能冷静控制自己情绪
- B. 感情用事，难以用理智控制情绪
- C. 遇事能坚持自己正确认识
- D. 已发展到不为小事而发怒和呕气

答案: B

在学完一篇逻辑结构严密的课文以后，勾画出课文的论点论据的逻辑关系图以帮助理解和记忆。这种学习方法属于\_\_\_\_\_。

- A. 精细加工策略
- B. 组织策略
- C. 复述策略
- D. 做笔记策略

答案: B

有学者强调，教育要根据一个民族固有的特征来定，这种观点体现了\_\_\_\_\_

- A. 生产力对教育的影响和制约
- B. 政治制度对教育的影响和制约
- C. 文化对教育的影响和制约
- D. 经济制度对教育的影响和制约

答案:

---

**OPTIONS**

- A
  - B
  - C
  - D
- 

Table 21 | An example of C-Eval.

---

**PROMPT**

女：这些药怎么吃？

男：一天三次，一次两片。

请根据上文回答问题：

他们在哪儿？

答案：

---

**OPTIONS**

- 商店
  - 饭店
  - 医院
  - 教室
- 

Table 22 | An example of C3.

---

**PROMPT**

以下是将某句古诗文翻译而成的现代表述：春天已至，万物复苏，春风如一位美丽而又心灵手巧的姑娘，迈着纤纤细步款款而来，她挥舞剪刀，尽情地展示那高超的女工技巧，她先裁出了柳叶，随着柳条袅袅依依地舞蹈，又裁出杏叶，桃叶。

该翻译所对应的古诗文是：

---

**OPTIONS**

- 春风骋巧如剪刀
  - 剪裁无巧似春风
  - 风吹怨恨快如刀
  - 春风欲擅秋风巧
- 

Table 23 | An example of CCPM.

---

**PROMPT**

Q: 某小学在“献爱心-为汶川地震区捐款”活动中，六年级五个班共捐款8000元，其中一班捐款1500元，二班比一班多捐款200元，三班捐款1600元，四班与五班捐款数之比是3: 5。四班捐款多少元？

A: 一班捐款1500元，而二班比一班多捐200元，所以二班捐款 $1500+200=1700$ 元，又知道六年级五个班一共捐款8000元，所以四班和五班捐款之和=一共捐款-一班和二班和三班捐款之和，即 $8000-1500-1700-1600=3200$ 元，而题目说四班与五班捐款数之比是3: 5，则四班捐款了 $3200/(3+5)*3=1200$ 元。所以答案是: 1200。

Q: 小俊在东西大道上跑步，若规定向东为正。他先向东跑了800米，然后又跑了一段之后，他位于出发点西边100米处，小俊第二段跑了多少米？

A: 小俊第二段跑完后位于出发点西边，所以第二段应该是向西跑，第二段跑的长度-第一段跑的长度=100，第二段跑了 $100+800=900$ 米。所以答案是: 900。

Q: A车和B车同时从甲、乙两地相向开出，经过5小时相遇。然后，它们又各自按原速原方向继续行驶3小时，这时A车离乙地还有135千米，B车离甲地还有165千米。甲、乙两地相距多少千米？

A: 假设A车的速度为x千米每小时，B车的速度为y千米每小时，根据而A、B相遇时A车行驶了5小时，A车行驶3小时后离乙地还有135千米，B车行驶3小时后距离甲地还有165千米，可以得到甲乙两地相距 $=5x+5y=135+8x=165+8y$ ，变换得到:  $10(x+y)=300+8(x+y)$ ，于是 $x+y=150$ ，甲乙两地相距 $5(x+y)=750$ 千米。所以答案是: 750。

Q: 在一个底面半径为10厘米的圆柱形容器内，倒入10厘米深的水，然后将一个底面直径4厘米，高6厘米的圆锥形铅锤放入水中，容器中水面上升多少厘米？

A:

---

Table 24 | An example of CMATH.

---

**PROMPT**

以下是关于解剖学的单项选择题，请直接给出正确答案的选项。

题目：壁胸膜的分部不包括

- A. 肋胸膜
- B. 肺胸膜
- C. 膈胸膜
- D. 胸膜顶

答案是：B

题目：属于蝶骨上的结构为

- A. 垂体窝
- B. 棘孔
- C. 破裂孔
- D. 视神经管

答案是：B

题目：属于右心房的结构是

- A. 肉柱
- B. 室上嵴
- C. 乳头肌
- D. 梳状肌

答案是：D

题目：咽的分部

- A. 咽隐窝
- B. 口咽部
- C. 鼻咽部
- D. 喉咽部

答案是：C

题目：舌下神经核位于

- A. 间脑
- B. 延髓
- C. 中脑
- D. 脑桥

答案是：B

题目：从脑干背侧出脑的脑神经是

- A. 副神经
- B. 三叉神经
- C. 舌下神经
- D. 滑车神经

答案是：

---

**OPTIONS**

- A
  - B
  - C
  - D
- 

Table 25 | An example of CMMLU.

---

**PROMPT**

Passage: The median age in the city was 22.1 years. 10.1% of residents were under the age of 18; 56.2% were between the ages of 18 and 24; 16.1% were from 25 to 44; 10.5% were from 45 to 64; and 7% were 65 years of age or older. The gender makeup of the city was 64.3% male and 35.7% female.

Answer the following questions based on the above passage, please calculate carefully if calculation is necessary.

Q: How many percent were not from 25 to 44?

A: The answer type is number. So according to above Passage, the answer is 83.9.

Q: How many in percent weren't 25 to 44?

A: The answer type is number. So according to above Passage, the answer is

---

Table 26 | An example of DROP.

---

**PROMPT**

中新网12月7日电综合外媒6日报道,在美国得克萨斯州,负责治疗新冠肺炎患者的医生约瑟夫·瓦隆(Joseph Varon)已连续上班超260天,每天只睡不超过2小时。瓦隆日前接受采访时呼吁,美国民众应遵从防疫规定,一线的医护人员“已

---

**OPTIONS**

- 神清气爽”。
  - 诡计多端”。
  - 精疲力竭”。
  - 分工合作”。
  - 寅吃卯粮”。
  - 土豪劣绅”。
  - 芸芸众生”。
- 

Table 27 | An example of CHID.

---

**PROMPT**

胡雪岩离船登岸，坐轿进城，等王有龄到家，他接着也到了他那里，脸上是掩抑不住的笑容，王有龄夫妇都觉得奇怪，问他什么事这么高兴。

上面的句子中的"他"指的是  
胡雪岩

渐渐地，汤中凝结出一团团块状物，将它们捞起放进盆里冷却，肥皂便出现在世上了。

上面的句子中的"它们"指的是  
块状物

“她序上明明引着JulesTellier的比喻，说有个生脱发病的人去理发，那剃头的对他说不用剪发，等不了几天，头毛压儿全掉光了；大部分现代文学也同样的不值批评。这比喻还算俏皮。”

上面的句子中的"他"指的是  
生脱发病的人

在洛伦佐大街的尽头处，矗立着著名的圣三一大教堂。它有着巨大的穹顶，还有明亮的彩色玻璃窗，上面描绘着《旧约》和《新约》的场景。

上面的句子中的"它"指的是  
圣三一大教堂

他伯父还有许多女弟子，大半是富商财主的外室；这些财翁白天忙着赚钱，怕小公馆里的情妇长日无聊，要不安分，常常叫她们学点玩艺儿消遣。

上面的句子中的"她们"指的是  
情妇

赵雨又拿出了一个杯子，我们热情地请老王入座，我边给他倒酒边问：1962年的哪次记得吗？“

上面的句子中的"他"指的是

---

Table 28 | An example of CLUEWSC.



---

**PROMPT**

Q: Max can mow the lawn in 40 minutes. If it takes him twice that long to fertilize the lawn, how long will it take him to both mow and fertilize the lawn?

A: Let's think step by step. It takes Max  $2 * 40$  minutes = 80 minutes to fertilize the lawn. In total, Max takes 80 minutes + 40 minutes = 120 minutes to both mow and fertilize the lawn. The answer is 120.

Q: The bagels cost \$2.25 each, or a dozen for \$24. How much is saved, per bagel, in cents, by buying a dozen at a time?

A: Let's think step by step. They cost  $2.25 * 100 = 225$  cents each. At the bulk rate, they are  $24 / 12 = 2$  dollar each. They cost  $2 * 100 = 200$  cents each.  $225 - 200 = 25$  cents are saved per bagel. The answer is 25.

Q: Tim is 5 years old. His cousin, Rommel, is thrice as old as he is. His other cousin, Jenny, is 2 years older than Rommel. How many years younger is Tim than Jenny?

A: Let's think step by step. Rommel is  $5 * 3 = 15$  years old. Jenny is  $15 + 2 = 17$  years old. So, Tim is  $17 - 5 = 12$  years younger than Jenny. The answer is 12.

Q: The school has 14 boys and 10 girls. If 4 boys and 3 girls drop out, how many boys and girls are left?

A: Let's think step by step. There are 14 boys - 4 boys = 10 boys left. There are 10 girls - 3 girls = 7 girls left. In total there are 10 boys + 7 girls = 17 boys and girls left. The answer is 17.

Q: Building one birdhouse requires 7 planks and 20 nails. If 1 nail costs 0.05, and one plank costs 3, what is the cost, in dollars, to build 4 birdhouses?

A: Let's think step by step. The cost of the planks for one birdhouse is  $7 * 3 = 21$ . And the nails are a cost of  $20 * 0.05 = 1$  for each birdhouse. So to build one birdhouse one will need  $21 + 1 = 22$ . So the cost of building 4 birdhouses is at  $4 * 22 = 88$ . The answer is 88.

Q: Danny brings 3 watermelons to his family picnic. He cuts each watermelon into 10 slices. His sister brings 1 watermelon to the family picnic, and she cuts the watermelon into 15 slices. How many watermelon slices are there in total at the picnic?

A: Let's think step by step. From Danny, there are  $3 * 10 = 30$  watermelon slices. From his sister, there are  $1 * 15 = 15$  watermelon slices. There are a total of  $30 + 15 = 45$  watermelon slices. The answer is 45.

Q: Angela is a bike messenger in New York. She needs to deliver 8 times as many packages as meals. If she needs to deliver 27 meals and packages combined, how many meals does she deliver?

A: Let's think step by step. Let  $p$  be the number of packages Angela delivers and  $m$  be the number of meals. We know that  $p + m = 27$  and  $p = 8m$ . Substituting the second equation into the first equation, we get  $8m + m = 27$ . Combining like terms, we get  $9m = 27$ . Dividing both sides by 9, we get  $m = 3$ . The answer is 3.

Q: Cori is 3 years old today. In 5 years, she will be one-third the age of her aunt. How old is her aunt today?

A: Let's think step by step. In 5 years, Cori will be  $3 + 5 = 8$  years old. In 5 years, Cori's aunt will be  $8 * 3 = 24$  years old. Today, her aunt is  $24 - 5 = 19$  years old. The answer is 19.

Q: Indras has 6 letters in her name. Her sister's name has 4 more letters than half of the letters in Indras' name. How many letters are in Indras and her sister's names?

A: Let's think step by step.

---

---

**PROMPT**

Playing piano: A man is seated at a piano. He

---

**OPTIONS**

- is playing the piano with his hands and his face.
  - begins to play a song by timbaland on the piano.
  - plays slowly, and pauses to snap his fingers.
  - is playing a song in front of him.
- 

Table 30 | An example of HellaSwag.

---

**PROMPT**

```
def starts_one_ends(n):
```

```
    """
```

Given a positive integer n, return the count of the numbers of n-digit positive integers that start or end with 1.

```
    """
```

---

Table 31 | An example of HumanEval.

---

**PROMPT**

Problem:

Find the domain of the expression  $\frac{\sqrt{x-2}}{\sqrt{5-x}}$ .

Solution:

The expressions inside each square root must be non-negative.

Therefore,  $x-2 \geq 0$ , so  $x \geq 2$ , and  $5-x \geq 0$ , so  $x \leq 5$ .

Also, the denominator cannot be equal to zero, so  $5-x > 0$ , which gives  $x < 5$ .

Therefore, the domain of the expression is  $\boxed{[2,5)}$ .

Final Answer: The final answer is  $[2,5)$ . I hope it is correct.

Problem:

If  $\det \mathbf{A} = 2$  and  $\det \mathbf{B} = 12$ , then find  $\det (\mathbf{A} \mathbf{B})$ .

Solution:

We have that  $\det (\mathbf{A} \mathbf{B}) = (\det \mathbf{A})(\det \mathbf{B}) = (2)(12) = \boxed{24}$ .

Final Answer: The final answer is  $24$ . I hope it is correct.

Problem:

Terrell usually lifts two 20-pound weights 12 times. If he uses two 15-pound weights instead, how many times must Terrell lift them in order to lift the same total weight?

Solution:

If Terrell lifts two 20-pound weights 12 times, he lifts a total of  $2 \cdot 12 \cdot 20 = 480$  pounds of weight. If he lifts two 15-pound weights instead for  $n$  times, he will lift a total of  $2 \cdot 15 \cdot n = 30n$  pounds of weight.

Equating this to 480 pounds, we can solve for  $n$ :

$$30n = 480$$

$$\Rightarrow n = 480/30 = \boxed{16}$$

Final Answer: The final answer is  $16$ . I hope it is correct.

Problem:

If the system of equations

$$\begin{aligned} 6x - 4y &= a \\ 6y - 9x &= b \end{aligned}$$

has a solution  $(x, y)$  where  $x$  and  $y$  are both nonzero, find  $\frac{a}{b}$ , assuming  $b$  is nonzero.

Solution:

If we multiply the first equation by  $-\frac{3}{2}$ , we obtain

$$-9x + 6y = -\frac{3}{2}a$$

Since we also know that  $6y - 9x = b$ , we have

$$-\frac{3}{2}a = b \Rightarrow \frac{a}{b} = \boxed{-\frac{2}{3}}$$

Final Answer: The final answer is  $-\frac{2}{3}$ . I hope it is correct.

Problem: Evaluate  $\log_2 21$ .

Solution:

---

---

**PROMPT**

You are an expert Python programmer, and here is your task: Write a function to find the similar elements from the given two tuple lists. Your code should pass these tests:

```
assert similar_elements((3, 4, 5, 6),(5, 7, 4, 10)) == (4, 5)
assert similar_elements((1, 2, 3, 4),(5, 4, 3, 7)) == (3, 4)
assert similar_elements((11, 12, 14, 13),(17, 15, 14, 13)) == (13, 14)
```

[BEGIN]

```
def similar_elements(test_tup1, test_tup2):
    res = tuple(set(test_tup1) & set(test_tup2))
    return (res)
```

[DONE]

You are an expert Python programmer, and here is your task: Write a python function to identify non-prime numbers. Your code should pass these tests:

```
assert is_not_prime(2) == False
assert is_not_prime(10) == True
assert is_not_prime(35) == True
```

[BEGIN]

```
import math
def is_not_prime(n):
    result = False
    for i in range(2,int(math.sqrt(n)) + 1):
        if n % i == 0:
            result = True
    return result
```

[DONE]

You are an expert Python programmer, and here is your task: Write a function to find the largest integers from a given list of numbers using heap queue algorithm. Your code should pass these tests:

```
assert heap_queue_largest( [25, 35, 22, 85, 14, 65, 75, 22, 58],3)==[85, 75, 65]
assert heap_queue_largest( [25, 35, 22, 85, 14, 65, 75, 22, 58],2)==[85, 75]
assert heap_queue_largest( [25, 35, 22, 85, 14, 65, 75, 22, 58],5)==[85, 75, 65, 58, 35]
```

[BEGIN]

```
import heapq as hq
def heap_queue_largest(nums,n):
    largest_nums = hq.nlargest(n, nums)
    return largest_nums
```

[DONE]

You are an expert Python programmer, and here is your task: Write a function to return the sum of all divisors of a number. Your code should pass these tests:

```
assert sum_div(8)==7
assert sum_div(12)==16
assert sum_div(7)==1
```

[BEGIN]

---

Table 33 | An example of MBPP.

---

**PROMPT**

The following are multiple choice questions (with answers) about miscellaneous.

How many axles does a standard automobile have?

- A. one
- B. two
- C. four
- D. eight

Answer: B

What place is named in the title of the 1979 live album by rock legends Cheap Trick?

- A. Budapest
- B. Budokan
- C. Bhutan
- D. Britain

Answer: B

Who is the shortest man to ever win an NBA slam dunk competition?

- A. Anthony 'Spud' Webb
- B. Michael 'Air' Jordan
- C. Tyrone 'Muggsy' Bogues
- D. Julius 'Dr J' Erving

Answer: A

What is produced during photosynthesis?

- A. hydrogen
- B. nylon
- C. oxygen
- D. light

Answer: C

Which of these songs was a Top 10 hit for the rock band The Police?

- A. 'Radio Ga-Ga'
- B. 'Ob-la-di Ob-la-da'
- C. 'De Do Do De Da Da Da'
- D. 'In-a-Gadda-Da-Vida'

Answer: C

Which of the Three Stooges was not related to the others?

- A. Moe
- B. Larry
- C. Curly
- D. Shemp

Answer:

---

**OPTIONS**

- A
  - B
  - C
  - D
- 

Table 34 | An example of MMLU.

---

**PROMPT**

Answer these questions:

Q: Who is hosting the fifa world cup in 2022?

A: Qatar

Q: Who won the first women 's fifa world cup?

A: United States

Q: When did miami vice go off the air?

A: 1989

Q: Who wrote the song shout to the lord?

A: Darlene Zschech

Q: Who was thrown in the lion 's den?

A: Daniel

Q: What is the meaning of the name habib?

A:

---

Table 35 | An example of NaturalQuestions.

---

**PROMPT**

A woman notices that she is depressed every autumn, and wonders why. A friend suggests to her that perhaps certain changes that take place as seasons move from warm to cold may be having an effect on her. When pressed for an example of these changes, the friend cites

---

**OPTIONS**

- flowers blooming
  - grass turning brown
  - trees growing
  - blossoms blooming
- 

Table 36 | An example of OpenBookQA.

---

**PROMPT**

To make it easier to push the reset button of the garbage disposable machine which is located underneath the machine,

---

**OPTIONS**

- place a wall mirror on the floor of the cabinet
  - hold a hand mirror under the garbage disposable machine
- 

Table 37 | An example of PIQA.

---

**PROMPT**

Article:

When you read an article you will understand and remember it better if you can work out how the writer has put the ideas together. Sometimes a writer puts ideas together by asking questions and then answering them. For example, if the article is about groundhogs, the set of questions in the writer's head might be:

What does a groundhog look like?

Where do groundhogs live?

What do they eat?...

In the article, the author might answer those questions.

Sometimes an author writes out her questions in the article. These questions give you signals. They tell you what the author is going to write next. Often an author has a question in her head but she doesn't write it out for you. You have to work out her question for yourself. Here's a sample reading for you to practice this method.

Earthworms

Do you know how many kinds of earthworms there are? There are about 1800 kinds in the world! They can be brown, purple, green. They can be as small as 3 cm long and as large as 3 m long.

The best time to see earthworms is at night, especially a cool, damp night. That's when they come up from their burrows to hunt for food. Earthworms don't like to be in the sun. That's because they breathe through their skin, and they can't breathe if their skin gets too dry. Earthworms must come out of the earth if it rains a lot, because they can't breathe in their flooded burrows. What a dangerous life!

Earthworms don't have eyes, so how can they tell when it's dark? They have special places on their skin that are sensitive to light. These spots tell whether it's light or dark. If you shine a flashlight on an earthworm at night, it will quickly disappear into the ground.

Earthworms don't have ears either, but they can hear by feeling movements in the earth. If you want to hear like an earthworm, lie on the ground with your fingers in your ears. Then have a friend stamp his or her feet near you. This is how earthworms feel birds and people walking, and moles digging, near them.

Earthworms are useful. Farmers and gardeners like having lots of earthworms in their land because the worms help to make better soil when they dig. That digging keeps the soil loose and airy. In one year earthworms can pile up as much as 23,000 kg of castings in an area about the size of a football field.

Q: What's the purpose of reading Earthworms?

A: To put the writer's idea into real use.

Q: Which question CANNOT be answered in the passage?

A: Why can human listen like earthworms?

Q: How can you understand Earthworms better according to this passage?

A: Read to work out all the questions in the writer's head while reading.

Q: What's the best title for the passage?

A:

---

**OPTIONS**

- One way to help with understanding
  - One way to practice with a new idea
  - One way to learn to be a wise writer
  - One way to be clearer about worms
-

---

**PROMPT**

Answer these questions:

Q: A Jayhawker was a term applied to anti-slavery militant bands from a certain US state that clashed with pro-slavery factions from Missouri. Which state is this, sometimes referred to as the Jayhawk State?

A: Kans.

Q: Which Swedish DJ and record producer had a UK Number One single in 2013 with 'Wake Me Up'?

A: Tim Bergling

Q: Who is the MP for Sheffield Hallam?

A: Nick clegg

Q: A case that riveted the nation, the case of The State of Tennessee v. John Thomas Scopes concluded on July 21, 1925, with the jury finding Mr. Scopes guilty of teaching what?

A: Survival of species

Q: What cartoon series featured a character called Little My?

A: Muumi

Q: "What English model, with her short-haired androgynous look, born Lesley Hornby, was discovered in 1966 by Nigel Davies when she was 16 and weighed 6 stone (41 kg, 91 lbs), and became ""The Face of '66"" with her high fashion mod look created by Mary Quant?"

A:

---

Table 39 | An example of TriviaQA.

---

**PREFIXES**

- So Monica

- So Jessica

---

**COMPLETION**

avoids eating carrots for their eye health because Emily needs good eyesight while Monica doesn't.

---

Table 40 | An example of WinoGrande. Note that there are multiple prefixes and only one completion for WinoGrande, and we choose the predicted prefix with the lowest perplexity of the completion.