

Instella: Fully Open Language Models with Stellar Performance

Jiang Liu, Jialian Wu, Xiaodong Yu, Yusheng Su, Prakamya Mishra, Gowtham Ramesh, Sudhanshu Ranjan, Chaitanya Manem, Ximeng Sun, Ze Wang, Pratik Prabhanjan Brahma, Zicheng Liu, Emad Barsoum

AMD



<https://huggingface.co/amd/Instella-3B>

<https://github.com/AMD-AGI/Instella>

Abstract

Large language models (LLMs) have demonstrated remarkable performance across a wide range of tasks, yet the majority of high-performing models remain closed-source or partially open, limiting transparency and reproducibility. In this work, we introduce Instella, a family of fully open three billion parameter language models trained entirely on openly available data and codebase. Powered by AMD Instinct™ MI300X GPUs, Instella is developed through large-scale pre-training, general-purpose instruction tuning, and alignment with human preferences. Despite using substantially fewer pre-training tokens than many contemporaries, Instella achieves state-of-the-art results among fully open models and is competitive with leading open-weight models of comparable size. We further release two specialized variants: Instella-Long, capable of handling context lengths up to 128K tokens, and Instella-Math, a reasoning-focused model enhanced through supervised fine-tuning and reinforcement learning on mathematical tasks. Together, these contributions establish Instella as a transparent, performant, and versatile alternative for the community, advancing the goal of open and reproducible language modeling research.

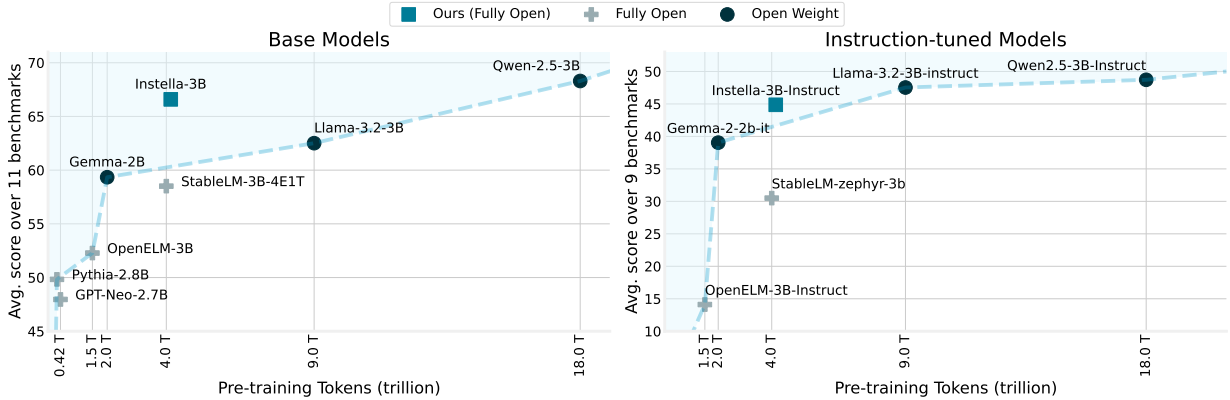


Figure 1: **Average Score versus Pre-training Tokens** for base (left) and instruction-tuned (right) models. Instella surpasses prior fully open models of comparable size and, despite being trained on substantially fewer pre-training tokens, achieves competitive performance with state-of-the-art open-weight models for both (left) base models (Table 4) and (right) instruction-tuned models (Table 6).

1 Introduction

The rapid advancement of artificial intelligence, driven in large part by large language models (LLMs) (Gemini Team, 2024; OpenAI, 2023; Dubey et al., 2024; Yang et al., 2025a), has accelerated progress toward artificial general intelligence and transformed society at large. However, much of this progress has been led by proprietary releases (e.g., GPT-4 (OpenAI, 2023), Claude (Anthropic, 2025), Gemini (Gemini Team, 2024)), where training data, methods, and evaluation details remain opaque. While these models have set new state-of-the-art performance, their closed nature hinders scientific understanding, reproducibility, and equitable access.

In response, the research community has placed increasing emphasis on open-weight models, where trained parameters are released. Projects such as LLaMA-3.2-3B (Dubey et al., 2024), Qwen-2.5-3B (Yang et al., 2024), and Gemma-2-2B (Team et al., 2024) have demonstrated competitive capabilities in relatively compact architectures. Yet most of these remain open-weight rather than fully open: their training data, preprocessing, and training recipes are either undisclosed or proprietary. As a result, researchers cannot fully reproduce the results, audit potential data contamination, or study the effects of data and training choices at scale.

To bridge this gap, we introduce Instella, a new family of fully open 3B-parameter language models. Instella makes available not only model weights, but also the complete training pipeline, datasets, and optimization details, thereby offering full transparency. Instead of solely relying on general-purpose corpora, Instella is pretrained in two distinct stages: an initial 4T-token general-domain pre-training stage, followed by a 57B-token second-stage emphasizing reasoning-heavy domains. To further enrich this stage, we introduce an in-house synthetic dataset for mathematics, constructed by abstracting GSM8K problems into symbolic Python programs and parameterizing them to generate diverse yet solvable variants. This approach expands mathematical coverage while maintaining the correctness of synthesized data, providing a principled way to inject reasoning signals into pre-training. In addition, we leverage weight ensembling across stochastic pre-training seeds by conducting multiple second-stage runs with different random seeds and merging their weights into the final checkpoint, which further enhances model performance. Following pre-training, Instella undergoes supervised fine-tuning (SFT) on a carefully curated mixture of 2.3 million high-quality instruction-response pairs drawn from diverse domains such as mathematics, coding, commonsense reasoning, and multi-turn dialogue. This step equips the model with the ability to follow user prompts, handle complex instructions, and generalize across a wide range of task formats, and is further refined through direct preference optimization (DPO) (Rafailov et al., 2023), aligning outputs with human expectations for helpfulness, safety, and factuality.

Building on this foundation, we extend Instella into the long-context regime with Instella-Long, capable of processing sequences up to 128K tokens. Instella-Long is trained in two stages of continued pre-training on 40B tokens, followed by long-context SFT and short-context DPO. Because of the limited availability of long-context SFT data, we synthesize long-context instruction-following examples directly from pre-training documents. Compared with other open-weight models, Instella-Long delivers competitive performance on the challenging Helmet benchmark (Yen et al., 2024), while fully releasing its training details and data to ensure transparency and reproducibility.

Finally, Instella advances reasoning-centric reinforcement learning at small scale through Instella-Math. Using only 3B parameters, Instella-Math is, to our knowledge, the first fully open model of this size to apply multi-stage group relative policy optimization (GRPO) (Shao et al., 2024) entirely on open datasets. By gradually increasing rollout lengths and incorporating Olympiad-level problems from DeepScaleR (Luo et al., 2025), the model demonstrates substantial improvements in mathematical and logical reasoning. Remarkably, Instella-Math performs strongly not only on benchmarks like GSM8K and OlympiadBench (He et al., 2024b) but also on TTT-Bench (Mishra et al., 2025), highlighting that reinforcement learning can meaningfully enhance reasoning even for compact models.

Despite being trained on significantly fewer tokens compared to some leading models, Instella achieves state-of-the-art results among fully open models and rivals the performance of stronger open-weight models. To summarize, our contributions are threefold:

- **Instella.** A 3B-parameter language transformer trained with a carefully staged pre-training process. Instella significantly outperforms prior fully open models of comparable size across diverse benchmarks.
- **Instella-Long.** A long-context variant extending sequence length to 128K tokens driven by continued pre-training and synthetic QA-based long-context instruction tuning. Instella-Long attains competitive performance on the challenging long-context benchmark Helmet.
- **Instella-Math.** A reasoning-centric variant fine-tuned with curated math datasets and reinforcement learning, delivering strong gains on AIME, OlympiadBench, and GSM8K while achieving the highest reported performance on the strategic reasoning benchmark TTT-Bench among fully open models.

Our work demonstrates that openness and competitiveness are not mutually exclusive. By releasing model weights, training code, data recipes, and evaluation protocols, Instella enables transparent benchmarking, reproducibility, and further research into the foundations of language modeling.

2 Background

2.1 Open-Weight versus Fully-Open Large Language Models

The release of open-weight large language models such as LLaMA (Touvron et al., 2023; Dubey et al., 2024) and Qwen (Bai et al., 2023; Yang et al., 2024; 2025a) series has significantly broadened community access to high-performing models. These systems are compact enough to be fine-tuned on modest hardware, enabling academic research and downstream applications. However, most such models are not completely transparent: their pre-training datasets, training pipelines, and optimization hyperparameters remain undisclosed. This opacity prevents reproducibility, makes data contamination difficult to audit, and constrains the ability to study scaling laws or understand how training data composition affects downstream performance.

In contrast, completely transparent models release not only weights but also data recipes, preprocessing scripts, and training code. Notable examples include OLMo (Groeneveld et al., 2024; OLMo et al., 2024) and SmoLLM (Allal et al., 2025), which provide comprehensive training pipelines and fully specified data mixtures. These initiatives enable researchers to systematically investigate questions such as how data diversity affects generalization, how alignment methods interact with model size, and how pre-training choices influence reasoning capabilities. However, prior fully open 3B models still underperform compared to state-of-the-art open-weight systems by a considerable margin on challenging benchmarks such as GSM8K (Cobbe et al., 2021b), BBH (Suzgun et al., 2023), and MMLU (Hendrycks et al., 2021b), motivating further work to bridge the gap between transparency and competitiveness. Instella addresses this gap by offering a fully open 3B-parameter model family with state-of-the-art results. We release not only weights but also training data recipes, preprocessing scripts, optimization settings, and evaluation pipelines, providing a truly reproducible foundation for scientific study.

2.2 Long-context Language Models

Many real-world applications demand reasoning over inputs significantly longer than the typical 2K–8K context windows used in base large language models. Tasks such as legal document analysis, multi-chapter summarization, and retrieval-augmented generation require context lengths exceeding 100K tokens. Recent advances including efficient attention mechanisms (Dao, 2024; Jacobs et al., 2023; Liu et al., 2023), rotary position embedding (RoPE) scaling (Gradient Team, 2024; emozilla, 2023; Ding et al., 2024), and specialized training strategies for long sequences (Gao et al., 2024) have enabled models to process extended sequences. Despite these developments, few transparent models provide both long-context support and strong performance. On the other hand, open-weight models such as Qwen2.5-1M (Yang et al., 2025b) offer extended context windows, but their training data remain proprietary, limiting reproducibility. Instella-Long contributes to this space by transparently extending the context length to 128K tokens through continued pre-training and post-training on the long-context data we release publicly. It achieves competitive results on the long-context benchmarks while establishing a transparent, reproducible long-context baseline.

2.3 Large Reasoning Models

The ability to perform multi-step reasoning represents a central goal for large language model development. Benchmarks such as MMLU, BBH, GSM8K, MATH (Hendrycks et al., 2021d) and AIME (AIME) measure a model’s capacity to perform structured, compositional thinking beyond surface-level pattern matching. Recent research demonstrates that high-quality reasoning data and post-training techniques such as reinforcement learning can dramatically improve performance. Models like DeepSeek-R1 (DeepSeek-AI et al., 2025) and DeepSeek-Math (Shao et al., 2024) show that incorporating step-by-step solutions and applying alignment methods like group relative policy optimization (GRPO) (Shao et al., 2024) can lead to substantial gains in reasoning capabilities.

However, most reasoning-focused models remain only partially open: either the reasoning datasets are proprietary, the reinforcement learning recipes are undisclosed, or the resulting models are released without reproducible training pipelines. This lack of transparency hinders systematic study of reasoning capabilities and prevents independent validation of methodological claims.

Instella-Math addresses this limitation by providing the first fully open 3B-parameter model trained with multi-stage reinforcement learning entirely on open data. We release not only the model weights but also

the reasoning datasets and training configurations, enabling reproducible research into reasoning emergence and reinforcement learning training for small-scale models.

3 Instella

3.1 Model Architecture

The Instella models are text-only, autoregressive transformer-based language models (Vaswani et al., 2017) with 3 billion parameters. Architecture-wise, Instella consists of 36 decoder layers, each having 32 attention heads with a hidden dimension of 2,560 and an intermediate dimension of 6,912. We use standard multi-head attention (Vaswani et al., 2017). For layer normalization, we employ RMSNorm (Zhang & Sennrich, 2019), which has been shown to provide better training stability and convergence properties compared to standard LayerNorm (Ba et al., 2016), particularly for large-scale language models (Takase et al., 2025; Touvron et al., 2023; Muennighoff et al., 2025).

In addition, we apply QK-Norm (Dehghani et al., 2023; Muennighoff et al., 2025; Naseer et al., 2021), where layer normalization is injected after the query and key projections within each attention head. QK-Norm normalizes the query and key vectors before computing attention scores, helping to maintain more balanced attention distributions throughout training. It has been shown to be effective in improving training stability by preventing attention weights from becoming overly extreme, which can lead to gradient instability and poor convergence.

Our model uses a standard causal attention mask. The feed-forward network within each transformer layer follows the standard architecture with SwiGLU activation function, which has demonstrated superior performance compared to ReLU-based activations in recent language models. We also employ rotary position embeddings (RoPE) (Su et al., 2024) to encode positional information, which provides better extrapolation to longer sequences compared to absolute positional embeddings.

The key hyperparameters of Instella-3B architecture are shown in Table 1. We use the OLMo tokenizer (Groeneveld et al., 2024) with a vocabulary size of 50,304 tokens. This vocabulary size strikes a balance between computational efficiency and representation capacity, allowing the model to handle diverse text while maintaining reasonable embedding and output layer sizes.

Table 1: Key hyper-parameters of Instella-3B architecture.

Number of transformer layers	Hidden dimension	Intermediate dimension	Number of attention heads	Number of KV heads	Sequence length	Vocabulary size
36	2560	6912	32	32	4096	50,304

3.2 Training Setup

Our training pipeline is based on the open-sourced OLMo codebase, adapted, and optimized for our hardware and model architecture. For pre-training we use a total of 128 Instinct MI300X GPUs distributed across 16 nodes. During both pre-training and post-training, we utilize FlashAttention 2 (Dao, 2024), Torch Compile, and bfloat16 mixed-precision training to reduce memory usage and speed up training. To balance inter-node memory efficiency and intra-node communication overhead within our cluster, we employ fully sharded data parallelism (FSDP) with hybrid sharding, with model parameters, gradients, and optimizer states sharded within a node and replicated across the nodes.

3.3 Pre-training

We pre-train the model using two stages with a sequence length of 4,096 tokens and a global batch size of 1,024. The Instella 3B pretraining pipeline is shown in Fig. 2. In the first pre-training stage, we train the model from scratch on 4.07 trillion tokens sourced from OLMoE-mix-0924 (Muennighoff et al., 2025), which is a diverse mix of two high-quality datasets DCLM-baseline (Li et al., 2024) and Dolma 1.7 (Soldaini et al., 2024) covering domains like coding, academics, mathematics, and general world knowledge from web crawl. This extensive first stage pre-training established a foundational understanding of general language

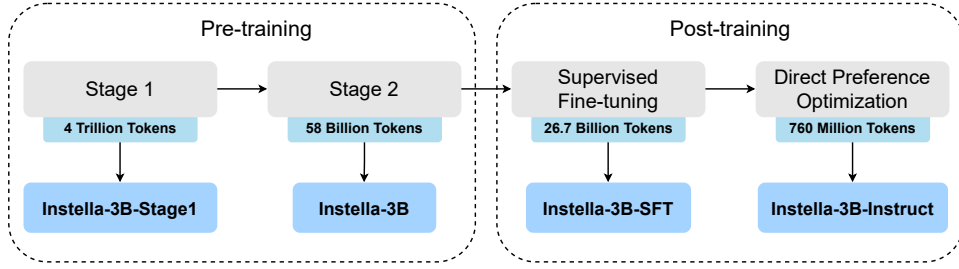


Figure 2: Instella-3B model training pipeline.

in our Instella model. We use the cosine decay learning rate schedule with a maximum learning rate of 4×10^{-4} and set the global batch size to 1024.

For our final pre-trained checkpoint, Instella-3B, we conduct a second stage pre-training on top of the first-stage Instella-3B-Stage1 model to further enhance its capabilities on MMLU (Hendrycks et al., 2021b), BBH (Suzgun et al., 2023), and GSM8K (Cobbe et al., 2021b). The model is trained three times with different random seeds, and the resulting weights are ensembled to obtain the final checkpoint. Specifically, the second-stage training uses 58 billion tokens sourced from diverse and high-quality datasets, including Dolmino-Mix-1124 (OLMo et al., 2024), SmolLM-Corpus (python-edu) (Ben Allal et al., 2024), Deepmind Mathematics (Saxton et al., 2019), and conversational datasets such as Tulu-3-SFT-Mixture (Lambert et al., 2024), OpenHermes-2.5 (Teknium, 2023), WebInstructSub (Yue et al., 2024), Code-Feedback (Zheng et al., 2024), and Ultrachat 200k (Ding et al., 2023). We use the linear decay learning rate schedule with a maximum learning rate of 4×10^{-5} and set the global batch size to 1024.

In addition to the publicly available datasets, 28.5 million tokens in the second-stage pre-training data mixture are derived from our in-house synthetic dataset focused on mathematical problems. This dataset is generated using the training set of GSM8k dataset, where we first use Qwen2.5-72B-Instruct (Yang et al., 2024) to 1) abstract numerical values as function parameters and generate a python program to solve the math question, 2) identify and replace numerical values in the existing question with alternative values that are still answerable with the same python program solution as the original question. Next, by assigning different new values to these python parameters and using the abstract solution program to compute the corresponding answers, we expand our synthetic dataset with new and reliable question-answer pairs (Yu et al., 2024).

3.4 Post-training

We first perform supervised finetuning (SFT) to enable the pre-trained model to follow instructions and respond effectively to user queries. We train for three epochs on 2.3 millions of high-quality instruction-response pairs, resulting in Instella-3B-SFT. During this phase, we utilize datasets spanning a broad spectrum of tasks and domains to ensure that the model generalizes across diverse instruction types. The mixture is selectively sourced from SmolTalk (Allal et al., 2025), OpenMathInstruct-2 (Toshniwal et al., 2024), Tulu-3 Instruction Following (Lambert et al., 2024), MMLU auxiliary train set (Hendrycks et al., 2021b), and o1-journey (Qin et al., 2024). We use the linear decay learning rate schedule with a maximum learning rate of 1×10^{-5} and set the global batch size to 128.

In the final training stage, we align Instella-3B-SFT with human preferences to ensure its outputs are helpful, accurate, and safe. Building on Instella-3B-SFT, Instella-3B-Instruct is trained with direct preference optimization (DPO) (Rafailov et al., 2023) on 0.76 billion tokens from the OLMo 2 1124 7B Preference Mix (OLMo et al., 2024). This alignment step tailors the model’s responses to better reflect human values and expectations, thereby improving the quality and reliability of its outputs. We use the linear decay learning rate schedule with a maximum learning rate of 5×10^{-7} and set the global batch size to 128.

4 Instella-Long

In this section, we introduce the long-context model of Instella, namely, Instella-3B-Long-Instruct, supporting 128K context length. To extend the context length, we continually train the model from Instella-

3B-Instruct through: 1. continued pre-training, 2. supervised finetuning (SFT), and 3. direct preference optimization (DPO), as shown in Fig. 3. We detail the training method and data in the following subsections.

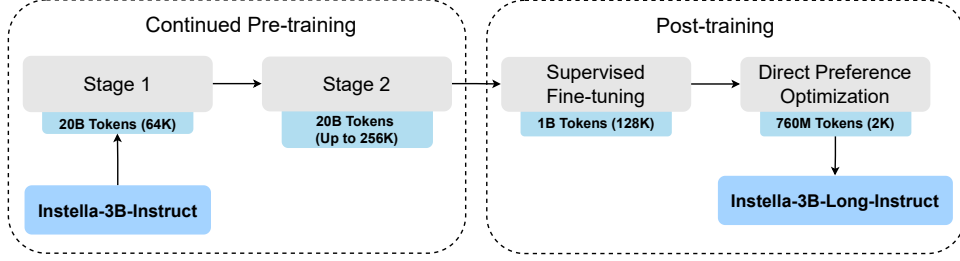


Figure 3: Instella-Long model training pipeline.

4.1 Continued Pre-training

The long context training is initialized from the short-context checkpoint, Instella-3B-Instruct, which has a context length of 4K. We conduct a two-stage continued pre-training to gradually increase the context length. **Stage 1:** We extend the context length from 4K to 64K and train the model using 20B tokens. The batch size is 4M tokens and the training steps are 5,000. We follow the RoPE scaling law (Gradient Team, 2024) to increase the base frequency of RoPE from 10,000 to 514,640. We also experiment with alternative RoPE scaling methods (emozilla, 2023; Gao et al., 2024) and observe only minor differences in performance. **Stage 2:** As indicated by (Gao et al., 2024), it is beneficial to train the model with the data whose context length is longer than the target context length. In this stage, we train the model on 20B tokens with a maximum context length of 256K - twice our target context length of 128K. Following the RoPE scaling law, we further increase the RoPE base frequency to 3,691,950. The batch size is 8M tokens and the training steps are 2,500. For both stages, we use the linear decay learning rate schedule and the maximum learning rate is 2×10^{-5} .

Table 2: Long-context continued pre-training data by source and portion. Each stage consists of 20 billion tokens in total.

Training Stage	64K Long Data	256K Long Data	Short Data
Stage 1	Code repos (30%) Books (30%) Textbooks (3%)	–	FineWeb-Edu (10%) FineWeb (10%) Wikipedia (5%) OpenWebMath (5%) StackExchange (4%) ArXiv (3%)
Stage 2	Code repos (10%) Books (15%)	Code repos (20%) Books (15%) Textbooks (2%)	FineWeb-Edu (10%) FineWeb (10%) Wikipedia (5%) OpenWebMath (5%) StackExchange (4%) ArXiv (4%)

The continued pre-training data originates from the data mixture created by Prolong (Gao et al., 2024). We use the raw text data curated by Prolong and process the data through tokenization, filtering, and packing. In each stage of the continued pre-training, we train on a 20B-token mixture of short- and long-context data with an approximate ratio of 4 to 6. The detailed data sources and portion are listed in Table 2. Let L be the maximum context length of the training stage. We pack both short- and long-context data into L -length sequences for training. For short-context data, we randomly select multiple documents and concatenate them into an L -length sequence. The extra texts beyond L in the last document are discarded. For long-context data, we filter out the documents that are shorter than L . We observe that the raw text data has some super long documents ($\gg L$). For these documents, we randomly sample a few segments from them to avoid producing an excessive number of training examples from a single document. We mix 64K data into

the long-context data in the second stage for improving training throughput, where we pack four different 64K documents into a 256K sequence. During data processing, we ensure that the documents used in the first and second stages are mutually exclusive. In training, we apply document masking so that different documents within the same sequence cannot attend to each other.

4.2 Post-training

After continued training on the long-context pre-training data, we perform supervised finetuning on a 1B-token mixture of short- and long-context instruction data. We use a batch size of 4M tokens and train for 250 steps. A linear decay learning rate schedule is employed, with a maximum learning rate of 4×10^{-5} . For the SFT data, we pack multiple samples into a 256K sequence with document masking applied during training. Padding tokens are added in order to reach exactly 256K tokens.

Similar to the continued pre-training, we train the model on a mixture of short- and long-context instructions data with a ratio of 4 to 6. For short-context instruction data, we use publicly available instruction-tuning datasets, some of which are also used in the post-training of Instella-3B-Instruct. Specifically, we use Ultrachat 200K (Ding et al., 2023), OpenMathinstruct-2 (Toshniwal et al., 2024), Tulu-3 Instruction Following (Lambert et al., 2024), and MMLU auxiliary train set (Hendrycks et al., 2021b).

Due to the lack of long-context SFT data, we construct a long-context instruction-following dataset where the context length is controlled to be between 8K and 128K tokens. Specifically, we make use of the long-context documents of Books from our continued pre-training data corpus. We use the documents that have at least 8K tokens and truncate the document to 128K tokens if it is over 128K. Then, we use Qwen2.5-14B-Instruct-1M (Yang et al., 2025b) as a teacher model to synthetically generate a question and an answer for the document. To speed up this process, we randomly choose a subpart of the document for the QA generation instead of using the whole document. The length of the subpart is randomly set to be between 2K and 8K tokens. We use NLTK Bird & Loper (2004) sentence tokenizer to divide documents into sentences to make sure that the selected subpart has complete sentences. The generated question and answer are appended to the end of the long document, serving as a complete single-round instruction-following data sample. Furthermore, we generate long-context instruction data from short-context documents, thereby enhancing dataset diversity with a broader range of sources. We use ArXiv from our continued pre-training corpus and the DCLM subset from Dolmino-Mix-1124 (OLMo et al., 2024). We first generate QA for each short-context document following the same pipeline aforementioned. Next, we iteratively concatenate different short-context documents into a long sequence until it reaches 128K tokens. Since we do not truncate the last document, the concatenated sequence may exceed 128K tokens. Lastly, we randomly choose one QA corresponding to one of the short-context documents and append it to the end of the concatenated sequence. Contrary to the findings by (Gao et al., 2024), we observe that our synthetic long-context instruction data notably improves performance on long-context tasks. The final SFT data mixture is shown in Table 3.

Table 3: Long-context supervised finetuning data by source and portion, totaling 1 billion tokens.

Short Data	Long Data
Ultrachat 200K (25%), OpenMathinstruct-2 (10%), MMLU auxiliary train set (3%), Tulu-3 Instruction Following (2%)	Books (44%), DCLM (10%), ArXiv (6%)

In the final training stage, we perform human preference alignment using DPO (Rafailov et al., 2023), employing the same training setting and dataset as Instella-3B-Instruct. Different from the previous long-context training stages, this DPO stage is trained on short-context data only with a maximum context length of 2K. Consistent with the findings of other open-weights models, we observe that applying DPO solely on short-context data continues to improve on long-context tasks.

4.3 Implementation Details

Sequence Parallelism. We implement sequence parallelism based on DeepSpeed Ulysses (Jacobs et al., 2023), which distributes the attention heads across GPUs during attention computation. Compared to Ring-Attention (Liu et al., 2023), this approach is more communication-efficient. For the second continued pre-training stage and SFT, we employ four GPUs as a sequence parallelism group to handle the long input

sequences. Sequence parallelism is not used in other stages, as the memory requirements fit within a single GPU.

Document Masking and Data Batching. We apply document masking during the continued pre-training and SFT, as each input sequence may contain multiple documents. Document masking is achieved through variable-length FlashAttention (Dao, 2023), which computes attention within each individual document rather than across the entire sequence. This design can also improve training throughput when combined with sorted data batching. Following Prolong (Gao et al., 2024), we sort microbatches at each training step by the sum of document lengths in the sequence. With gradient accumulation, later microbatches benefit from faster processing when they consist of shorter documents.

5 Instella-Math

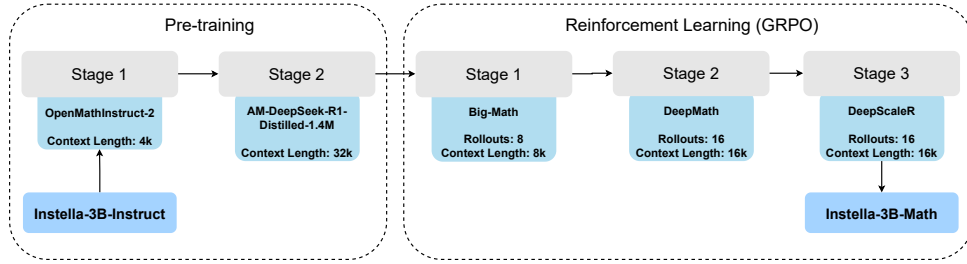


Figure 4: Instella-Math model training pipeline.

In this section, we introduce Instella-Math, a reasoning-centric language model trained with long chain-of-thought reinforcement learning. To enhance the model’s mathematical and logical reasoning capabilities, we continually train Instella-3B-Instruct through two stages of supervised finetuning and three stages of reinforcement learning, as shown in Figure 4. We detail the training procedure and datasets below.

5.1 Supervised Finetuning

As a cold start, we perform a two-stage supervised finetuning process to enhance the reasoning capabilities of Instella-3B-Instruct:

Stage 1: Instruction Tuning with OpenMathInstruct-2 for Mathematical Coverage. In the first SFT stage, we begin with instruction tuning, following instructions or prompts properly, especially in a question-answer or problem-solution format. Using the OpenMathInstruct-2 dataset (Toshniwal et al., 2024), which consists of 14 million problem-solution pairs generated from the GSM8K (Cobbe et al., 2021b) and MATH (Hendrycks et al., 2021d) training sets, the model is trained to solve mathematical questions covering a diverse range of topics from arithmetic and algebra to probability and calculus.

Stage 2: Deep Reasoning with Long-Context Training on AM-DeepSeek-R1-Distilled. In the second SFT stage, we further improve the model’s reasoning capability by training on AM-DeepSeek-R1-Distilled-1.4M (Zhao et al., 2025), which is a large-scale general reasoning dataset containing high-quality and challenging problems. In this stage, we increase the context length of the model from 4K to 32K to allow the model to learn from the long chain-of-thought responses distilled from large reasoning models such as DeepSeek-R1 (DeepSeek-AI et al., 2025).

5.2 Reinforcement Learning

Following supervised finetuning, we apply three stages of reinforcement learning using the group relative policy optimization (GRPO) algorithm (Shao et al., 2024) to further strengthen the model’s mathematical reasoning abilities. Training is orchestrated with verl (Sheng et al., 2024) and vLLM (Kwon et al., 2023) for efficient rollout collection, reward scoring, and policy updates.

Stage 1: GRPO on Big-Math-RL-Verified (8 Rollouts @ 8K Tokens). In the first stage of reinforcement learning, we apply the GRPO algorithm to train the model on Big-Math-RL-Verified (Albalak et al., 2025), a collection of curated, complex, multi-step math problems. We generate 8 rollouts per prompt, each with up

Table 4: **Base model performance.**

Models	ARC-C	ARC-E	BoolQ	HS.	PiQA	SciQ	WG.	OBQA	MMLU	BBH	GSM8K	Avg.
<i>Open Weight Models</i>												
Gemma2-2B	39.5	59.3	74.5	70.5	76.4	96.6	69.8	44.8	53.3	40.8	27.4	59.3
Llama-3.2-3B	47.2	64.9	74.8	73.1	75.9	95.3	70.3	51.2	57.8	47.0	30.1	62.5
Qwen2.5-3B	51.5	67.2	79.1	72.1	77.4	95.5	69.3	51.4	67.2	56.7	63.8	68.3
<i>Fully Open Models</i>												
Pythia-2.8B	40.5	60.7	64.8	60.1	72.5	89.7	60.8	42.6	26.1	27.7	2.7	49.8
GPTNeo-2.7B	38.5	54.6	62.7	55.2	70.8	88.0	58.3	40.8	27.8	27.3	3.7	48.0
OpenELM-3B	37.5	58.4	68.6	71.7	75.6	92.5	65.4	46.4	26.7	29.4	3.0	52.3
StableLM-3B	44.8	67.0	75.4	74.2	78.4	93.4	68.4	48.6	45.2	37.3	10.8	58.5
Instella-3B-Stage1	53.9	73.2	78.7	74.2	77.5	94.9	71.2	51.4	54.7	34.3	10.8	61.3
Instella-3B	52.8	70.5	76.5	75.0	77.8	96.4	73.1	52.4	58.3	39.7	59.8	66.6

Table 5: **Instella 3B base model performance.** We report the model performance after stage 1 and stage 2 pretraining. For stage 2, we run the training for three times with different random seeds and merge model weights to obtain the final stage 2 model.

Models	ARC-C	ARC-E	BoolQ	HS.	PiQA	SciQ	WG.	OBQA	MMLU	BBH	GSM8K	Avg.
Stage1	53.9	73.2	78.7	74.2	77.5	94.9	71.2	51.4	54.7	34.3	10.8	61.3
Stage2-seed1	51.2	68.8	76.2	73.8	77.3	96.6	72.1	52.0	57.7	38.5	56.1	65.5
Stage2-seed2	50.8	68.4	77.8	74.3	77.2	96.6	71.8	51.4	58.2	38.5	58.8	65.8
Stage2-seed3	49.8	68.8	73.5	75.6	77.2	96.7	72.8	52.0	58.0	38.6	58.3	65.6
Stage2	52.8	70.5	76.5	75.0	77.8	96.4	73.1	52.4	58.3	39.7	59.8	66.6

to 8K output tokens, to explore diverse reasoning trajectories. The model is trained for 1,200 GRPO steps using rule-based reward signals provided by Prime-RL (Cui et al., 2025), which incentivize correctness and well-structured outputs.

Stage 2: GRPO on DeepMath (16 Rollouts @ 16K Tokens). To push the limits of long-form reasoning, we conduct a second GRPO stage on DeepMath (He et al., 2025) using 16 rollouts per prompt with up to 16K output tokens. This stage is designed to maximize the model’s capacity for deep mathematical reasoning, enabling it to solve problems that require extended derivations, multiple nested logical steps, or structured proof-like outputs. In this stage, the model is trained for 600 GRPO steps.

Stage 3: GRPO on DeepScaleR (16 Rollouts @ 16K Tokens). In the final GRPO stage, we finetune the model on DeepScaleR (Luo et al., 2025), which includes original Olympiad math problems (e.g., AIME and AMC). Similar to Stage 2, this training uses 16 rollouts and a 16K token limit. We run 740 GRPO steps in this phase to improve performance on competition-style reasoning tasks.

6 Evaluation

6.1 Base Model

We evaluate the pre-trained base models on ARC-Challenge (ARC-C) (Clark et al., 2018), ARC-Easy (ARC-E) (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (HS) (Zellers et al., 2019), PiQA (Bisk et al., 2019), SciQ (Welbl et al., 2017), WinoGrande (WG) (Sakaguchi et al., 2019), OpenBookQA (OBQA) (Mihaylov et al., 2018), BBH (Suzgun et al., 2022), MMLU (Hendrycks et al., 2021a), and GSM8k (Cobbe et al., 2021a). All the benchmarks use a zero-shot evaluation setting, except BBH, MMLU, and GSM8k, which are evaluated using 3-shot, 5-shot, and 8-shot prompting, respectively.

As shown in Table 4, both Instella-3B-Stage1 and Instella-3B models outperform all the other fully open models over all the benchmarks individually (except PIQA). Our final pre-trained checkpoint Instella-3B outperforms the prior top performant fully open pre-trained models by a lead of 8.1% on average, with

Table 6: **Instruction-tuned model performance.**

Models	MMLU	TQA	BBH	GPQA	GSM8K	MATH	IFEval	AE 2	MT	Avg.
<i>Open Weight Models</i>										
Gemma-2-2B-Instruct	58.4	55.8	43.0	25.2	53.5	22.5	55.6	29.4	8.1	39.0
Llama-3.2-3B-Instruct	61.5	50.2	61.5	29.7	77.0	46.0	75.4	19.3	7.1	47.5
Qwen-2.5-3B-Instruct	66.9	57.2	57.3	28.1	76.0	60.4	62.5	22.1	8.0	48.7
<i>Fully Open Models</i>										
StableLM-zephyr-3B	45.1	47.9	39.3	25.7	58.4	10.4	34.2	7.5	6.0	30.5
OpenELM-3B-Instruct	27.4	38.1	24.2	18.1	1.6	0.4	16.1	0.2	1.0	14.1
Instella-3B-SFT	58.8	52.5	46.0	28.1	71.7	40.5	66.2	7.6	7.1	42.1
Instella-3B-Instruct	58.9	55.5	46.8	30.1	73.9	42.5	71.4	17.6	7.2	44.9

significant improvements in ARC Challenge (+8%), ARC Easy (+3.5%), Winnograde (+4.7%), OpenBookQA (+3.9%), MMLU (+13.1%) and GSM8K (+49%).

Second stage pre-training elevates the overall average performance relative to stage-1 by 5.3%, substantially narrowing the performance gap between Instella-3B model and the prior open weight models, and outperforming Llama-3.2-3B by 4.1% on average (+5.7% ARC-Challenge, +5.6% ARC-Easy, and +29.7% GSM8k), Gemma-2-2B by 7.3% on average (+13.4% ARC-Challenge, +11.2% ARC-Easy, +4.5% HellaSwag, +7.6% OpenBookQA, +5.0% MMLU, and +32.5% GSM8k), and is competitive with Qwen-2.5-3B on the majority of the benchmarks. As shown in Table 5, the Instella-3B checkpoint, obtained by merging the weights of three independently trained models with different random seeds during second stage pretraining, achieves an average performance of 66.6%, surpassing all individual seed runs.

The multi-stage pre-training with diverse and high-quality data mixture significantly enhances Instella-3B’s capabilities, establishing it as a competitive and open alternative in the landscape of comparable size language models.

6.2 Instruction-tuned Model

The instruction-tuned models are evaluated on MMLU (Hendrycks et al., 2021a), TruthfulQA (TQA) (Lin et al., 2022), BBH (Suzgun et al., 2022), GPQA (Rein et al., 2023), GSM8K (Cobbe et al., 2021a), Minerva Math (Lewkowycz et al., 2022) (MATH), IFEval (Zhou et al., 2023), Alpaca Eval V2 (AE2) (Dubois et al., 2025), and MT-Bench (MT) (Zheng et al., 2023). Here, GPQA, Minerva Math, IFEval, and Alpaca V2 use a zero-shot evaluation setting, whereas MMLU, TQA, BBH, and GSM8k use few-shot prompting using 5-shots, 6-shots, 3-shots, and 8-shots, respectively.

Instella-3B-Instruct model consistently outperforms other fully open models across all evaluated benchmarks with a significant average score lead of 14.37% with respect to the next top performing fully open instruction-tuned models (Table 6). With substantial margins across all the chat benchmarks (+13% MMLU, +7.57% TruthfulQA, +7.43% BBH, +4.46% GPQA, +37.15% IFEval, +10.08% Alpaca 2, and +1.2% MT-Bench).

Instella-3B-Instruct narrows the performance gap with leading open-weight models. Instella-3B-Instruct performs on par with or slightly surpasses existing state-of-the-art open weight instruction-tuned models such as Llama-3.2-3B-Instruct (+5.24% TruthfulQA, +0.45% GPQA, and +0.1% MT-Bench), and Qwen2.5-3B-Instruct (+2.01% GPQA and +8.87% IFEval), while significantly outperforming Gemma-2-2B-Instruct with an average score lead of +5.83% (+0.55% MMLU, +3.79% BBH, +4.91% GPQA, +20.47% GSM8k, +19.98% Minerva MATH, and +15.17% IFEval).

Overall, Instella-3B-Instruct excels in instruction following tasks and multi-turn QA tasks like TruthfulQA, GPQA, IFEval and MT-Bench, while being highly competitive compared to existing state-of-the-art open weight models on other knowledge recall and math benchmarks, while being trained on significantly fewer training tokens.

Table 7: **Long-context evaluation on the Helmet benchmark.** NQ: Natural Question. Inf: InfiniteBench. NarrQA: NarrativeQA. The NIAH-MV task and RAG task (NQ, TriviaQA, and HotpotQA) are evaluated at five context lengths: 8K, 16K, 32K, 64K, and 128K, and the number is reported by averaging across the five context lengths. The InfQA, InfMC, and NarrQA are evaluated at 128K context length.

Models	NQ	TriviaQA	HotpotQA	InfQA	InfMC	NarrQA	NIAH-MV	Avg.
<i>Open Weight Models</i>								
Llama-3.2-3B-Instruct	51.8	86.2	56.4	38.7	56.0	26.0	99.2	59.2
Phi-3.5-Mini-Instruct	41.2	78.6	48.6	24.0	55.0	27.7	87.0	51.7
Gemma-3-4B-it	47.2	76.8	45.2	21.0	49.0	20.7	74.0	47.7
Qwen-2.5-3B-Instruct	34.6	65.8	41.8	14.7	35.0	21.0	80.4	41.9
MiniCPM-2B-128k	28.4	61.6	30.8	3.7	22.0	3.3	46.6	28.1
<i>Fully Open Models</i>								
Instella-3B-Long-Instruct	43.6	73.0	51.6	30.7	54.0	32.3	84.0	52.7

6.3 Instella-Long

We evaluate the long-context performance on Helmet (Yen et al., 2024), a recent comprehensive long-context evaluation benchmark encompassing diverse categories. Helmet demonstrates more consistent alignment with human judgment. We evaluate three main tasks across seven datasets: multi-value needle-in-a-haystack (NIAH-MV), retrieval augmented generation (Natural Questions (Kwiatkowski et al., 2019), TriviaQA (Joshi et al., 2017), HotpotQA (Yang et al., 2018)), and long-document QA (InfiniteBench MC/QA (Zhang et al., 2024), NarrativeQA (Kočísky et al., 2018)). We use substring exact match (SubEM) for the RAG task, recall for NIAH-MV, and exact match for InfiniteBench MC. For InfiniteBench QA and NarrativeQA, which involve open-ended answers, we rely on gpt-4o-mini to evaluate model responses against the ground truth, following the prompt and metric provided by Helmet. As shown in Table 7, Instella-3B-Long-Instruct outperforms open weights models including Phi-3.5-mini-instruct (Abdin et al., 2024), Gemma-3-4B-it (Gemma Team, 2025), Qwen2.5-3B-Instruct (Yang et al., 2024), and MiniCPM-2B-128k (Hu et al., 2024) on most tasks of the Helmet benchmark. Since the context length of Qwen2.5-3B-Instruct is 32K, we also conduct a side-by-side comparison at 8K, 16K, and 32K context lengths, as shown in Table 8. Instella-3B-Long-Instruct outperforms Qwen2.5-3B-Instruct by 2.8% on average.

Table 8: **Comparison with Qwen2.5-3B-Instruct at 8K, 16K, 32K context lengths.**

Model	NIAH-MV			NQ			TriviaQA			HotpotQA			Avg.
	8K	16K	32K	8K	16K	32K	8K	16K	32K	8K	16K	32K	
Qwen2.5-3B-Instruct	95	94	95	48	42	39	77	78	74	51	50	48	65.9
Instella-3B-Long-Instruct	98	95	87	53	49	46	79	73	75	59	59	51	68.7

We also evaluate the short-context performance as shown in Table 9. We observe performance drops on some short-context benchmarks compared to Instella-3B-Instruct. Interestingly, TruthfulQA remains stable, Crows-Pairs shows a slight improvement, and the reduction in Toxigen (57.02 \rightarrow 42.34, lower is better) suggests improved toxicity avoidance, together indicating potential gains in responsible AI benchmarks. We hypothesize that these results reflect a trade-off between optimizing for longer context lengths and retaining short-context performance, which may be more pronounced at the 3B parameter scale compared to larger models.

Table 9: **Evaluation of Instella-Long on general benchmarks.**

Models	MMLU	IFEval	MT-Bench	TruthfulQA	Toxigen (↓)	Crows-Pair
Instella-3B-Instruct	58.9	71.4	7.2	55.5	57.0	58.9
Instella-3B-Long-Instruct	57.4	68.8	6.8	55.5	42.3	60.1

6.4 Instella-Math

Following the same evaluation settings as DeepScaleR-1.5B (Luo et al., 2025), we report Pass@1 accuracy over AIME 2024/25 (AIME), MATH500 (Hendrycks et al., 2021c), AMC (AMC), Mnerva MATH (Lewkowycz et al., 2022), OlympiadBench (He et al., 2024a), GSM8k (Cobbe et al., 2021b), and GPQA-Diamond (Rein et al., 2023). Table 10 reports the Pass@1 rate for the above benchmarks, calculated based on 16 responses per question. Instella-Math delivers competitive performance when compared to leading small-scale open-weight models such as Deepseek-R1-Distilled-Qwen-1.5B, Still-3-1.5B, DeepScaleR-1.5B and SmoLLM3-3B. In addition to achieving competitive average performance across all benchmarks, Instella-Math demonstrates the effectiveness of our RL training recipe—improving over its supervised finetuned variant (Instella-Math-SFT) by 10.81 points, compared to a 6.22-point improvement seen in DeepScaleR over its base model (Deepseek-R1-Distilled-Qwen-1.5B).

Table 10: Evaluation of Instella-Math on Reasoning Benchmarks

Models	AIME 2024	AIME 2025	MATH500	AMC	Minerva	OlympiadBench	GSM8K	GPQA-D	Avg.
Pass@1									
<i>Open-Weight Models</i>									
Qwen2.5-Math-1.5B	7.7	4.0	57.8	35.8	15.7	26.0	66.3	15.4	28.6
DeepSeek-R1-Distill-Qwen-1.5B	27.5	22.5	82.6	63.5	26.5	43.0	84.1	16.5	45.8
STILL-3-1.5B-preview	30.6	25.2	84.6	66.7	28.6	45.3	86.6	19.5	48.4
DeepScaleR-1.5B-Preview	40.6	30.8	87.4	73.2	30.1	49.9	87.3	16.5	52.0
<i>Fully-Open Models</i>									
OLMo-2-1124-7B-Instruct	1.3	0.2	32.6	12.3	10.3	8.5	80.9	11.1	19.6
SmoLLM3-3B	52.5	35.8	90.2	78.7	31.8	55.4	92.3	44.9	60.2
Instella-Math SFT	20.0	19.0	77.6	53.9	18.8	43.3	88.0	23.4	43.0
Instella-Math RL Stage 1	27.9	22.5	82.2	58.8	25.1	49.2	90.9	34.2	48.8
Instella-Math RL Stage 2	29.6	22.9	85.8	66.7	27.5	52.7	91.7	37.4	51.8
Instella-Math RL Stage 3	35.6	27.7	86.5	69.7	27.7	53.1	92.5	37.6	53.8
Pass@16									
<i>Open-Weight Models</i>									
Qwen2.5-Math-1.5B	36.7	20.0	87.6	71.1	48.5	53.8	96.0	71.7	60.7
DeepSeek-R1-Distill-Qwen-1.5B	73.3	46.7	95.0	89.2	54.4	63.9	97.0	46.5	70.7
STILL-3-1.5B-preview	70.0	46.7	95.8	89.2	56.6	65.2	96.7	45.5	70.7
DeepScaleR-1.5B-Preview	70.0	53.3	95.2	91.6	54.0	66.2	96.5	39.9	70.9
<i>Fully-Open Models</i>									
OLMo-2-1124-7B-Instruct	13.3	3.3	66.6	50.6	35.1	23.2	97.3	49.0	42.3
SmoLLM3-3B	76.7	77.3	96.6	94.0	54.4	72.4	98.1	90.9	82.1
Instella-Math SFT	50.0	40.0	94.8	89.2	44.9	64.0	97.7	83.8	70.6
Instella-Math RL Stage 1	53.3	43.3	94.6	88.0	51.5	68.6	97.6	90.9	73.5
Instella-Math RL Stage 2	46.7	43.3	95.6	89.2	51.1	68.3	97.7	89.4	72.7
Instella-Math RL Stage 3	63.3	50.0	95.8	86.8	50.4	68.2	97.4	88.9	75.1

Table 11: Evaluation of Instella-Math on TTT-Bench

Models	oTTT	dTTT	cTTT	sTTT	Avg.
<i>Open-weight models</i>					
Qwen2.5-Math-1.5B	12.5	10.0	18.9	7.5	12.2
DeepSeek-R1-Distill-Qwen-1.5B	22.9	10.1	18.2	3.5	13.7
STILL-3-1.5B-preview	24.5	12.3	19.8	3.2	14.9
DeepScaleR-1.5B-Preview	23.0	16.5	23.0	8.2	17.7
<i>Fully-open models</i>					
SmoLLM3-3B	51.2	40.1	41.3	42.3	43.7
Instella-Math RL Stage 1	56.3	31.4	39.7	41.9	42.3
Instella-Math RL Stage 2	66.2	37.3	39.2	44.5	46.8
Instella-Math RL Stage 3	70.3	39.6	40.3	49.0	49.8

Additionally, we test Instella-Math on TTT-Bench (Mishra et al., 2025), a new benchmark targeting strategic, spatial, and logical reasoning. Remarkably, without any exposure to TTT-Bench-style or similar strategic

gaming data during any stage of training, Instella-Math achieves the best performance among all evaluated models (as shown in Table 11).

More importantly, like OLMo2 and SmoLLM-3B, Instella-Math is a fully-open language model, with fully-open training data for the base model (Instella-3B), reasoning SFT, and reinforcement learning stages. In contrast, many competing models are only open-weight releases; their base model training (e.g., Qwen-1.5B) and reasoning distillation processes (e.g., DeepSeek-R1) remain closed.

7 Conclusion

We present Instella, a family of fully open three billion parameter language models that are trained entirely on openly available data and codebase. The Instella model family consists of a strong base pre-trained model, a supervised finetuned instruct model, an 128k token context length long-context model, and a reasoning-centric model. Powered by AMD Instinct™ MI300X GPUs, Instella models attain state-of-the-art performance among fully open models of similar scale and remains competitive with leading open-weight systems despite using notably fewer pre-training tokens. Instella-Long demonstrates strong long-context capabilities, and Instella-Math delivers impressive gains on mathematical and strategic reasoning benchmarks. Alongside model weights, we release the training code, data recipes, and evaluation protocols to support complete reproducibility and transparent benchmarking to foster open-source innovation. Instella models offers a transparent, performant, and extensible foundation for research and application, supporting the community in building more capable and reproducible language models.

References

- Marah Abidin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*, 2024.
- AIME. Aime problems and solutions, 2025. URL https://artofproblemsolving.com/wiki/index.php/American_Invitational_Mathematics_Examination.
- Alon Albalak, Duy Phung, Nathan Lile, Rafael Rafailov, Kanishk Gandhi, Louis Castricato, Anikait Singh, Chase Blagden, Violet Xiang, Dakota Mahan, and Nick Haber. Big-math: A large-scale, high-quality math dataset for reinforcement learning in language models, 2025. URL <https://arxiv.org/abs/2502.17387>.
- Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarán, Vaibhav Srivastav, Joshua Lochner, Caleb Fahlgren, Xuan-Son Nguyen, Clémentine Fourrier, Ben Burtenshaw, Hugo Larcher, Haojun Zhao, Cyril Zakka, Mathieu Morlon, Colin Raffel, Leandro von Werra, and Thomas Wolf. SmolLM2: When smol goes big – data-centric training of a small language model, 2025. URL <https://arxiv.org/abs/2502.02737>.
- AMC. American mathematics contest 12 (amc 12), 2022. URL https://artofproblemsolving.com/wiki/index.php/AMC_12.
- Anthropic. System card: Claude opus 4 & claude sonnet 4. Technical report, Anthropic, AI, 2025. URL <https://www-cdn.anthropic.com/4263b940cabb546aa0e3283f35b686f4f3b2ff47.pdf>.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Loubna Ben Allal, Anton Lozhkov, Guilherme Penedo, Thomas Wolf, and Leandro von Werra. SmolLM-corpus, July 2024. URL <https://huggingface.co/datasets/HuggingFaceTB/smolLM-corpus>.
- Steven Bird and Edward Loper. NLTK: The natural language toolkit. In *Proceedings of the ACL Interactive Poster and Demonstration Sessions*, pp. 214–217, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL <https://aclanthology.org/P04-3031/>.

-
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language, 2019. URL <https://arxiv.org/abs/1911.11641>.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2924–2936, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1300. URL <https://aclanthology.org/N19-1300/>.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018. URL <https://arxiv.org/abs/1803.05457>.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021a. URL <https://arxiv.org/abs/2110.14168>.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021b.
- Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu Yu, Qixin Xu, Weize Chen, et al. Process reinforcement through implicit rewards. *arXiv preprint arXiv:2502.01456*, 2025.
- Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691*, 2023.
- Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. In *ICLR*, 2024.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaoqun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qishi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhua Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shutong Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yuxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, Rodolphe Jenatton, Lucas Beyer, Michael Tschannen, Anurag Gritsenko, Joan Puigcerver, Matthias Minderer, Filip Pavetic, Francesco Locatello, Thomas Kipf, Sylvain Gelly, Andrew Brock, Alec Radford, Mario Lucic, and Neil Houlsby. Scaling vision transformers to 22 billion parameters. *arXiv preprint arXiv:2302.05442*, 2023.

-
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations. In *EMNLP*, 2023.
- Yiran Ding, Li Lyna Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan Yang, and Mao Yang. Longrope: Extending llm context window beyond 2 million tokens. *arXiv preprint arXiv:2402.13753*, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv-2407, 2024.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. Length-controlled alpaca-eval: A simple way to debias automatic evaluators, 2025. URL <https://arxiv.org/abs/2404.04475>.
- emozilla. Dynamically scaled rope further increases performance of long context llama with zero fine-tuning, 2023. URL https://www.reddit.com/r/LocalLLaMA/comments/14mrgpr/dynamically_scaled_rope_further_increases/.
- Tianyu Gao, Alexander Wettig, Howard Yen, and Danqi Chen. How to train long-context language models (effectively). *arXiv preprint arXiv:2410.02660*, 2024.
- Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. Technical report, Google, 2024. URL <https://storage.googleapis.com/deepmind-media/gemini/gemini-v1.5-report.pdf>.
- Gemma Team. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025.
- Gradient Team. Scaling rotational embeddings for long-context language models, 2024. URL <https://www.gradient.ai/blog/scaling-rotational-embeddings-for-long-context-language-models>.
- Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, William Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah Smith, and Hannaneh Hajishirzi. OLMo: Accelerating the science of language models. In *ACL*, 2024.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems, 2024a. URL <https://arxiv.org/abs/2402.14008>.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint arXiv:2402.14008*, 2024b.
- Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning. 2025. URL <https://arxiv.org/abs/2504.11456>.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding, 2021a. URL <https://arxiv.org/abs/2009.03300>.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *ICLR*, 2021b.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021c. URL <https://arxiv.org/abs/2103.03874>.

-
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021d.
- Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, et al. Minicpm: Unveiling the potential of small language models with scalable training strategies. *arXiv preprint arXiv:2404.06395*, 2024.
- Sam Ade Jacobs, Masahiro Tanaka, Chengming Zhang, Minjia Zhang, Shuaiwen Leon Song, Samyam Rajbhandari, and Yuxiong He. Deepspeed ulysses: System optimizations for enabling training of extreme long sequence transformer models. *arXiv preprint arXiv:2309.14509*, 2023.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
- Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. The narrativeqa reading comprehension challenge. *Transactions of the Association for Computational Linguistics*, 6:317–328, 2018.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in neural information processing systems*, 35:3843–3857, 2022.
- Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik Bansal, Etash Kumar Guha, Sedrick Keh, Kushal Arora, Saurabh Garg, Rui Xin, Niklas Muennighoff, Reinhard Heckel, Jean Mercat, Mayee F Chen, Suchin Gururangan, Mitchell Wortsman, Alon Albalak, Yonatan Bitton, Marianna Nezhurina, Amro Kamal Mohamed Abbas, Cheng-Yu Hsieh, Dhruva Ghosh, Joshua P Gardner, Maciej Kilian, Hanlin Zhang, Rulin Shao, Sarah M Pratt, Sunny Sanyal, Gabriel Ilharco, Giannis Daras, Kalyani Marathe, Aaron Gokaslan, Jieyu Zhang, Khyathi Chandu, Thao Nguyen, Igor Vasiljevic, Sham M. Kakade, Shuran Song, Sujay Sanghavi, Fartash Faghri, Sewoong Oh, Luke Zettlemoyer, Kyle Lo, Alaaeldin El-Nouby, Hadi Pouransari, Alexander T Toshev, Stephanie Wang, Dirk Groeneveld, Luca Soldaini, Pang Wei Koh, Jenia Jitsev, Thomas Kollar, Alex Dimakis, Yair Carmon, Achal Dave, Ludwig Schmidt, and Vaishaal Shankar. Datacomp-LM: In search of the next generation of training sets for language models. In *NeurIPS Datasets and Benchmarks Track*, 2024.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods, 2022. URL <https://arxiv.org/abs/2109.07958>.
- Hao Liu, Matei Zaharia, and Pieter Abbeel. Ring attention with blockwise transformers for near-infinite context. *arXiv preprint arXiv:2310.01889*, 2023.
- Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl. <https://pretty-radio-b75.notion.site/DeepScaler-Surpassing-O1-Preview-with-a-1-5B-Model-by-Scaling-RL-19681902c1468005bed8ca303013a4e2>, 2025. Notion Blog.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering, 2018. URL <https://arxiv.org/abs/1809.02789>.
- Prakamya Mishra, Jiang Liu, Jialian Wu, Xiaodong Yu, Zicheng Liu, and Emad Barsoum. Ttt-bench: A benchmark for evaluating reasoning ability with simple and novel tic-tac-toe-style games, 2025. URL <https://arxiv.org/abs/2506.10209>.

-
- Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Evan Pete Walsh, Oyvind Tafjord, Nathan Lambert, Yuling Gu, Shane Arora, Akshita Bhagia, Dustin Schwenk, David Wadden, Alexander Wettig, Binyuan Hui, Tim Dettmers, Douwe Kiela, Ali Farhadi, Noah A. Smith, Pang Wei Koh, Amanpreet Singh, and Hannaneh Hajishirzi. OLMoE: Open mixture-of-experts language models. In *ICLR*, 2025.
- Muzammal Naseer, Kanchana Ranasinghe, Salman Khan, Munawar Hayat, Fahad Khan, and Ming-Hsuan Yang. Intriguing properties of vision transformers. In *NeurIPS*, 2021.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Michal Guerquin, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2 olmo 2 furious, 2024. URL <https://arxiv.org/abs/2501.00656>.
- OpenAI. GPT4 technical report. *CoRR*, abs/2303.08774, 2023.
- Yiwei Qin, Xuefeng Li, Haoyang Zou, Yixiu Liu, Shijie Xia, Zhen Huang, Yixin Ye, Weizhe Yuan, Zhengzhong Liu, Yuanzhi Li, and Pengfei Liu. O1 replication journey: A strategic progress report – part 1. <https://github.com/GAIR-NLP/O1-Journey>, 2024.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. In *NeurIPS*, 2023.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 2023. URL <https://arxiv.org/abs/2311.12022>.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale, 2019. URL <https://arxiv.org/abs/1907.10641>.
- David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. Analysing mathematical reasoning abilities of neural models. In *ICLR*, 2019.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Jha, Sachin Kumar, Li Lucy, Xinxu Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Evan Walsh, Luke Zettlemoyer, Noah Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. Dolma: an open corpus of three trillion tokens for language model pretraining research. In *ACL*, 2024.
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V. Le, Ed H. Chi, Denny Zhou, and Jason Wei. Challenging big-bench tasks and whether chain-of-thought can solve them, 2022. URL <https://arxiv.org/abs/2210.09261>.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, and Jason Wei. Challenging BIG-bench tasks and whether chain-of-thought can solve them. In *ACL Findings*, 2023.

-
- Sho Takase, Shun Kiyono, Sosuke Kobayashi, and Jun Suzuki. Spike no more: Stabilizing the pre-training of large language models. In *COLM*, 2025.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- Teknium. Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants, 2023. URL <https://huggingface.co/datasets/teknium/OpenHermes-2.5>.
- Shubham Toshniwal, Wei Du, Ivan Moshkov, Branislav Kisacarin, Alexan Ayrapetyan, and Igor Gitman. Openmathinstruct-2: Accelerating ai for math with massive open-source instruction data. *arXiv preprint arXiv:2410.01560*, 2024.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions, 2017. URL <https://arxiv.org/abs/1707.06209>.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv:2412.15115*, 2024.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025a.
- An Yang, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoyan Huang, Jiandong Jiang, Jianhong Tu, Jianwei Zhang, Jingren Zhou, Junyang Lin, Kai Dang, Kexin Yang, Le Yu, Mei Li, Minmin Sun, Qin Zhu, Rui Men, Tao He, Weijia Xu, Wenbiao Yin, Wenyuan Yu, Xiafei Qiu, Xingzhang Ren, Xinlong Yang, Yong Li, Zhiying Xu, and Zipeng Zhang. Qwen2.5-1m technical report. *arXiv preprint arXiv:2501.15383*, 2025b.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.
- Howard Yen, Tianyu Gao, Minmin Hou, Ke Ding, Daniel Fleischer, Peter Izsak, Moshe Wasserblat, and Danqi Chen. Helmet: How to evaluate long-context language models effectively and thoroughly. *arXiv preprint arXiv:2410.02694*, 2024.
- Xiaodong Yu, Ben Zhou, Hao Cheng, and Dan Roth. Reasonagain: Using extractable symbolic programs to evaluate mathematical reasoning. *arXiv preprint arXiv:2410.19056*, 2024.
- Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhui Chen. MAMmoTH2: Scaling Instructions from the Web. *NeurIPS*, 2024.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1472. URL <https://aclanthology.org/P19-1472/>.
- Biao Zhang and Rico Sennrich. Root Mean Square Layer Normalization. In *NeurIPS*, 2019.
- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Hao, Xu Han, Zhen Thai, Shuo Wang, Zhiyuan Liu, et al. ∞ bench: Extending long context evaluation beyond 100k tokens. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15262–15277, 2024.

-
- Han Zhao, Haotian Wang, Yiping Peng, Sitong Zhao, Xiaoyu Tian, Shuaiting Chen, Yunjie Ji, and Xiangang Li. 1.4 million open-source distilled reasoning dataset to empower large language model training, 2025. URL <https://arxiv.org/abs/2503.19633>.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL <https://arxiv.org/abs/2306.05685>.
- Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and Xiang Yue. OpenCodeInterpreter: Integrating code generation with execution and refinement. In *ACL Findings*, 2024.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023. URL <https://arxiv.org/abs/2311.07911>.