

Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone

Microsoft

Abstract

We introduce **phi-3-mini**, a 3.8 billion parameter language model trained on 3.3 trillion tokens, whose overall performance, as measured by both academic benchmarks and internal testing, rivals that of models such as Mixtral 8x7B and GPT-3.5 (e.g., **phi-3-mini** achieves 69% on MMLU and 8.38 on MT-bench), despite being small enough to be deployed on a phone. The innovation lies entirely in our dataset for training, a scaled-up version of the one used for **phi-2**, composed of heavily filtered publicly available web data and synthetic data. The model is also further aligned for robustness, safety, and chat format. We also provide some initial parameter-scaling results with a 7B and 14B models trained for 4.8T tokens, called **phi-3-small** and **phi-3-medium**, both significantly more capable than **phi-3-mini** (e.g., respectively 75% and 78% on MMLU, and 8.7 and 8.9 on MT-bench). Moreover, we also introduce **phi-3-vision**, a 4.2 billion parameter model based on **phi-3-mini** with strong reasoning capabilities for image and text prompts.

1 Introduction

The striking progress of AI in the last few years can be largely attributed to major efforts throughout the world towards *scaling-up* to ever-larger models and datasets. Large Language Models (LLMs) have steadily increased in size from a mere billion parameters just five years ago (GPT-2 had 1.5 billion parameters [RWC⁺19]) to trillion parameters today. The impetus for this effort originates in the seemingly predictable improvement one obtains by training large models, the so-called *scaling laws* [KMH⁺20, HBM⁺22, MRB⁺23]. However these laws assume a “fixed” data source. This assumption is now significantly disrupted by the existence of frontier LLMs themselves, which allow us to interact with data in novel ways. In our previous works on the phi models [GZA⁺23, LBE⁺23, JBA⁺23] it was shown that a combination of LLM-based filtering of publicly available web data, and LLM-created synthetic data, enable performance in smaller language models that were typically seen only in much larger models. For example our previous model trained on this data recipe, **phi-2** (2.7B parameters), matched the performance of models 25 times larger trained on regular data. In this report we present a new model, **phi-3-mini** (3.8B parameters), trained for 3.3T tokens on larger and more advanced versions of the datasets used in **phi-2**. With its small size, **phi-3-mini** can easily be inferenced locally on a modern phone (see Figure 2), yet it achieves a quality that seems on-par with models such as Mixtral 8x7B [JSR⁺24] and GPT-3.5.

arXiv:2404.14219v3 [cs.CL] 23 May 2024

User: Explain why it is surprising that one can build a language model small enough to fit on a phone, yet almost as powerful as ChatGPT. Just use one funny sentence.

phi-3-mini: It's like fitting a supercomputer in a flip phone, but instead of breaking the phone, it just breaks the internet with its tiny, yet mighty, linguistic prowess!

User: Okay now more serious answer, and note that this was achieved solely by changing the training data.

phi-3-mini: The development of a compact language model that rivals the capabilities of ChatGPT, while fitting on a phone, is a testament to the power of data-driven machine learning. By meticulously curating and optimizing the training dataset, researchers can significantly reduce the model's size without compromising its performance. [...]

2 Technical Specifications

The **phi-3-mini** model is a transformer decoder architecture [VSP⁺17], with default context length 4K. We also introduce a long context version via LongRope [DZZ⁺24a] that extends the context length to 128K, called **phi-3-mini-128K**.

To best benefit the open source community, **phi-3-mini** is built upon a similar block structure as Llama-2 [TLI⁺23] and uses the same tokenizer with vocabulary size of 32064¹. This means that all packages developed for Llama-2 family of models can be directly adapted to **phi-3-mini**. The model uses 3072 hidden dimension, 32 heads and 32 layers. We trained using bfloat16 for a total of 3.3T tokens. The model is already chat-finetuned, and the chat template is as follows:

```
<|user|>\n Question <|end|>\n <|assistant|>
```

The **phi-3-small** model (7B parameters) leverages the tiktoken tokenizer (for better multilingual tokenization) with a vocabulary size of 100352² and has default context length 8192. It follows the standard decoder architecture of a 7B model class, having 32 heads, 32 layers and a hidden size of 4096. We switched from GELU activation to GEGLU and used Maximal Update Parametrization (muP) [?] to tune hyperparameters on a small proxy model and transfer them to the target 7B model. Those helped ensure better performance and training stability. Also, the model leverages a grouped-query attention, with 4 queries sharing 1 key. To optimize the training and inference speed, we design a novel blocksparse attention module. For each attention head, the blocksparse attention enforces different sparsity patterns over KV cache. This ensures that all tokens are attended to on different heads for the given choice of sparsity. As illustrated in Figure 1, the context is then efficiently divided and conquered among attention heads, with significant KV cache reduction. To achieve actual deployment speed-up from the blocksparse design, we implemented highly efficient, yet flexible kernels for both training and inference. For training, we build a triton kernel based on Flash Attention [DFE⁺22]. For inference, we implemented a kernel for the prefilling phase and extended the paged attention kernel in vLLM for the decoding phase [KLZ⁺23]. Lastly, in **phi-3-small** architecture, we alternate dense attention layers and blocksparse attention layers to optimize KV cache savings while maintaining long context retrieval performance. An additional 10% multilingual data was also used for this model.

Highly capable language model running locally on a cell-phone. Thanks to its small size, **phi-3-mini** can be quantized to 4-bits so that it only occupies ≈ 1.8 GB of memory. We tested the quantized model by deploying **phi-3-mini** on iPhone 14 with A16 Bionic chip running natively on-device and fully offline achieving more than 12 tokens per second.

¹We remove BoS tokens and add some additional tokens for chat template.

²We remove unused tokens from the vocabulary.

	Block1	Block2	Block3	Block4	Block5	Block6	Block7	Block8
Head1	Orange	Gray	Gray	Orange	Gray	Gray	Blue	Blue
Head2	Gray	Orange	Gray	Gray	Orange	Gray	Blue	Blue
Head3	Gray	Gray	Orange	Gray	Gray	Orange	Blue	Blue
Head4	Orange	Gray	Gray	Orange	Gray	Gray	Blue	Blue

Figure 1: Toy illustration of the blocksparse attention in phi-3-small with 2 local blocks and vertical stride of 3. The table shows the Keys/values a query token in block 8 attended to. Blue=local blocks, orange=remote/vertical blocks, gray=blocks skipped.

Training Methodology. We follow the sequence of works initiated in “Textbooks Are All You Need” [GZA⁺23], which utilize high quality training data to improve the performance of small language models and deviate from the standard *scaling-laws*. In this work we show that such method allows to reach the level of highly capable models such as GPT-3.5 or Mixtral with only 3.8B total parameters (while Mixtral has 45B total parameters for example). Our training data consists of heavily filtered publicly available web data (according to the “educational level”) from various open internet sources, as well as synthetic LLM-generated data. Pre-training is performed in two disjoint and sequential phases; phase-1 comprises mostly of web sources aimed at teaching the model general knowledge and language understanding. Phase-2 merges even more heavily filtered webdata (a subset used in Phase-1) with some synthetic data that teach the model logical reasoning and various niche skills.

Data Optimal Regime. Unlike prior works that train language models in either “compute optimal regime” [HBM⁺22] or “over-train regime”, we mainly focus on the quality of data for a *given scale*.³ We try to calibrate the training data to be closer to the “data optimal” regime for small models. In particular, we filter the publicly available web data to contain the correct level of “knowledge” and keep more web pages that could potentially improve the “reasoning ability” for the model. As an example, the result of a game in premier league in a particular day might be good training data for frontier models, but we need to remove such information to leave more model capacity for “reasoning” for the mini size models. We compare our approach with Llama-2 in Figure 3.

To test our data on larger size of models, we also trained **phi-3-medium**, a model with 14B parameters using the same tokenizer and architecture of **phi-3-mini**, and trained on the same data for slightly more epochs (4.8T tokens total as for **phi-3-small**). The model has 40 heads and 40 layers, with embedding dimension 5120. We observe that some benchmarks improve much less from 7B to 14B than they do from 3.8B to 7B, perhaps indicating that our data mixture needs further work to be in the “data optimal regime” for 14B parameters model.

Post-training. Post-training of **phi-3-mini** went through two stages, including supervised finetuning (SFT) and direct preference optimization (DPO). SFT leverages highly curated high-quality data across diverse domains, e.g., math, coding, reasoning, conversation, model identity, and safety. The SFT data mix starts with using English-only examples. DPO data covers chat format data, reasoning, and responsible AI (RAI) efforts. We use DPO to steer the model away from unwanted behavior, by using those outputs as “rejected” responses. Besides improvement in math, coding, reasoning, robustness, and safety, post-training transforms a language model to an AI assistant that users can efficiently and safely interact with.

³Just like for “compute optimal regime”, we use the term “optimal” in an aspirational sense for “data optimal regime”. We are not implying that we actually found the provably “optimal” data mixture for a given scale.

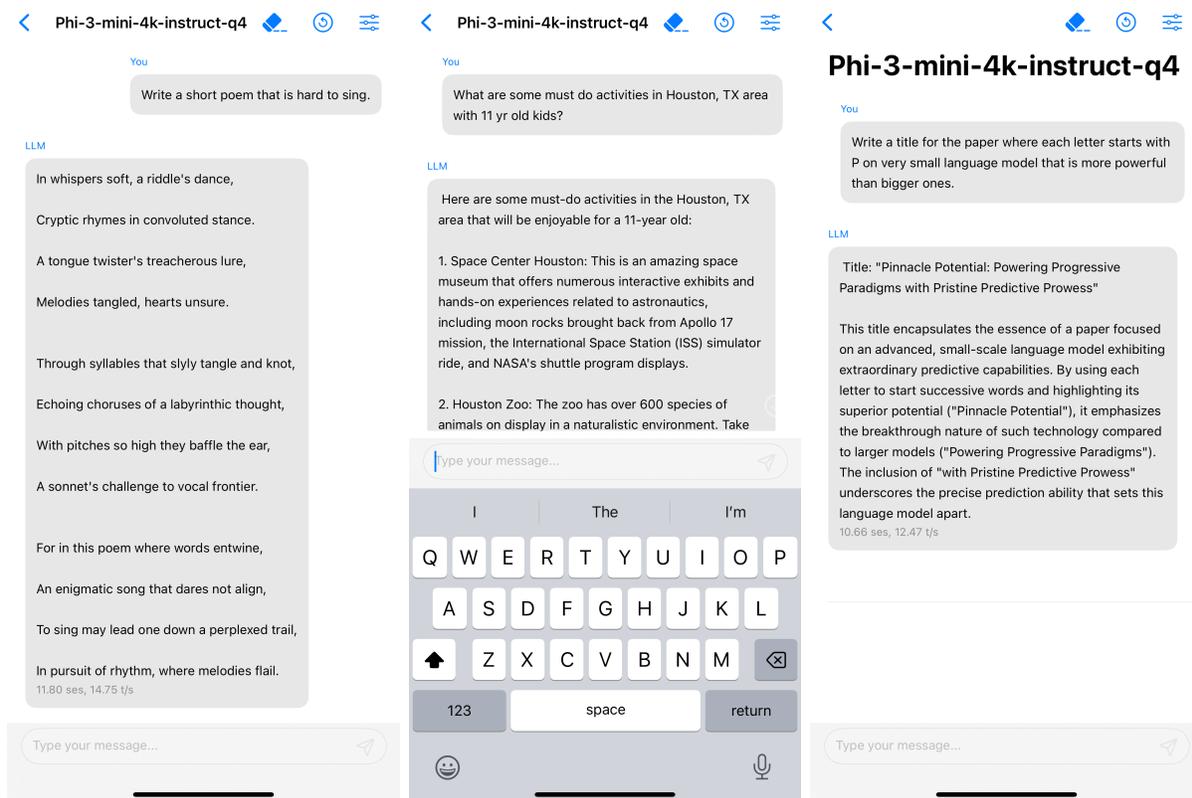


Figure 2: 4-bit quantized **phi-3-mini** running natively on an iPhone with A16 Bionic chip, generating over 12 tokens per second.

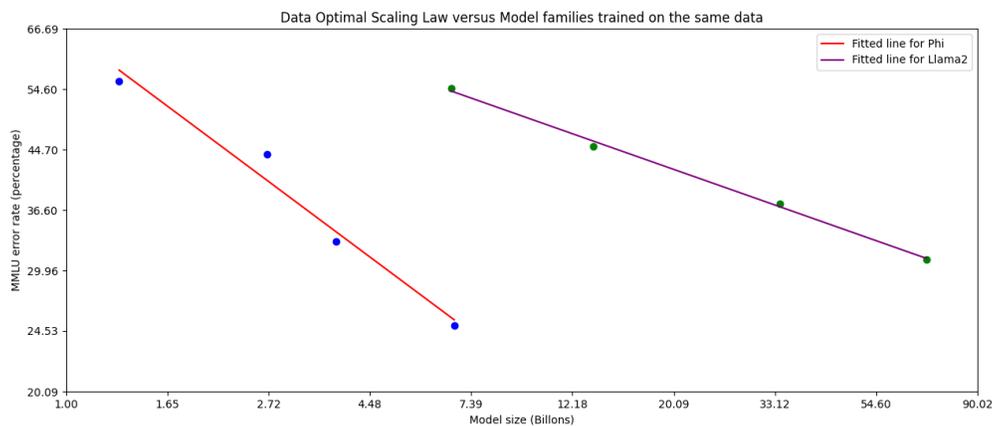


Figure 3: Scaling law close to the “Data Optimal Regime” (from left to right: phi-1.5, phi-2, phi-3-mini, phi-3-small) versus Llama-2 family of models (7B, 13B, 34B, 70B) that were trained on the same fixed data. We plot the log of MMLU error versus the log of model size.

As part of the post-training process, we developed a long context version of **phi-3-mini** with context length limit enlarged to 128K instead of 4K. Across the board, the 128K model quality is on par with the 4K length version, while being able to handle long context tasks. Long context extension has been done in two stages, including long context mid-training and long-short mixed post-training with both SFT and DPO.

3 Academic benchmarks

On the next page we report the results for **phi-3-mini** on standard open-source benchmarks measuring the model’s reasoning ability (both common sense reasoning and logical reasoning). We compare to phi-2 [JBA⁺23], Mistral-7b-v0.1 [JSM⁺23], Mixtral-8x7b [JSR⁺24], Gemma 7B [TMH⁺24], Llama-3-instruct-8b [AI], and GPT-3.5. All the reported numbers are produced with the exact same pipeline to ensure that the numbers are comparable. These numbers might differ from other published numbers due to slightly different choices in the evaluation. As is now standard, we use few-shot prompts to evaluate the models, at temperature 0. The prompts and number of shots are part of a Microsoft internal tool to evaluate language models, and in particular we did no optimization to the pipeline for the **phi-3** models.⁴ The number of k -shot examples is listed per-benchmark. An example of a 2-shot prompt is described in Appendix A.

⁴For example, we found that using `##` before the Question can lead to a noticeable improvement to **phi-3-mini**’s results across many benchmarks, but we did not do such changes in the prompts.

	Phi-3-mini 3.8b	Phi-3-small 7b	Phi-3-medium 14b	Phi-2 2.7b	Mistral 7b	Gemma 7b	Llama-3-In 8b	Mixtral 8x7b	GPT-3.5 version 1106
MMLU (5-Shot) [HBK ⁺ 21]	68.8	75.7	78.0	56.3	61.7	63.6	66.5	70.5	71.4
HellaSwag (5-Shot) [ZHB ⁺ 19]	76.7	77.0	82.4	53.6	58.5	49.8	71.1	70.4	78.8
ANLI (7-Shot) [NWD ⁺ 20]	52.8	58.1	55.8	42.5	47.1	48.7	57.3	55.2	58.1
GSM-8K (8-Shot; CoT) [CKB ⁺ 21]	82.5	89.6	91.0	61.1	46.4	59.8	77.4	64.7	78.1
MedQA (2-Shot) [JPO ⁺ 20]	53.8	65.4	69.9	40.9	50.0	49.6	60.5	62.2	63.4
AGIEval (0-Shot) [ZCG ⁺ 23]	37.5	45.1	50.2	29.8	35.1	42.1	42.0	45.2	48.4
TriviaQA (5-Shot) [JCWZ17]	64.0	58.1	73.9	45.2	75.2	72.3	67.7	82.2	85.8
Arc-C (10-Shot) [CCE ⁺ 18]	84.9	90.7	91.6	75.9	78.6	78.3	82.8	87.3	87.4
Arc-E (10-Shot) [CCE ⁺ 18]	94.6	97.0	97.7	88.5	90.6	91.4	93.4	95.6	96.3
PIQA (5-Shot) [BZGC19]	84.2	86.9	87.9	60.2	77.7	78.1	75.7	86.0	86.6
SociQA (5-Shot) [BZGC19]	76.6	79.2	80.2	68.3	74.6	65.5	73.9	75.9	68.3
BigBench-Hard (3-Shot; CoT) [SRR ⁺ 22, SSS ⁺ 22]	71.7	79.1	81.4	59.4	57.3	59.6	51.5	69.7	68.32
WinoGrande (5-Shot) [SLBBC19]	70.8	81.5	81.5	54.7	54.2	55.6	65.0	62.0	68.8
OpenBookQA (10-Shot) [MCKS18]	83.2	88.0	87.4	73.6	79.8	78.6	82.6	85.8	86.0
BoolQ (2-Shot) [CLC ⁺ 19]	77.2	84.8	86.5	–	72.2	66.0	80.9	77.6	79.1
CommonSenseQA (10-Shot) [THLB19]	80.2	80.0	82.8	69.3	72.6	76.2	79.0	78.1	79.6
TruthfulQA (10-Shot; MC2) [LHE22]	65.0	70.2	75.1	–	53.0	52.1	63.2	60.1	85.8
HumanEval (0-Shot) [CTJ ⁺ 21]	58.5	61.0	62.2	59.0	28.0	34.1	60.4	37.8	62.2
MBPP (3-Shot) [AON ⁺ 21]	70.0	71.7	75.2	60.6	50.8	51.5	67.7	60.2	77.8
Average	71.2	75.7	78.5	–	61.2	61.7	69.4	69.8	74.3
GPQA (2-Shot; CoT) [RHS ⁺ 23]	32.8	34.3	–	–	–	–	–	–	29.0
MT Bench (2 round ave.) [ZCS ⁺ 23]	8.38	8.70	8.91	–	–	–	–	–	8.35

4 Safety

Phi-3-mini was developed in accordance with Microsoft’s responsible AI principles. The overall approach consisted of safety alignment in post-training, red-teaming, automated testing and evaluations across dozens of RAI harm categories. Helpfulness and harmlessness preference datasets [BJN⁺22, JLD⁺23] with modifications inspired by [BSA⁺24] and multiple in-house generated datasets were leveraged to address the RAI harm categories in safety post-training. An independent red team at Microsoft iteratively examined **phi-3-mini** to further identify areas of improvement during the post-training process. Based on their feedback, we curated additional datasets tailored to address their insights, thereby refining the post-training dataset. This process resulted in significant decrease of harmful response rates, as shown in Figure 4.

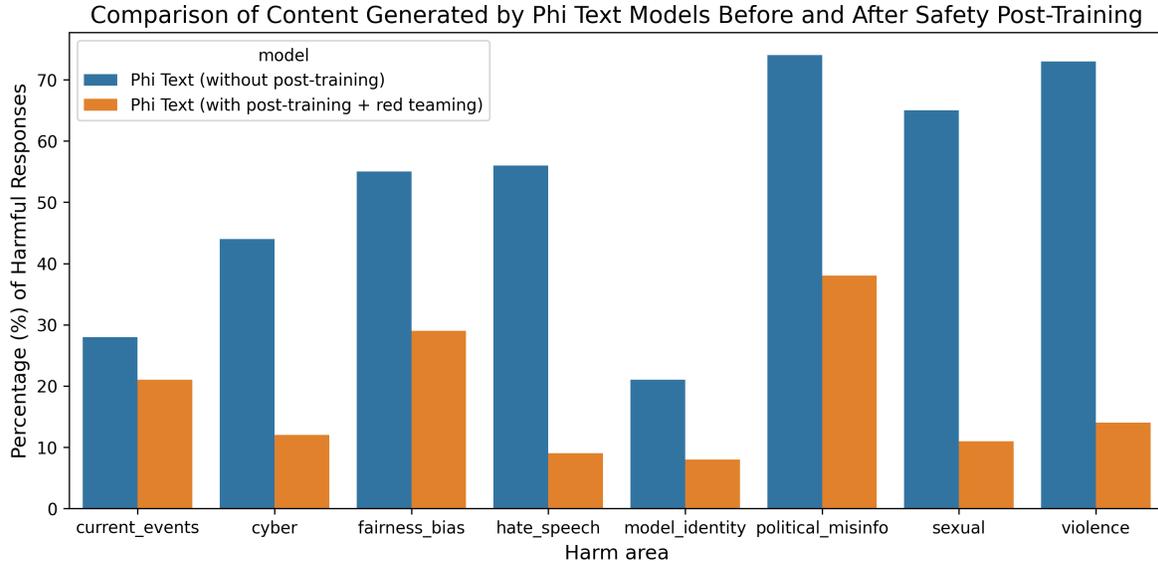


Figure 4: Comparison of harmful response percentages by Microsoft AI Red Team between **phi-3-mini** before and after the safety alignment. Note that the harmful response percentages in this chart are inflated numbers as the red team tried to induce **phi-3-mini** in an adversarial way to generate harmful responses through multi-turn conversations.

	Phi-3-mini 3.8b	Phi-3-small 7b	Phi-3-medium 14b	Phi-2 2.7b	Mistral 7b	Gemma 7b	Llama-3-In 8b
Ungroundedness	0.603	0.299	0.213	1.481	0.935	0.679	0.328
Third Party Harm (DR-1)	0.240	0.253	0.251	0.240	0.562	0.383	0.373
Harmful Content Continuation (DR-3)	0.007	0.003	0.010	0.029	0.026	0.013	0.013
Harmful Content Summarization (DR-3)	0.100	0.110	0.112	0.144	0.223	0.103	0.082
Jailbreak (DR-1)	0.123	0.107	0.111	0.150	0.156	0.114	0.130

Table 1: Comparison of Microsoft internal multi-turn conversation RAI benchmark results of **phi-3** models and other models. Note that a lower value indicates a better performance for all metrics in the table.

The safety alignment of **phi-3-small** and **phi-3-medium** was conducted by undergoing the same red-teaming process, utilizing identical datasets, and incorporating a slightly larger number of samples. Table 1 shows the results of in-house RAI benchmarks [MHJ⁺23] for **phi-3** models compared to phi-2 [JBA⁺23], Mistral-7b-v0.1 [JSM⁺23], Gemma 7b [TMH⁺24], and Llama-3-instruct-8b [AI]. This benchmark utilized GPT-4 to simulate multi-turn conversations in five different categories and to evaluate the model responses. Ungroundedness between 0 (fully grounded) and 4 (not grounded) measures if the information in a response is based on a given prompt. In other categories, responses were evaluated in terms of the severity of harmfulness from 0 (no harm) to 7 (extreme harm) and the defect rates (DR- x) were computed as the percentage of samples with the severity score being greater than or equal to x .

5 Weakness

In terms of LLM capabilities, while **phi-3-mini** model achieves similar level of language understanding and reasoning ability as much larger models, it is still fundamentally limited by its size for certain tasks.

Suggest an itinerary for a 3-day trip during the 2026 Winter Olympic games

Day 1: Arrival and Surrounding Neighborhoods

Morning:

- Arrive in Pyeongchang, South Korea, and check-in to your accommodations.
- Enjoy a light breakfast at your hotel.

Mid-morning:

- Take a leisurely walk around the town and explore the local shops and cafes.

Lunch:

- Try local specialties such as kimbab (Korean rice rolls) or mandu (Korean dumplings) at a nearby restaurant.

Afternoon:

- Visit Pyeongchang Olympic Cauldron Park, the site of the 2018 Winter Olympics, to take photos and admire the beautiful scenery.
- Rent a bike and explore the surrounding neighborhoods, such as Pyeongchang Village, where you can find souvenir shops, outdoor activities, and various dining options.

Evening:

- Dine at a local restaurant and enjoy traditional Korean cuisine such as samgyeopsal (grilled pork belly) or bulgogi (marinated beef).

Day 2: Olympic Sports and Cultural Activities

Morning:

- Visit the Alpensia Sports Park, where the majority of the 2026 Winter Olympic events will be held.
- See the ski slopes, ice rinks, and indoor arenas, and watch a glimpse of the opening and closing ceremonies.

Mid-morning:

- Participate in a guided tour to learn about the history of the 2026 Winter Olympics and the unique features of the Pyeongchang venues.

Search web

Ask anything

Suggest an itinerary for a 3-day trip during the 2026 Winter Olympic games

Web Search Completed

Day 1: Arrival and Exploring Milan

- Arrive in Milan on February 6, 2026, the day of the Opening Ceremony.
- Spend the day exploring Milan's city center, including Piazza del Duomo, Galleria Vittorio Emanuele II, and the Sforza Castle.
- Enjoy the vibrant Milanese nightlife, with options like the city's many trendy cafes and restaurants.
- Visit the new Fiera Milano expo center for a unique and inspiring experience.

Day 2: Cortina d'Ampezzo and Outdoor Activities

- Take the train from Milan to Cortina d'Ampezzo, arriving in the afternoon.
- Check into your hotel and spend the evening relaxing and enjoying the local cuisine, such as risotto and polenta, at traditional pizzerias or trattorias.
- On Day 2, take part in various outdoor activities, such as skiing, snowboarding, or tubing. Many of these activities are available at nearby resorts like Madonna di Campiglio, Kronplatz, or Passo San Pellegrino.
- In the evening, explore Cortina d'Ampezzo's charming streets, visit the historic St. Felix Castle, and enjoy the beautiful mountain scenery.

Day 3: Paralympic Games and Departure

- On Day 3, attend the Paralympic Games, held in Turin on March 6, 2026.
- Take the train from Cortina d'Ampezzo to Turin, arriving in time for the Paralympic events.
- Explore Turin's attractions, such as the Mole Antonelliana, the Egyptian Museum, and Piazza Vittorio Veneto.
- Return to Milan in the evening for your departure, concluding your 3-day itinerary during the Milano Cortina 2026 Winter Olympic Games.

Throughout your trip, make sure to book your accommodations, transportation, and tickets in advance, as these events are highly popular and may sell out quickly. Enjoy your time in Italy and the exciting Winter Olympic Games!

Sources: olympics.com olympics.com milanocortina2026.olympics.com milanwintertravel.com

Search web

Ask anything

Figure 5: Left: **phi-3-mini**'s completion without search. Right: **phi-3-mini**'s completion with search, using the default HuggingFace Chat-UI search ability. For reference, the 2026 Winter Olympic Games are scheduled to be held in Milano and Cortina in Italy, while the 2022 and 2018 Winter Olympic Games were held in Beijing, China and PyeongChang, Korea, respectively. Without the search results, the response is incorrect, while with the web search, not only does the response become accurate, but also gets more specific with suggestions.

The model simply does not have the capacity to store too much “factual knowledge”, which can be seen for example with low performance on TriviaQA. However, we believe such weakness can be resolved by augmentation with a search engine. We show an example using the HuggingFace default Chat-UI with **phi-3-mini** in Figure 5. Another weakness related to model’s capacity is that we mostly restricted the language to English. Exploring multilingual capabilities for Small Language Models is an important next step, with some initial promising results on **phi-3-small** by including more multilingual data.

Despite our diligent RAI efforts, as with most LLMs, there remains challenges around factual inaccuracies (or hallucinations), reproduction or amplification of biases, inappropriate content generation, and safety issues. The use of carefully curated training data, and targeted post-training, and improvements from red-teaming insights significantly mitigates these issues across all dimensions. However, there is significant work ahead to fully address these challenges, and downstream use of the models should be evaluated for the specific use cases and safety considerations for that context.



extract all displayed texts and keep original displayed format. For the left section, please use markdown table with column "coffee", "small", "medium" and "large". Instead, for the right section, please use JSON format.

Certainly! Below is the markdown table converted from the image:

Coffee	Small	Medium	Large
Espresso	1.30	1.60	2.00
Americana	1.40	1.75	2.10
Cappuccino	2.25	2.90	3.25
Latte	2.35	3.00	3.25
Mocha	2.35	3.00	3.50
Chai	2.25	2.40	3.25
Chai Bomb	2.75	3.25	3.75
House Cafe	1.50	1.75	2.00
Cafe au Lait	1.75	2.50	3.25

And here is the JSON representation of the teas offered:

```

{
  "Teas": {
    "CUP": {
      "GREEN": "$1.50",
      "BLACK": "$1.50",
      "WHITE": "$1.50",
      "OOLONG": "$1.50",
      "HERBAL": "$1.50"
    },
    "POT": {
      "GREEN": "$4.00",
      "BLACK": "$4.00",
      "WHITE": "$4.00",
      "OOLONG": "$4.00",
      "HERBAL": "$4.00"
    }
  },
  "**ADD .50 FOR HONEY**",
  "**ADD $1.00 FOR ICE**"
}

```

What is the price of a Cappuccino with large size?

- The price for a Cappuccino in large size is \$3.25.

How much does it cost to add ice to a tea?

- It costs an additional \$1.00 to add ice to any tea.

If someone wants to buy a pot of tea, how much would it cost?

- A pot of tea would cost \$4.00.

Figure 6: The demo case shows Phi-3-Vision’s capability in natural image understanding and reasoning.

6 Phi-3-Vision

6.1 Technical Specifications

Architecture The **Phi-3-Vision** (4.2B parameters) is a multimodal model designed to process an image and a textual prompt as inputs, and subsequently generate textual outputs. This model is composed of two primary components: an image encoder, *i.e.*, CLIP ViT-L/14 [RKH⁺21] and a transformer decoder, *i.e.*, phi-3-mini-128K-instruct. The visual tokens, once extracted by the image encoder, are then combined with text tokens in an interleaved way (no particular order for image and text tokens). To accommodate high-resolution images and various aspect ratios, a dynamic cropping strategy [DZZ⁺24b] is utilized to split the input image into a 2d array of blocks, where the tokens of the blocks are concatenated to represent the whole image.

Pre-training The **Phi-3-Vision** model undergoes a pre-training phase using a diverse dataset, which consists of a combination of interleaved image-text documents (*e.g.*, [LST⁺24]), image-text pairs from FLD-5B [XWX⁺24], synthetic data derived from Optical Character Recognition (OCR) of PDF files, datasets for chart/table comprehension, and text-only data. The objective of predicting the next token is employed specifically on text tokens, while any loss associated with image tokens is disregarded during this phase. The pre-training process involves a total of $0.5T$ tokens that encompass both visual and text elements. During the pre-training phase, the maximum image resolution is capped at 1344×1344 as the majority of the training images are smaller than this resolution.

Post-training. The **Phi-3-Vision** model contains two post-training stages: supervised finetuning (SFT) and direct preference optimization (DPO). For SFT, we leveraged text SFT dataset, public multimodal instruct tuning datasets along with large-scale multimodal instruct tuning datasets that we built ourselves, covering diverse domains and tasks such as general natural image understanding, chart/table/-diagram understanding/reasoning, PowerPoint understanding, and model safety. The multimodal SFT data has about a total of 15B tokens. For DPO we mainly use a text DPO dataset and a relatively smaller-scale multimodal DPO dataset. For these two stages, we jointly train multimodal tasks and text-only tasks so that the model can achieve multi-modal reasoning while maintaining language capabilities as much as possible.

6.2 Academic benchmarks

We report in Table 2 the evaluation results of Phi-3-Vision on nine open-source academic benchmarks. These benchmarks evaluate reasoning and perceptual capabilities on visual and text inputs and can be grouped in three categories: Science, Charts, and Generic knowledge. We compare Phi-3-Vision with the following baselines: MM1-3B-Chat [MGF⁺24], MM1-7B-Chat [MGF⁺24], Llava-1.6 Vicuna 7B [LLLL23], Llava-1.6 Llama3-8B [LLL⁺24], Qwen-VL-Chat [BBY⁺23], Claude 3 Haiku [Ant24], Gemini 1.0 Pro V [TAB⁺23], and GPT-4V-Turbo. Our performance quality assessment setup used the same evaluation pipeline for all the baselines to ensure a fair comparison, with the exception of MM1-3B-Chat. We just copied and pasted their published numbers since the model is not publicly available.

Our evaluation setup aimed to mimic scenarios where regular users interact with a multi-modal model, i.e., users who are not experts in prompt engineering or know special techniques that can improve performance. For this reason, we adopted the evaluation setting used in Llava-1.5 [LLLL23]. In this setup, the prompts include instructions to select a single letter corresponding to an answer from a list of given options, or answer with a single word or phrase. In our prompts, we did not use specific tokens for multiple-choice questions. Moreover, we did not scale or pre-process any image in our benchmarking system. We placed the images as the first item in the prompts, except on the MMMU dataset where the prompts interleave the images anywhere in the question or the answers. Lastly, our evaluation setup only considered a 0-shot format. Because of these evaluation parameters, our reported numbers can differ from the published numbers of the considered baselines.

6.3 Safety

To ensure the integration of **Phi-3-Vision** aligns with Microsoft’s Responsible AI (RAI) principles, we involved safety post-training in both Supervised Fine-Tuning (SFT) stage and Direct Preference Optimization (DPO) stage. In creating the safety training datasets, we utilized not only the text-only RAI datasets, but also a variety of in-house Multi-Modal (MM) RAI datasets that cover various harm categories identified in both public and internal MM RAI benchmarks. For the purpose of RAI evaluation, we performed a rigorous quantitative assessment on both public and internal benchmarks, this was done in conjunction with a human evaluation conducted by Microsoft’s internal red team.

In Table 3, we present the evaluation outcomes of Phi-3-Vision on three MM RAI benchmarks: one internal and two public benchmarks (specifically, RTVLM [LLY⁺24] and VGuard [ZBY⁺24]). We juxtapose these results with those of other open-source models such as Llava-1.5 [LLLL23], Llava-1.6 [LLL⁺24], Qwen-VL-Chat [BBY⁺23], and GPT4-V[Ope23]. The results clearly indicate that safety post-training notably enhances the RAI performance of Phi-3-Vision across all RAI benchmarks. In Figure 7, we further breakdown the performance across different RAI categories of the VGuard and Internal benchmarks, demonstrating that safety post-training can aid Phi-3-Vision in improving RAI performance in nearly all categories.

	Phi-3-Vision 4.2b	MM1-3B-Chat 3.6b [MGF ⁺ 24]	MM1-7B-Chat 7.6b [MGF ⁺ 24]	LLaVA-1.6 Vicuna-7b [LLLL23]	LLaVA-Next LLama3-8b [LLL ⁺ 24]	Qwen-VL-Chat 9.6b [BBY ⁺ 23]	Claude 3 haiku [Ant24]	Gemini 1.0 Pro V [TAB ⁺ 23]	GPT-4V-Turbo turbo-2024-04-09
MMMU (val) [YNZ ⁺ 23]	40.4	33.9	37.0	34.2	36.4	39.0	40.7	42.0	55.5
ScienceQA (test) [LMX ⁺ 22]	90.8	69.4	72.6	70.6	73.7	67.2	72.0	79.7	75.7
MathVista (testmini) [LBX ⁺ 24]	44.5	32.0	35.9	31.5	34.8	29.4	33.2	35.0	47.5
Inter-GPS (test) [LGJ ⁺ 21]	38.1	-	-	20.5	24.6	22.3	32.1	28.6	41.0
MMBench (dev-en) [LDZ ⁺ 24]	80.5	75.9	79.0	76.3	79.4	75.8	62.4	80.0	86.1
POPE (test) [LDZ ⁺ 23]	85.8	87.4	86.6	87.2	87.0	82.6	74.4	84.2	83.7
AI2D (test) [KSK ⁺ 16]	76.7	-	-	63.1	66.9	59.8	60.3	62.8	74.7
ChartQA (test) [MLT ⁺ 22]	81.4	-	-	55.0	65.8	50.9	59.3	58.0	62.3
TextVQA (test) [SNS ⁺ 19]	70.9	71.9	72.8	64.6	55.7	59.4	62.7	64.7	68.1

Table 2: Comparison results on public MLLM benchmarks. All the reported numbers are produced with the exact same pipeline to ensure that the numbers are comparable except for MM1-3B-Chat [MGF⁺24] and MM1-7B-Chat [MGF⁺24], which are not publicly available. We adopted the evaluation setting used in Llava-1.5 [LLLL23], without any specific prompt or pre-processing image for all results. These numbers might differ from other published numbers due to slightly different prompts.

6.4 Weakness

Regarding the multi-modal LLM capabilities of our Phi-3-Vision, it performs admirably across various fields. However, we have identified certain limitations, particularly with questions necessitating high-level reasoning abilities. Additionally, the model has been observed to occasionally generate ungrounded outputs, making it potentially unreliable in sensitive areas, such as finance. To mitigate these issues, we will incorporate more reasoning-focused and hallucination-related DPO data into post-training in the future.

From a responsible AI standpoint, whilst safety post-training has made significant strides, our Phi-3-Vision occasionally fails to refrain from answering harmful or sensitive inquiries. Examples of such occasions include deciphering particular types of captcha and describing scam images containing disinformation or hallucination. We find that this issue partly arises from the capabilities, such as OCR, acquired during the training process with normal instruct tuning datasets, which can be regarded as the trade-off between helpfulness and harmlessness. Moving forward, we need to further explore this area to achieve a better balance.

	Phi-3-Vision 3.8b+0.3b	Phi-3-Vision w/o safety 3.8b+0.3b	Llava-1.6 Vicuna 7b+0.3b	Qwen-VL-Chat 7.7b+1.9b	GPT4-V N/A
Internal (private)	8.30	7.06	5.44	7.27	8.55
RTVLM (public)	4.64	3.56	3.86	4.78	6.81
VLGuard (public)	9.12	4.66	5.62	8.33	8.90

Table 3: Comparison results on public and private multi-modal RAI benchmarks. Note that all metrics in the table are [0,10] and a higher value indicates a better performance.

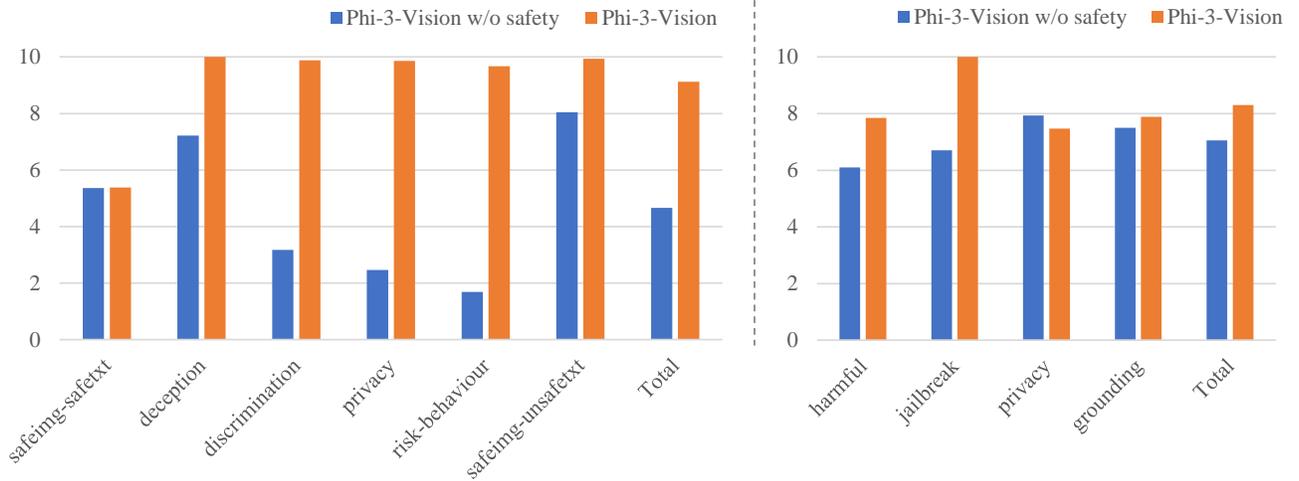


Figure 7: Comparison of categorized RAI performance of Phi-3-Vision with and without the safety post-training on the VLGuard (left) and Internal (right) benchmark, respectively. It clearly indicates that safety post-training can enhance the RAI performance across nearly all the RAI categories.

References

- [AI] Meta AI. Introducing meta llama 3: The most capable openly available llm to date.
- [Ant24] AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 2024.
- [AON⁺21] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- [BBY⁺23] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023.
- [BJN⁺22] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022.

- [BSA⁺24] Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Röttger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. Safety-tuned llamas: Lessons from improving the safety of large language models that follow instructions, 2024.
- [BZGC19] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. *arXiv preprint arXiv:1911.11641*, 2019.
- [CCE⁺18] 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.
- [CKB⁺21] 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. *arXiv preprint arXiv:2110.14168*, 2021.
- [CLC⁺19] 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 *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2924–2936, 2019.
- [CTJ⁺21] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021.
- [DFE⁺22] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.
- [DZZ⁺24a] 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, 2024.
- [DZZ⁺24b] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Songyang Zhang, Haodong Duan, Wenwei Zhang, Yining Li, et al. Internlm-xcomposer2-4khd: A pioneering large vision-language model handling resolutions from 336 pixels to 4k hd. *arXiv preprint arXiv:2404.06512*, 2024.
- [GZA⁺23] Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Gustavo de Rosa Piero Kauffmann, Olli Saarikivita,

- Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*, 2023.
- [HBK⁺21] 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, 2021.
- [HBM⁺22] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Eliza Rutherford Trevor Cai, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [JBA⁺23] Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, Jyoti Aneja, Caio César Teodoro Mendes, Weizhu Chen, Allie Del Giorno, Ronen Eldan, Sivakanth Gopi, Suriya Gunasekar, Piero Kauffmann, Yin Tat Lee, Yuanzhi Li, Anh Nguyen, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Michael Santacrose, Harkirat Singh Behl, Adam Tauman Kalai, Xin Wang, Rachel Ward, Philipp Witte, Cyril Zhang, and Yi Zhang. Phi-2: The surprising power of small language models. *Microsoft Research Blog*, 2023.
- [JCWZ17] Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension, 2017.
- [JLD⁺23] Jiaming Ji, Mickel Liu, Juntao Dai, Xuehai Pan, Chi Zhang, Ce Bian, Chi Zhang, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a human-preference dataset, 2023.
- [JPO⁺20] Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams, 2020.
- [JSM⁺23] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023.
- [JSR⁺24] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Léo Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024.
- [KLZ⁺23] 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.

- [KMH⁺20] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [KSK⁺16] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images, 2016.
- [LBE⁺23] Yuezhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023.
- [LBX⁺24] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts, 2024.
- [LDZ⁺23] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models, 2023.
- [LDZ⁺24] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your multi-modal model an all-around player?, 2024.
- [LGJ⁺21] Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning, 2021.
- [LHE22] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods, 2022.
- [LLL⁺24] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024.
- [LLLL23] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023.
- [LLY⁺24] Mukai Li, Lei Li, Yuwei Yin, Masood Ahmed, Zhenguang Liu, and Qi Liu. Red teaming visual language models. *arXiv preprint arXiv:2401.12915*, 2024.
- [LMX⁺22] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- [LST⁺24] Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander Rush, Douwe Kiela, et al. Obelics: An open web-scale filtered dataset of interleaved image-text documents. *Advances in Neural Information Processing Systems*, 36, 2024.
- [MCKS18] 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.

- [MGF⁺24] Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, Anton Belyi, Haotian Zhang, Karanjeet Singh, Doug Kang, Ankur Jain, Hongyu Hè, Max Schwarzer, Tom Gunter, Xiang Kong, Aonan Zhang, Jianyu Wang, Chong Wang, Nan Du, Tao Lei, Sam Wiseman, Guoli Yin, Mark Lee, Zirui Wang, Ruoming Pang, Peter Grasch, Alexander Toshev, and Yinfei Yang. Mm1: Methods, analysis & insights from multimodal llm pre-training. *arXiv preprint arXiv:2403.09611*, 2024.
- [MHJ⁺23] Ahmed Magooda, Alec Helyar, Kyle Jackson, David Sullivan, Chad Atalla, Emily Sheng, Dan Vann, Richard Edgar, Hamid Palangi, Roman Lutz, Hongliang Kong, Vincent Yun, Eslam Kamal, Federico Zarfati, Hanna Wallach, Sarah Bird, and Mei Chen. A framework for automated measurement of responsible ai harms in generative ai applications, 2023.
- [MLT⁺22] Ahmed Masry, Do Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [MRB⁺23] Niklas Muennighoff, Alexander M Rush, Boaz Barak, Teven Le Scao, Aleksandra Piktus, Nouamane Tazi, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. Scaling data-constrained language models. *arXiv preprint arXiv:2305.16264*, 2023.
- [NWD⁺20] Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial nli: A new benchmark for natural language understanding, 2020.
- [Ope23] OpenAI. Gpt-4v(ision) system card, 2023. https://cdn.openai.com/papers/GPTV_System_Card.pdf.
- [RHS⁺23] 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.
- [RKH⁺21] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [RWC⁺19] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [SLBBC19] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *arXiv preprint arXiv:1907.10641*, 2019.
- [SNS⁺19] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read, 2019.
- [SRR⁺22] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*, 2022.

- [SSS⁺22] 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.
- [TAB⁺23] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [THLB19] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question answering challenge targeting commonsense knowledge, 2019.
- [TLI⁺23] 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. *arXiv preprint arXiv:2302.13971*, 2023.
- [TMH⁺24] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology, 2024.
- [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- [XWX⁺24] Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. Florence-2: Advancing a unified representation for a variety of vision tasks. 2024.
- [YNZ⁺23] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhua Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi, 2023.
- [ZBY⁺24] Yongshuo Zong, Ondrej Bohdal, Tingyang Yu, Yongxin Yang, and Timothy Hospedales. Safety fine-tuning at (almost) no cost: A baseline for vision large language models. *arXiv preprint arXiv:2402.02207*, 2024.
- [ZCG⁺23] Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models, 2023.
- [ZCS⁺23] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*, 2023.
- [ZHB⁺19] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, 2019.

A Example prompt for benchmarks

Question:

Solve for x : $(-\frac{1}{3})(-4-3x) = \frac{1}{2}$

Options:

- A. $-\frac{5}{6}$
- B. $\frac{7}{6}$
- C. $\frac{5}{3}$
- D. $\frac{1}{6}$

Answer: A

Question:

Which of the following is the body cavity that contains the pituitary gland?

Options:

- A. Abdominal
- B. Cranial
- C. Pleural
- D. Spinal

Answer: B

Question:

Where was the most famous site of the mystery cults in Greece?

Options:

- A. Ephesus
- B. Corinth
- C. Athens
- D. Eleusis

Answer:

B Authors (alphabetical)

Marah Abdin	Jamie Huynh	Olli Saarikivi
Sam Ade Jacobs	Dan Iter	Amin Saied
Ammar Ahmad Awan	Russell J. Hewett	Adil Salim
Jyoti Aneja	Mojan Javaheripi	Michael Santacroce
Hany Awadalla	Xin Jin	Shital Shah
Ahmed Awadallah	Nikos Karampatziakis	Ning Shang
Nguyen Bach	Piero Kauffmann	Hiteshi Sharma
Amit Bahree	Mahoud Khademi	Swadheen Shukla
Arash Bakhtiari	Dongwoo Kim	Xia Song
Jianmin Bao	Lev Kurilenko	Yin Tat Lee
Harkirat Behl	Yuanzhi Li	Andrea Tupini
Alon Benhaim	Yunsheng Li	Xin Wang
Misha Bilenko	Chen Liang	Guanhua Wang
Johan Bjorck	Lars Liden	Lijuan Wang
Sébastien Bubeck	Eric Lin	Chunyu Wang
Caio César Teodoro Mendes	Zeqi Lin	Yu Wang
Martin Cai	Weishung Liu	Rachel Ward
Qin Cai	Mengchen Liu	Wen Wen
Vishrav Chaudhary	Ce Liu	Philipp Witte
Weizhu Chen	Chong Luo	Haiping Wu
Yi-Ling Chen	Li Lyna Zhang	Michael Wyatt
Yen-Chun Chen	Piyush Madan	Bin Xiao
Dongdong Chen	David Majercak	Jiahang Xu
Dong Chen	Matt Mazzola	Can Xu
Parul Chopra	Arindam Mitra	Weijian Xu
Xiyang Dai	Hardik Modi	Sonali Yadav
Gustavo de Rosa	Anh Nguyen	Fan Yang
Allie Del Giorno	Brandon Norick	Ziyi Yang
Matthew Dixon	Barun Patra	Jianwei Yang
Ronen Eldan	Daniel Perez-Becker	Yifan Yang
Victor Fragoso	Thomas Portet	Donghan Yu
Mei Gao	Reid Pryzant	Lu Yuan
Min Gao	Heyang Qin	Chengruidong Zhang
Jianfeng Gao	James R. Lee	Cyril Zhang
Abhishek Goswami	Marko Radmilac	Jianwen Zhang
Suriya Gunasekar	Corby Rosset	Yi Zhang
Emman Haider	Sambudha Roy	Yunan Zhang
Junheng Hao	Olatunji Ruwase	Xiren Zhou

C Acknowledgements

We would like to thank Zhuohan Li, Simon Mo, and Kaichao You from UC Berkeley for sharing their insights on the vLLM kernel.