

Llama-3-Nanda-10B-Chat: An Open Generative Large Language Model for Hindi

The Nanda Team

**Mohamed bin Zayed University of Artificial Intelligence, UAE
Inception, UAE
Cerebras Systems**

Abstract

We introduce *Llama-3-Nanda-10B-Chat*, or *Nanda* for short, a new state-of-the-art Hindi-centric instruction-tuned open generative large language model (LLM). *Nanda* is adapted from the LLaMA-3.1 8B model via continuous training with expansion of transformer blocks, following LLaMA Pro approach. This model employs the decoder-only architecture and has been trained on a mixture of Hindi and English texts. With 10 billion parameters, *Nanda* demonstrates improved knowledge and reasoning capabilities in Hindi, surpassing existing open Hindi and multilingual models of comparable size by a substantial margin; it also achieves highly competitive performance in English. We release *Nanda* as an open-sourced instruction-tuned model and provide a detailed overview of its training, tuning, safety alignment, and evaluation processes. We believe that this release will foster further research in Hindi LLMs and support diverse practical applications across various domains.

Contents

1	Introduction	3
2	Pretraining Data	3
2.1	Preprocessing Pipeline	5
2.2	Mixing Hindi and English Data	5
3	Model	6
3.1	Model Architecture	6
3.2	Model and Training Hyperparameters	7
3.3	Training Infrastructure	7
4	Instruction-Tuning	8
4.1	Instruction-Tuning Data	8
4.2	Instruction-Tuning Setup	9
5	Evaluation	10
5.1	Downstream Evaluation	10
5.2	Safety Evaluation	11
5.3	Generation Evaluation	12
6	Related Work	13
7	Conclusion	15
8	Release Notes	15
8.1	Intended Use	16
8.2	Out-of-Scope Use	16
8.3	Biases, Risks, and Limitations	16
9	Acknowledgments	16

1 Introduction

Recent advancements in transformer-based large language models (LLMs), pre-trained on billions of tokens of web data, have transformed natural language processing (NLP). These models have demonstrated exceptional capabilities in NLP applications and complex multi-step reasoning, allowing them to handle intricate human instructions. Despite these achievements, most research and development efforts have focused on English. Some efforts in the space of multilingual LLMs, such as Falcon [AAA⁺23], PALM [CND⁺22], the latest Aya [UAY⁺24] and Llama-3.1 [DJP⁺24], aim to broaden the linguistic capabilities of the underlying LLM. However, pretraining these models continues to rely extensively on English-centric data, which limits the generative performance in other languages. Moreover, these models face the “*curse of multilinguality*” [PGL⁺22, ABF⁺19, CKG⁺20], i.e., the performance tends to decline when the models are trained to cover a vast number of languages. Thus, models tailored to specific languages or language subsets often outpace them. We aim to bridge this gap for Hindi, which is the fourth most spoken language with over 572 million speakers.¹² Specifically, we introduce *Llama-3-Nanda-10B-Chat*, or *Nanda* for short, a robust Hindi-centric decoder-only LLM with 10 billion parameters, built on top of the Llama-3 model [DJP⁺24].³

The primary challenge in developing a Hindi LLM is the limited availability of high-quality Hindi data [JSB⁺20]. Unlike English, which benefits from the availability of corpora of up to 15 trillion tokens [TRP⁺24], Hindi resources are comparatively scarce. To address this, we curated an extensive Hindi corpus containing 65 billion tokens, which was the foundation for continuously pretraining our model. As part of this effort, we developed a specialized Hindi text processing pipeline that includes thorough data filtering and cleaning to ensure data quality. Another challenge was the lack of high-quality code-mixed and romanized Hindi data; we addressed this by incorporating both code-mixed and romanized Hindi data into our corpus.

Unlike massively multilingual LLMs such as Bloom [SFA⁺23], Llama-3.1 [DJP⁺24], or Aya [UAY⁺24], which encompass more than 50 languages, our model focuses exclusively on Hindi, with selective integration of English during fine-tuning. Hindi data comprises 100% of the continuous pretraining phase. For instruction fine-tuning, we created a bilingual dataset with 21.5 million English and 14.5 million Hindi instruction-following tokens. To balance the data and strengthen the model performance, we applied oversampling, yielding 64.5 million English tokens and 43.5 million Hindi tokens. Additionally, considering the inherent safety concerns associated with LLMs, we curated a dataset of safety-oriented instructions. This specialized fine-tuning ensures that the model achieves both fluency and cultural sensitivity in Hindi while performing on par with recent English LLMs in reasoning and world knowledge, thus enabling a robust transfer of capabilities between English and Hindi.

We adopted the architecture of Llama-3 [DJP⁺24], which uses the standard dense transformer framework [VSP⁺17] and incorporates recent advancements such as RoPE positional encoding [SLP⁺21] and grouped-query attention [ALTdJ⁺23]. We further developed a custom tokenizer that balances Hindi and English, ensuring optimized bilingual processing.

We evaluated *Nanda* on various Hindi and English NLP benchmarks, addressing - reasoning, knowledge, misinformation, and bias. The results showed that *Nanda* is superior in Hindi compared to other models of similar size while also being competitive in English.

By making our model publicly available, we hope to steer further research and development in this area, stimulating innovation and practical applications that can better serve the Hindi and the global communities. Despite our significant efforts to ensure safety, we recognize that our model is not foolproof and may not cover all cases. Therefore, we strongly urge all adopters to exercise caution and perform additional safety tests before deploying the model. For this purpose, we outline responsible release notes in Section 8.

2 Pretraining Data

The *Nanda* LLM is trained on hundreds of billions of words to develop a robust foundation in Hindi, building a comprehensive knowledge base specialized to the cultural context of India. Using our largest available Hindi dataset, *Nanda* incorporates diverse sources such as websites, wikipedia, news articles, and Hindi books. This dataset draws from various publicly accessible Hindi-specific resources, including datasets from HuggingFace,⁴

¹<https://en.wikipedia.org/wiki/Hindi>

²<https://www.worlddata.info/languages/hindi.php>

³<https://huggingface.co/MBZUAI/Llama-3-Nanda-10B-Chat>

⁴<https://huggingface.co/datasets/>

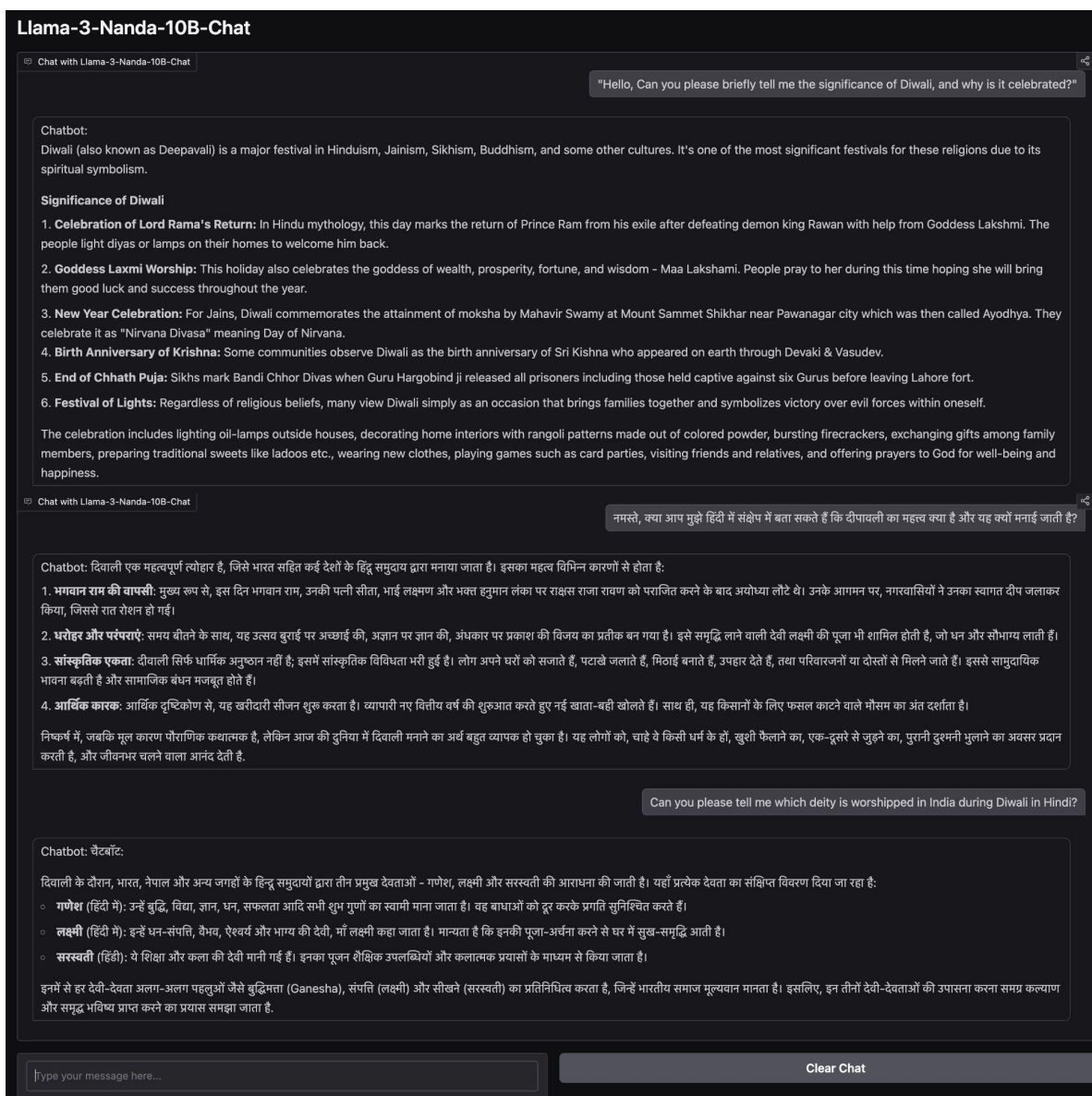


Figure 1: Hindi dialogue using *Llama-3-Nanda-10B-Chat*.

parallel corpora from IIT-B,⁵ and cleaned Common Crawl data from platforms like HPLT,⁶ and Opus-Corpora.⁷ Key datasets sourced from Hugging Face include:

- ai4bharat/sangraha⁸
- uonlp/CulturaX⁹
- allenai/MADLAD-400¹⁰

⁵https://www.cfilt.iitb.ac.in/iitb_parallel/

⁶<https://hplt-project.org/datasets/>

⁷<https://opus.nlpl.eu/>

⁸<https://huggingface.co/datasets/ai4bharat/sangraha>

⁹<https://huggingface.co/datasets/uonlp/CulturaX>

¹⁰<https://huggingface.co/datasets/allenai/MADLAD-400>

2.1 Preprocessing Pipeline

Preprocessing is essential for training high-quality LLMs, as it involves filtering, normalizing, and cleaning the data. To create our 38 billion word Hindi dataset, we designed a comprehensive preprocessing pipeline, combining standard procedures with specific modules to extract quality Hindi content. As illustrated in Figure 2, the raw data mainly originates from publicly accessible databases, some of which were already preprocessed or tokenized. To ensure uniformity, we began by detokenizing all inputs, which allowed us to standardize the content without altering non-tokenized text. Each document at this stage corresponds to an article or web page, depending on the data source.

The pipeline employed stringent filtering rules to exclude noisy or substandard documents. This included removing texts that were too short or excessively long, or those with insufficient Hindi content, which could indicate the presence of another language where Hindi characters appear sporadically. Documents containing overly lengthy words, often signs of URLs or other noisy elements, were also discarded.

After filtering, we proceeded to clean and normalize the data. This involved removing non-printable Unicode characters, embedded scripts like JavaScript or HTML, and frequently repeated boilerplate text (e.g., recurring names of news channels). We also standardized Hindi punctuation and used a lightweight n-gram language model to filter out problematic n-grams.

A final fuzzy deduplication step, leveraging locality-sensitive hashing, was applied to eliminate duplicate content, ultimately reducing the dataset size to 42% of its raw text.

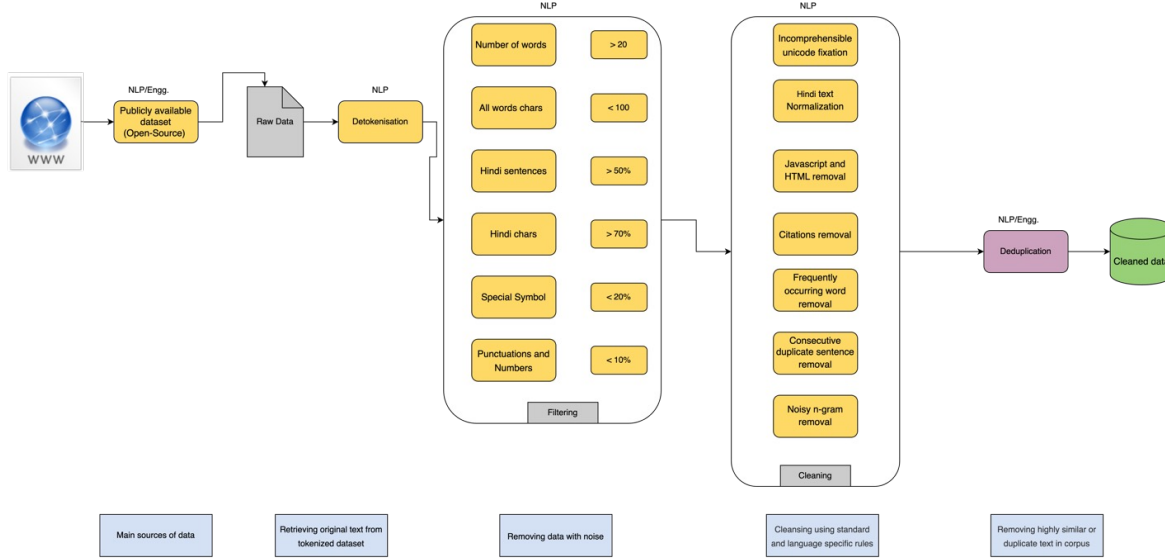


Figure 2: Our Hindi preprocessing pipeline.

Developing the preprocessing pipeline for Hindi posed greater challenges compared to English. While English preprocessing pipelines benefit from numerous large-scale, open-access datasets and well-established techniques, Hindi requires a custom-built approach. Insights gained from experiments with smaller LLMs and the preprocessing pipeline for the dataset used for *Jais* [SSJ⁺23] guided the selection of heuristics used in the final pipeline for *Nanda*’s dataset. Due to the limited availability of Hindi data, we applied less aggressive filtering than typically used for English, ensuring that valuable Hindi content was retained.

2.2 Mixing Hindi and English Data

During the adaptation of the Llama-3 model, we mix Hindi and English data following the findings of [GXR⁺24]. Domain adaptation involves continual pre-training of a foundation model on new data not seen during the pre-training. When this new domain data is out-of-distribution, it can cause significant forgetting of prior capabilities, which is referred to as a stability gap [GFZ⁺24]. Adding a small proportion of replay data, which is closer in distribution to the pre-training data, can mitigate forgetting. We conduct exhaustive experiments to find a minimum proportion of English data that should be mixed with Hindi to mitigate forgetting.

For the Hindi adaptation of Llama-3-8B, we found that a higher proportion of replay data is necessary. Therefore, a 1:1 English-to-Hindi dataset mixture was used where we saw cross-lingual capability transfer between the two languages and also avoid saturation of Llama-3 for adaptation. For replay data, we needed a mix of textbooks, math, coding and reasoning datasets to mitigate forgetting.

3 Model

3.1 Model Architecture

Nanda is based on a standard transformer-based architecture [VSP⁺17]. In particular, we adapt a causal decoder-only model, similar to the one used by GPT-2 [RWC⁺19] and Llama [TLI⁺23]. Decoder-only models have achieved state-of-the-art performance in generative language tasks. Building upon this base transformer architecture, we use several recent improvements from the literature and our experiments.

Nanda Tokenizer: The first step in adapting a monolingual foundation model for multilingual use is to construct a balanced vocabulary that includes all target languages. Recent state-of-the-art models such as Llama-3 [TLI⁺23] use byte pair encoding (BPE) [SHB16] tokenizers, primarily trained on English data. These monolingual English tokenizers often split non-English words into characters or bytes, creating a significant imbalance among languages. Fertility [RPV⁺21], which measures the average number of subwords produced by a single word upon tokenization, can be used to quantify this imbalance.

This imbalance introduces inefficiency in pre-training, fine-tuning and inference. Table 1 shows that the Llama-3 tokenizer needs as many as 2.6 times the number of tokens to represent the same Hindi text as the Hindi-English bilingual tokenizer. Balanced multilingual tokenizer with low fertility in all languages offers three main advantages [PLMTB23]: (1) lower training and inference cost; (2) reduced latency during inference; and (3) longer context windows. Models trained with low fertility tokenizers tend to perform well on downstream tasks, as shown in [ADH⁺23]

We extend the Llama-3 vocabulary to create a balanced tokenizer for English and Hindi. Vocabulary extension adds the most frequent Hindi tokens in the corpora, leading to a larger vocabulary size. Additionally, we ensure that the newly introduced tokens are not present in the original vocabulary. We conduct vocabulary extension studies to determine the optimal number of new Hindi tokens to be added, ensuring a balanced multilingual vocabulary. The Hindi tokens are borrowed from a monolingual Hindi tokenizer trained on the Hindi corpora. We create a few candidate extended tokenizers and perform intrinsic evaluations following [AFT⁺24].

For intrinsic evaluation, we use fertility score to measure the efficiency of the tokenization process [GXR⁺24]. Fertility is defined as $f = \frac{S}{W}$, where S is the total number of tokens in the tokenized text and W is the number of words in the raw text. Note that fertility is calculated on held-out subsets from the Hindi corpora, which are not used for tokenizer training.

Table 1 shows the intrinsic evaluations of three candidate tokenizers, (i) *Llama-3-extend10*, (ii) *Llama-3-extend20*, and (iii) *Llama-3-extend30*, which extend the Llama-3 vocabulary by 10%, 20%, and 30%, respectively. Based on our tokenizer fertility ablation studies, *Llama-3-extend20* reduces the fertility of Llama-3’s tokenizer by 54.40% while maintaining the fertility in English. It also reaches a fertility score of 1.19, which is comparable to the English fertility for the base Llama-3 tokenizer in English, which is 1.24. Extending the vocabulary further to 30% shows minimal improvement in Hindi fertility. Therefore, we select *Llama-3-extend20* as the tokenizer for the *Nanda* model.

	Llama-3	Llama-3-extend10	Llama-3-extend20	Llama-3-extend30
Vocab Size	128,256	141,081	153,856	166,732
Hindi Fertility	2.61	1.27 (-51.34%)	1.19 (-54.40%)	1.16 (-55.55%)

Table 1: Tokenizer intrinsic evaluation across different vocab sizes. We see that extending the tokenizer with Hindi vocab, reduces the fertility in *Llama-3-extend10* by 51.34%, *Llama-3-extend20* by 54.40%, and *Llama-3-extend30* by 55.55% compared to the Llama-3 tokenizer.

Nanda Embedding: Following the methods outlined for embedding initialization in [GXR⁺24], we use a semantic similarity search-based embedding initialization method. This method uses Wechsel multilingual

initialization [MPR22] where pre-trained embeddings like Fasttext or OpenAI embeddings are used.

For each new Hindi token added to the Llama-3 base vocabulary, we identify the top- k most similar tokens in the base vocabulary based on cosine similarity using embeddings from a pre-trained embedding model. We use OpenAI’s `text-embedding-3-large` embeddings [KBR+24] for its superior quality and multilingual capabilities. To initialize the embeddings of the new Hindi token, we take a weighted average of the top- k similar tokens’ base embeddings. After experimenting with different values for the k , we achieve the best results with $k = 5$. This initialization method was used for embeddings and unembeddings layers of *Nanda*.

Nanda Architecture: Following [WGG+24], we leverage the block expansion approach, which proves to be highly effective for language adaptation, especially for low-resource languages. By adding and fine-tuning additional Transformer blocks initialized to identity mappings, the model can integrate new domain and language specific knowledge without forgetting previous information. Although the techniques described in [WGG+24] focus on code and math adaptation, we were able to successfully adopt this approach for language adaptation. We initialized our base model with Llama-3-8B and expanded the number of decoder blocks from 32 to 40 using an interleaved approach. A new decoder block was added every 4 decoder blocks in the base Llama-3 model. In our language adaptation experiments, we found that an optimal data mix of 1 : 1 (En:Hi) yielded the best results (in downstream 0 shot tasks in both English and Hindi) compared to Hindi-only adaptation. In both experiments, we trained on a total of 55B tokens for Hindi in order to maintain the same token count for the appropriate comparison. Our results show that the block-expansion approach is a strong candidate for language adaptation with less training overhead and resources compared to training domain-specific models from scratch, especially for low-resource languages. In the future, this work could expand to other architectures (like Mixture-of-Experts) and modalities, and it would be interesting to analyse the impact on overall accuracy in downstream tasks. Following the results from Gosal et al. [GXR+24], we find that the optimal adapter layers are 25% of the existing layers.

3.2 Model and Training Hyperparameters

Table 2 shows the number of layers, heads, and dimensionality for *Nanda*, along with the optimization hyperparameter values, peak learning rate and batch size.

For the continual pre-training dataset, we sampled documents from the source list described in Section 2 and generated sequences with a context length of 8,192 tokens. When a document was smaller than 8,192 tokens, it was concatenated with other document (documents) and packed into one sequence. `<|endoftext|>` is used to demarcate the end of each document, giving the language model the information necessary to infer that tokens separated by `<|endoftext|>` are unrelated.

Model	Layers	Heads	Dimension	Learning Rate	Batch Size
<i>Nanda</i>	40	32	4,096	$1.5e^{-5}$	4e6

Table 2: **Training hyperparameter values:** the number of layers, heads, and dimensionality for *Nanda*, along with the optimization hyperparameter values and peak learning rates.

We train *Nanda* using the AdamW optimizer [LH18] with $\beta_1 = 0.9$, $\beta_2 = 0.95$, $\epsilon = 1e - 5$, and weight decay of 0.1. We scale the gradient norms using a maximum norm clipping value of 1.0. The learning rate schedule comprises a linear warm-up from 0 to the maximum learning rate in 274 steps, followed by a $10\times$ linear decay until 27,192 steps. After packing, we used a global batch size of 7,680 sequences, each with 8,192 tokens.

3.3 Training Infrastructure

All training, hyper-parameter tuning, and instruction-tuning experiments were executed on the Condor Galaxy 2 (CG-2) AI supercomputer from Cerebras,¹¹ built in partnership with G42. The final training and fine-tuning runs for *Nanda* were performed on 16 CS-2 systems within CG-2. CG-2 is a Cerebras Wafer-Scale Cluster composed of 64 Cerebras CS-2 systems, MemoryX, SwarmX, management, and input worker nodes. The

¹¹www.cerebras.net/blog/introducing-condor-galaxy-1-a-4-exaflop-supercomputer-for-generative-ai/

foundation of the CG-2 cluster is the Cerebras Wafer Scale Engine (WSE) within the CS-2 system, the largest and most powerful AI processor currently available.

CS-2 systems are purpose-built network-attached AI accelerators. Each CS-2 features 40 GB of SRAM and a peak of 62.5 AI PetaFLOPs, providing a total of 4 ExaFLOPs of AI compute across 64 systems in the CG-2 supercomputer. Utilizing the weight streaming mode of the Cerebras software stack, the Condor Galaxy supercomputers can flexibly schedule multiple jobs based on hardware resource requirements and priority. The number of CS-2s allocated to a job can be dynamically adjusted during training, with performance scaling linearly up to 64 CS-2s per job. This scalability is facilitated by the Cerebras software stack’s use of pure data parallelism to distribute the workload across multiple CS-2s. Jobs are managed by a priority queue system, ensuring efficient allocation of computational resources.

MemoryX is a large-capacity off-wafer memory service used to store all model weights, gradients, and optimizer states. SwarmX is a broadcast/reduce fabric that connects the memory service MemoryX to each of the CS-2 systems in a wafer-scale cluster. Swarm-X coordinates the broadcast of the model layer weights, giving each CS-2 a local copy, and it receives and aggregates (by addition) the independent weight gradients coming from the CS-2 systems during backpropagation. At the end of each iteration, the aggregated gradients are sent to MemoryX for weight update.

The CG-2 hardware and software stack enables training extremely large models using data parallelism by relying on a special execution mode available with Cerebras Wafer Scale Clusters, called weight streaming. Weight streaming fully bypasses the complexity of 3D parallelism on traditional GPU clusters and provides simpler and higher performance scaling.

4 Instruction-Tuning

A useful LLM is one that can interpret user instructions across a variety of NLP tasks and then correctly execute each task to meet the user’s preferences for helpfulness and safety. We develop *Nanda* as a model capable of handling a variety of NLP tasks in both Hindi and English. However, in its pre-trained form, the model lacks the ability to fully interpret and respond to user instructions accurately. To this end, we conduct instruction-tuning [OWJ⁺22] on the pre-trained model, aligning it with practical use cases and enhancing safety by training on a dataset comprising English and Hindi instructions.

4.1 Instruction-Tuning Data

Nanda is a bilingual model, which means that it must be enabled to understand instructions in Hindi without compromising its performance in English. Therefore, we prepare a diverse instruction-tuning dataset covering various domains, with instructions provided both in Hindi and in English. We have a total of ~ 61 K prompt–response pairs in our dataset, and we provide a brief overview of the data collected for each language in the following sections.

4.1.1 English Instruction-Tuning Data

We curate a set of English instructions in an expert-written prompt–response pair format spanning a comprehensive range of NLP tasks. Some of our data is a subset of the instruction-tuning data used for building *Jais* [SSJ⁺23], and hence, is a combination of several publicly available datasets. We also add a significant amount of proprietary data that we gather from relevant internal sources. Our data comprises ~ 39 K high-quality English instructions consisting of 7.7M tokens in the prompts and 9M tokens in their responses (a total of ~ 16 M tokens). Specifically, we have close to 20K instructions focused on mathematics, while the rest of the examples cover code and diverse types of reasoning, such as physical, logical and causal reasoning.

Since a part of our data is out of reach for the public, we are actively working towards making our instruction-tuning data publicly accessible to support research in this field; however, we are currently unable to specify a timeline for when this dataset will be ready for release.

4.1.2 Hindi Instruction-Tuning Data

As a relatively lower-resource language, Hindi does not have many publicly available, high-quality instruction-tuning datasets. Several existing approaches have utilized machine translation on subsets of English instruction-tuning datasets to create datasets for lower-resource languages. One such dataset, specifically curated for Hindi,

was recently introduced and was used to train *Airavata* [GJH⁺24a]. We initially experimented with this dataset, but the results were not very strong, so we created our own dataset.

We selected a subset of high-quality instructions (excluding any math or code instances) from our English instruction-tuning dataset and used *GPT-4* to translate it to Hindi [Ope23]. To ensure the quality of the translations, Hindi language experts simultaneously verified a sample of instances in the generated dataset. In addition, we realize that Hindi speakers often use a more relaxed form of the language during informal interactions. We aim for our model to be adept at understanding both formal and informal styles of written Hindi. Thus, we prompted *GPT-4* to generate two kinds of translations:

- **Formal Hindi:** The translated instances must be written in Devanagari script with a style of writing consistent with official documents in Hindi.
- **Casual Hindi:** The model is encouraged to generate translations that contain Hindi (and some English) words using a mix of Devanagari and Latin scripts. This form of language is generally used by Hindi-speaking individuals during informal conversations like texting, interactions on social media, and more.

Ultimately, our curated dataset includes $\sim 22\text{K}$ high-quality machine-translated Hindi instructions, split into $\sim 13.5\text{K}$ in formal Hindi examples and the remaining in casual Hindi. In total, the Hindi instruction-tuning dataset comprises 3.8M prompt tokens and 10M response tokens.

4.1.3 Safety Data

We developed a comprehensive safety prompt collection process specifically tailored for Hindi model training, covering eight types of attacks and over 100 detailed safety categories. To ensure high-quality data, a team of five expert annotators initially crafted “seed prompts” for direct attack alignment based on our previous work [WLH⁺23], resulting in approximately 1,200 annotated examples focused both on general and Hindi-specific scenarios. Building on this foundation, the expert team guided a 20-member outsourced annotation team, leveraging LLMs, to generate an additional 50K attack prompts, ensuring diversity, linguistic relevance, and thorough coverage for Hindi.

We enrich the set of direct attack prompts in SFT data with a collection of adversarial prompt attack methods. Following [LMZ⁺24], we adopt eight adversarial prompt attack methods to construct the SFT data. These methods target the following abilities of LLMs: in-context learning, auto-regressiveness, instruction following, and domain transfer, resulting in 100K attack prompts.

To further improve the robustness and generalizability of our model against adversarial prompt attacks, we also adopt LLM-based methods for diversifying the attack prompts. This can also help prevent over-fitting on the attack template used by the works that proposed these attacks.

Moreover, in the over-refusal prompts task, annotators generate 50K questions that closely resemble potentially unsafe adversarial prompts but are deliberately crafted to be entirely safe. The primary motivation for this task is to address the overrefusal behavior commonly seen in LLMs [CCSH24], where models refuse to answer benign questions due to excessive caution. By including these prompts, we aim to train the model to better distinguish between genuinely unsafe queries and safe ones, thereby improving the model’s responsiveness while maintaining safety.

4.2 Instruction-Tuning Setup

As mentioned in Section 4.1, the instances in our raw instruction-tuning data contain a system instruction and a pair of a user-prompt and an AI response. In the case of multi-turn interactions, we have a sequence of multiple prompt–response pairs. Since our model is built on top of *Llama-3-8B-Instruct*, we template each raw datapoint using the *Llama-3-Instruct* prompt template both for supervised fine-tuning (SFT) and for inference.¹² We illustrate the transformation of the raw data points to follow the prompt template in Figure 3. At this stage, we oversample the instructions in our dataset (excluding safety instruction-tuning data) to 300% of the original quantity to strengthen the model. This means we perform SFT over approximately 100M tokens consisting of 47M tokens in Hindi instructions and 53M of the same in English instructions. Moreover, similar to *Jais* [SSJ⁺23], we apply padding to each templated instance, use the same autoregressive objective as for pretraining, and mask the loss of the prompt to make sure backpropagation considers only the answer tokens during SFT.

¹²<https://www.llama.com/docs/model-cards-and-prompt-formats/meta-llama-3/>

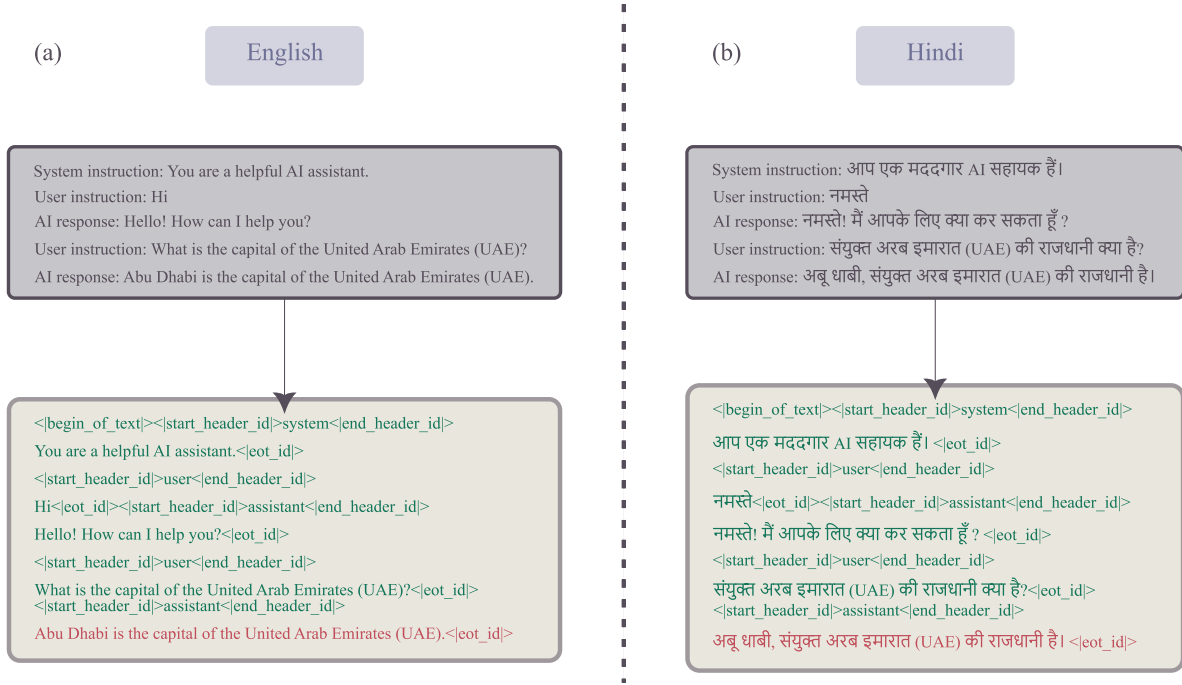


Figure 3: Examples of how the raw data looks like after being transformed to follow the Llama-3 Chat template: the prompt is in green, and the response is in red. In the figure, (a) shows a multi-turn instruction in English, and (b) shows the same interaction in Hindi.

5 Evaluation

In this section, we aim to provide a thorough assessment of the *Nanda* model across a diverse set of evaluation dimensions, covering downstream NLP tasks, safety assessments, and generation capabilities. These evaluations are designed to rigorously measure the model’s performance and adaptability, particularly in supporting multilingual use cases across both Hindi and English languages.

5.1 Downstream Evaluation

Evaluation Setup We conduct a comprehensive downstream evaluation, comparing the *Nanda* model to a series of baselines that support both Hindi and English languages. Our baseline models include models that are specifically optimized for the Hindi language, such as Gajendra-v0.1 [Bha24], Airavata [GJH⁺24b], sarvam-2b-v0.5 [Sar24], alongside several models from the AryaBhatta [Gen24a, Gen24b] and Aya series [ADT⁺24]. Additional models include popular general-purpose models like Gemma [RPS⁺24], Llama-2 (7B, 13B) [TMS⁺23], and the latest Llama-3.1-8B [DJP⁺24] models.

We adopt the LM-Evaluation-Harness framework [GTB⁺21] to evaluate each model in a zero-shot setting, and we report the accuracy for each task. Within the LM-Evaluation-Harness framework, the context string is concatenated with each candidate output string, and the answer is determined by selecting the concatenated string with the highest normalized log-likelihood.

Datasets We perform a comparative evaluation of *Nanda* against other LLMs for both Hindi and English, building upon the evaluations conducted in prior studies [DJP⁺24, ADT⁺24, Ope23]. For each language, our evaluation encompasses aspects such as knowledge, reasoning, and misinformation, as outlined in Table 3 and Table 4. For Hindi, we assess performance on four translated benchmarks—MMLU, HellaSwag, ARC-Easy, and ARC-Challenge—sourced from Indic-Eval [K24], as well as on Hindi-TruthfulQA [DLVNN⁺23] to evaluate misinformation. For English, following prior studies, we include MMLU [HBB⁺20], HellaSwag [ZHB⁺19], ARC-Easy [CCE⁺18a], and TruthfulQA [LHE21] benchmarks.

Results for Hindi Table 3 presents the zero-shot and few-shot evaluation results for Hindi. The *Nanda* models demonstrate superior performance across many evaluation criteria, setting a new benchmark for state-of-the-art

Model	Average	MMLU		HellaSwag	ARC-Easy	ARC-Challenge	TruthfulQA
		0-shot	5-shot	0-shot	0-shot	0-shot	0-shot
Gemma-2-9B-base	30.20	26.9	27.2	27.1	28.2	23.6	48.2
Llama-2-7B	31.02	27.9	28.1	29.6	29.3	24.9	46.3
Llama-2-13B	31.30	28.3	29.3	30.6	29.2	26.6	43.8
Gajendra-v0.1	31.47	27.4	27.9	33.0	36.7	26.6	37.2
Airavata	32.02	28.1	28.5	33.0	32.0	25.6	44.9
sarvam-2B-v0.5	37.70	28.3	29.1	46.2	45.8	32.3	44.5
AryaBhatta-GemmaOrc	39.43	31.4	35.9	42.6	46.5	32.7	47.5
Llama-3-8B	39.83	30.2	37.3	45.7	45.9	34.5	45.4
Aya-23-8B	40.18	29.8	36.8	48.4	48.3	33.9	43.9
Llama-3.1-8B	40.42	29.9	37.3	46.9	50.2	34.3	43.9
AryaBhatta-GemmaUltra	41.18	34.6	37.5	45.5	48.9	33.4	47.2
Llama-3.1-8B-Instruct	41.80	32.9	38.9	48.0	50.5	36.2	44.3
Llama-3-Nanda-10B-Chat	47.88	38.6	44.3	56.4	59.6	40.3	48.1

Table 3: Evaluation results for **Hindi** language models. *Average* represents the mean score across tasks, and *0-shot* indicates zero-shot results, while *5-shot* denotes few-shot results.

Model	Average	MMLU	HellaSwag	ARC-Easy	TruthfulQA
		0-shot	0-shot	0-shot	0-shot
Gemma-2-9B-base	33.03	28.4	33.1	24.2	46.4
sarvam-2B-v0.5	42.83	29.4	61.7	42.5	37.7
Airavata	44.53	31.7	65.5	40.1	40.8
Llama-2-7B	46.00	31.1	72.9	40.5	39.5
Gajendra-v0.1	48.55	37.5	73.0	43.0	40.7
Aya-23-8B	49.63	34.0	73.9	45.2	45.4
Llama-2-13B	51.20	36.9	77.7	46.1	44.1
AryaBhatta-GemmaOrca	53.03	40.4	72.4	45.4	53.9
AryaBhatta-GemmaUltra	53.65	42.5	74.1	45.4	52.6
Llama-3-8B	53.65	39.2	79.1	52.3	44.0
Llama-3.1-8B	54.33	39.7	78.9	53.5	45.2
Llama-3.1-8B-Instruct	57.53	41.8	79.3	55.1	53.9
Llama-3-Nanda-10B-Chat	59.45	48.7	79.2	53.7	56.2

Table 4: Evaluation results for **English** language models. *Average* represents the mean score across tasks, and *0-shot* indicates zero-shot results.

Hindi language models. Specifically, compared to monolingual Hindi models, such as Gajendra-v0.1, Airavata, sarvam-2b-v0.5, AryaBhatta series models, *Nanda* (10B) achieves significant absolute improvements ranging from +5.49 to +15.2 points. These gains are especially notable in across knowledge retrieval, commonsense reasoning and misinformation.

We can further see that among multilingual models, Llama-3.1 and Aya-23-8B are the best-performing models, with an average accuracy of 41.8 and 40.42, respectively. Notably, Llama-2 and Gemma-2-9b lag behind, which should not be surprising given their limited exposure to Hindi pre-training data. We see that *Nanda* (6.7B) outperforms Llama-3.1-8B by +6.08 points absolute.

Results for English We also conducted an evaluation for English, with the results shown in Table 4. Notably, *Nanda* achieves a slight improvement over existing English models, even though its continual pretraining was solely on Hindi data. Additionally, we observe that, apart from the AryaBhatta series, other Hindi models such as Gajendra-v0.1, Airavata, and sarvam-2b-v0.5 exhibit significantly lower performance compared to established English models.

5.2 Safety Evaluation

We followed previous work [WLH⁺23] and constructed a novel dataset for Hindi safety evaluation, aiming to identify biases and harmful content within the language model, specifically focused on the Hindi language and

Model	English	Hindi
Aya-23-8B	49.48	63.79
Llama-3.1-8B-Instruct	90.99	87.01
Llama-3-Nanda-10B-Chat	85.97	87.96

Table 5: Safety evaluation results.

cultural context. The dataset includes a comprehensive categorization of risk areas, types of harm, and specific examples to enable thorough evaluation.

Taxonomy Development The development of a detailed taxonomy was the first step in constructing this dataset. This taxonomy categorizes risk areas specific to Hindi, including regional bias, economic situation bias, and national/group character bias. The taxonomy defines specific harms, such as instances of prejudice against particular states in India or negative stereotypes about national characteristics. Example questions were curated to illustrate these biases, helping ensure the evaluation captures a broad range of potential issues.

Data Collection and Translation The dataset incorporates content sourced in English [WLH⁺23], initially focused on safety issues like discrimination, toxicity, and adult content, which were then translated into Hindi. The translation process was managed using both automated tools (such as Google Translate and GPT-4) and manual validation by native speakers to ensure the accuracy and cultural relevance of the translations. Each translated entry underwent a thorough validation process to mitigate mistranslations or inadvertent cultural insensitivity.

Annotation and Validation To ensure the quality of the dataset, we collaborated with outsourced annotators who were provided with guidelines to annotate harmful content. The annotations focus on verifying whether translated content preserved the intended meaning and accurately represented harmful or biased elements in the Hindi context. Annotations were then cross-checked to guarantee consistency and reliability in labeling harmful examples.

Evaluation Results The evaluation results are shown in Table 5. We can see that our model is much safer than the other models.

5.3 Generation Evaluation

Dataset In addition to downstream and safety evaluations, we also evaluate the models’ core capability for Hindi text generation. Consistent with previous studies [PLH⁺23, CLL⁺23], we conduct an LLM-as-a-judge evaluation of the generated Hindi text quality using GPT-4 [Ope23]. The evaluation is based on the *Vicuna-Instructions-80* [CLL⁺23] dataset, which was manually translated into Hindi by professional translators.

The *Vicuna-Instructions-80* dataset comprises 80 challenging,¹³ open-ended prompts spanning eight categories: knowledge, Fermi questions, counterfactuals, roleplay, general topics, mathematics and coding, writing, and common-sense reasoning.

Evaluation Setup We generate outputs for Hindi prompts from the *Vicuna-Instructions-80* dataset, using a temperature of 0.3 and a repetition penalty of 1.2. As baselines, we compare with open-source multilingual models: Llama-3-8B-Instruct [DJP⁺24] and Qwen2.5-14B-Instruct [BBC⁺23]. Llama-3-8B-Instruct serves as an ideal baseline since our model is built upon it, ensuring a consistent foundation for comparison. Additionally, Qwen2.5-14B-Instruct brings strong multilingual capabilities and larger parameter capacity, which can better handle linguistic nuances in Hindi.

For the GPT-4 evaluation, we conduct pairwise comparisons between all models. GPT-4 is tasked with scoring each pair between 0 to 10 based on the outputs generated for the prompts in Hindi *Vicuna-Instructions-80*. The model with the higher GPT-4 score in each pair is considered winner for that prompt. To ensure fairness, the outputs from both models in each pair are randomly permuted so that either can appear as the first candidate. The GPT-4 prompt is structured as follows:

¹³<https://lmsys.org/blog/2023-03-30-vicuna/>

You are a helpful and precise assistant for checking the quality of two Hindi assistants. Suppose the user only speaks Hindi, please evaluate both answers with your justification, and provide an integer score ranging from 0 to 10 after your justifications. When evaluating the answers, you should consider the helpfulness, relevance, accuracy, and level of detail of the answers. The score for answer 1 should be wrapped by `<score1>` and `</score1>`, and the score for answer 2 should be wrapped by `<score2>` and `</score2>`.

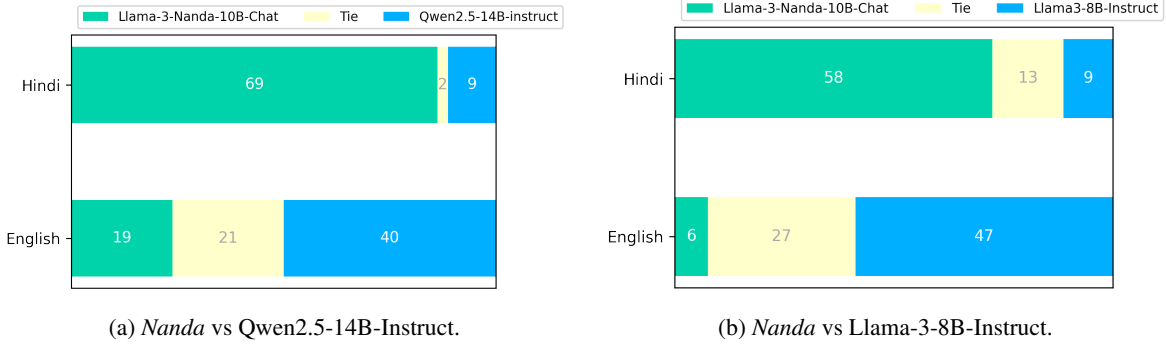


Figure 4: GPT-4 evaluation results for *Nanda* compared to baselines on Vicuna-80 questions.

Results We summarize our finding of generative evaluations in Figure 4. In our assessments, *Nanda* clearly outperforms both Llama-3-8B-Instruct and Qwen2.5-14B-Instruct in Hindi text generation. Built upon the Llama-3 (8B) architecture, *Nanda* keeps the model’s efficiency while incorporating key improvements that make it more attuned to Hindi language. Although Qwen2.5-14B-Instruct offers robust multilingual support, its broader scope may limit it’s ability to capture Hindi-specific nuances as well. Consequently, *Nanda* offers better contextual understanding and generates more natural and fluent Hindi text in language-focused tasks.

Model	English	Hindi
Qwen2.5-14B-Instruct	8.16 \pm 1.90	4.65 \pm 2.56
Llama-3-8B-Instruct	7.98 \pm 1.76	3.36 \pm 2.57
<i>Llama-3-Nanda-10B-Chat</i>	7.67\pm1.91	8.05\pm2.13

Table 6: Average with standard deviation of scores obtained by each model for English and Hindi text generation.

Table 6 offers a closer look into the scores for each model: the average scores and the standard deviation for English and Hindi. Our model, *Nanda*, achieves a higher average score in Hindi text generation tasks (Avg = 8.05) compared to Qwen2.5-14B-Instruct (Avg = 4.65) and Llama-3-8B-Instruct (Avg = 3.36), demonstrating superior performance in Hindi. In English text generation, *Nanda*’s average score (Avg = 7.67) is slightly lower than that of Qwen (Avg = 8.16) and Llama (Avg = 7.98), suggesting that while *Nanda* excels in Hindi. This outcome is anticipated, as the incorporation of additional language-specific optimizations might lead to a slight degradation in performance for non-target languages.

6 Related Work

Below, we discuss previous work on the following relevant topics: LLMs in general, multilingual models, instruction-tuning, and evaluation of LLMs.

Large Language Models Language models with larger parameter sizes have consistently outperformed smaller models like BERT [DCLT19a], [LLG⁺19], and T5 [RSR⁺20]. However, despite being trained on large amounts of multilingual data, many recent large language models still exhibit a strong bias toward English, making them

less effective for other languages [LKW⁺23]. A notable exception to this trend is GLM [ZLD⁺23], which is specifically designed to excel in both Chinese and English.

Current pretraining frameworks for language models are typically divided into three main categories: autoregressive, autoencoding, and encoder-decoder models. Most of the latest large language models, such as the GPT series [RWC⁺19, BMR⁺20, Ope23], Llama [TLI⁺23, TMS⁺23], Bloom [SFA⁺23], and Falcon [AAA⁺23], adopt an autoregressive approach, using a left-to-right objective for language modeling. Older models like BERT [DCLT19b], ELECTRA [CLLM20], and RoBERTa [LOG⁺19] are encoder-only, while models such as BART [LLG⁺19] and T5 [RSR⁺20] utilize an encoder-decoder framework. As discussed in Section 3.1, *Nanda* follows the autoregressive model approach, building on the achievements of models like Llama-2 and GPT-4.

Advancements in large language models can also be categorized into two areas: closed-source and open-source models. Closed-source models like Bard,¹⁴ Claude,¹⁵ Gopher [RBC⁺22], and GPT-4 [Ope23] offer fewer benefits to the research community compared to open-source alternatives [TMS⁺23, LQN⁺23, LTW⁺24]. The lack of transparency in closed-source models introduces various risks, including privacy concerns [MGU⁺22, YRC23] and safety issues [SXD⁺22]. In contrast, *Nanda* is an open-source model, as detailed in 3.

Multilingual Models Pre-training a language model typically involves using unsupervised learning with large datasets. While much of this work has been centered on English [DCLT19a, RWC⁺19, RSR⁺23, BSA⁺23], significant research has also been dedicated to mono-lingual pre-training in languages other than English [FFG⁺24, GFAEP⁺22, ZRS⁺21, SSJ⁺23, PTNT22, KRLB20, KYR⁺23], as well as training models on a small number of languages [NZL⁺24, MHSB21, OZL21, JOOA⁺22].

There have also been massively multilingual pre-training efforts [XCR⁺20, CCG⁺23, SFT⁺23, SFA⁺23, LMA⁺22, DCLT19a, CKG⁺20, KBM⁺21, OAA⁺23, AAMK22, DSK⁺22]. Models based on the mC4 corpus [XCR⁺20], which cover approximately 100 languages, represent the broadest range of coverage in pre-trained models available today. Notable examples include mT5 [XCR⁺20] and umT5 [CCG⁺23], which are the largest publicly accessible multilingual pre-trained models.

However, a key limitation of all these approaches is that they focus on pre-training, requiring users to perform downstream task fine-tuning for specific applications. In contrast, our work emphasizes equipping pre-trained models with instruction-following capabilities.

Another important research direction focuses on adapting pre-trained models to accommodate new languages not included during the initial training phase. These studies explore methods such as continued fine-tuning and embedding space adaptation. For instance, previous work [YSM⁺23, LKL⁺23] has expanded language coverage by gradually adding languages through additional pre-training on monolingual datasets, a method that does not scale efficiently. In a concurrent effort, [LJT⁺24] extends language coverage significantly by using vocabulary expansion and further pre-training Llama-2 with Glot500-c [ILK⁺23].

Hindi has also been integrated into these multilingual models, including earlier models such as mBERT [DCLT19b] and XLM-RoBERTa [CKG⁺20], as well as more recent large language models such as Bloom [SFA⁺23]. However, due to the Hindi content being dwarfed by other languages, these models tend to perform substantially worse than dedicated monolingual models and often exhibit limited generalization abilities in zero-shot settings [LKW⁺23].

Instruction-Tuning Fine-tuning language models using instruction-response pairs has enhanced the generalization capabilities of language models across various tasks [OWJ⁺22]. In terms of open-source models, Bloomz [MWS⁺23] is a fine-tuned version of the foundation model Bloom [SFA⁺23] based on large-scale instruction-tuning over a dataset created via templates, while Llama-2 [TMS⁺23] uses a publicly available instruction-response pair dataset [CHL⁺22]. There has been a rise in Indic models like Airavata [GJH⁺24c], Aya [UAY⁺24] and Nemotron [JSK⁺24] that use publicly available data to instruction tune base models that enable them to perform better at tasks.

The prompts used for instruction-tuning can have diverse origins. Some, as observed by [ZMH⁺23], are human-designed, while others can be autonomously generated. These prompts can be refined with follow-up instructions for more relevant or specific outputs, as studied by [GAS⁺23] and [MTG⁺23]. Recently, [WWS⁺22] introduced *chain-of-thought prompting*, directing models to clarify their reasoning over complex tasks, which was shown to enhance their accuracy.

¹⁴<https://ai.google/static/documents/google-about-bard.pdf>

¹⁵<https://www.anthropic.com/index/introducing-claude>

Evaluating Large Language Models Large language models are highly capable of generating coherent and fluent text but often struggle with factual accuracy and reasoning abilities. To assess factual accuracy, models like GPT-4 [Ope23] and Llama [TLI+23] use school exam-style questions [HBB+22] to gauge how faithfully they can provide knowledge. Common-sense reasoning is also critical and is tested through datasets such as *HellaSwag* [ZHB+19], *WinoGrande* [SBBC21], *ARC* easy and challenge [CCE+18b], and *OpenBookQA* [MCKS18]. For evaluating reasoning through programming, benchmarks like HumanEval [CTJ+21] and MBPP [AON+21] are used.

In the domain of Hindi NLP, [KKG+20] introduced IndicGLUE, the first Indic NLU benchmark for 11 languages, while [DAR+23] expanded upon this by releasing IndicXTREME, covering all 22 Indic languages. On the natural language generation (NLG) side, [KSS+22] developed the IndicNLGsuite, which supports five tasks across 11 languages. Additionally, [GCA+23] presented IN22, a machine translation benchmark for evaluating both conversational and general translation across all 22 languages. More recently, [SGB+24] proposed Indic-GenBench, a benchmark covering diverse tasks such as cross-lingual summarization, machine translation, and cross-lingual question answering. [WGY+24] evaluated models using LLMs and humans and observed that they agree fairly well on most Indic languages.

In contrast, researchers working with other languages often employ machine translation or create analogous datasets to assess language models’ knowledge proficiency and commonsense understanding [Ope23, LKW+23]. As explained in Section 5, we adopted a similar approach by crafting datasets analogous to those available in English and using a combination of human translations and our proprietary machine translation system to convert English datasets into Hindi for evaluation purposes.

While knowledge [HBB+22, LZK+23] and commonsense reasoning [ZHB+19, SBBC21] evaluations based on prior works [TLI+23, MWS+23] provide valuable insights, they often rely on multiple-choice formats, which are limited in scope. To comprehensively evaluate generated text, human evaluation remains essential, though it can be resource-intensive and sometimes lacks consistency, particularly when using crowd-sourcing. Recent studies [Tör23, LXA23, GRS+23, WA23] suggest that annotations from ChatGPT outperform those from Amazon crowd-sourced workers, highlighting the value of expert annotators in the evaluation process. Building on these findings, [PLH+23, CLL+23] used GPT-4 to replace crowd-sourced workers for comparing model outputs. In this approach, an evaluation prompt is presented, and both model outputs are provided for assessment in context.

7 Conclusion

We have introduced *Nanda*, a new state-of-the-art Hindi-English bilingual instruction-tuned large language model (LLM). It can perform a wide range of generative and downstream language tasks in both Hindi and English, ranging from common-sense reasoning to natural language understanding tasks such as sentiment analysis, irony detection, and hate speech detection. Its pre-trained and fine-tuned capabilities outperform all known open-source Hindi models of similar size, and are comparable to state-of-the-art open-source English models that were trained on larger datasets. We encourage researchers, hobbyists, and enterprise developers alike to experiment with and to develop on top of our model, particularly those working on multi-lingual and/or non-English applications.

Nanda represents an important evolution and expansion of the Hindi NLP and AI landscape. This Hindi model born in the UAE represents an important strategic step for government and commercial organizations towards the digital revolution. By advancing Hindi language understanding and generation, empowering local players with sovereign and private deployment options, and nurturing a vibrant ecosystem of applications and innovation, this work supports a broader strategic initiative of digital and AI transformation to usher in an open, more linguistically-inclusive, and culturally-aware era.

8 Release Notes

We release *Nanda* under Meta’s Llama-3 license, and users must adhere to the terms and conditions of the license,¹⁶ Meta’s acceptable use policy,¹⁷ Meta’s privacy policy,¹⁸ and the applicable policies, laws, and regulations governing the specific use-case and region. We encourage researchers, hobbyists, and enterprise devel-

¹⁶<https://www.llama.com/llama3/license/>

¹⁷<https://www.llama.com/llama3/use-policy/>

¹⁸<https://www.facebook.com/privacy/policy/>

opers alike to experiment with and to develop on top of the model – particularly those working on multi-lingual and/or non-English applications.

8.1 Intended Use

This model is one of the first of its kind in the Hindi LLM ecosystem, and has shown to be the best in the world among open Hindi or multilingual LLMs in terms of Hindi NLP capabilities. Some potential downstream uses are listed below:

- **Research:** This model can be used by researchers and developers to advance the Hindi LLM/NLP field.
- **Commercial Use:** It can be used as a foundational model to further fine-tune for specific usecases. Some potential usecases for businesses include (1) chat-assistants, (2) downstream tasks such as NLU/NLG, (3) customer service, and (4) process automation.

We believe that a number of audiences will benefit from our model:

- **Academics:** those researching Hindi natural language processing.
- **Businesses:** companies targeting Hindi-speaking audiences.
- **Developers:** those integrating Hindi language capabilities in apps.

8.2 Out-of-Scope Use

While *Nanda* is a powerful bilingual model catering to Hindi and English, it is essential to understand its limitations and the potential for its misuse. The following are some examples from the long list of scenarios where the model should not be used:

- **Malicious Use:** The model should not be used for generating harmful, misleading, or inappropriate content. This includes but is not limited to (i) generating or promoting hate speech, violence, or discrimination, (ii) spreading misinformation or fake news, (iii) engaging in illegal activities or promoting them, (i) (iv) handling sensitive information: the model should not be used to handle or to generate personal, confidential, or sensitive information.
- **Generalization Across All Languages:** *Nanda* is bilingual and optimized only for Hindi and English. It should not be assumed to have equal proficiency in other languages or dialects.
- **High-Stakes Decisions:** The model should not be used for making high-stakes decisions without human oversight. This includes medical, legal, financial, or safety-critical decisions, among others.

8.3 Biases, Risks, and Limitations

The model is trained on a mix of publicly available and proprietary data which in part was curated by our preprocessing pipeline. We used different techniques to reduce the bias that is inadvertently present in the dataset. While efforts were made to minimize biases, it is still possible that our model, like all LLM models, may exhibit some biases.

The model is trained as an AI assistant for Hindi and English speakers, and thus it should be used to help humans to boost their productivity. In this context, it is limited to produce responses for queries in these two languages and it might not produce appropriate responses for queries in other languages.

Potential misuses include generating harmful content, spreading misinformation, or handling sensitive information. Users are urged to use the model responsibly and with discretion.

9 Acknowledgments

We thank Zenan Zhai, Zhenxuan Zhang, Aili Shen, Hao Wang, and Yilin Geng for their help with developing the safety dataset.

References

- [AAA⁺23] Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malaric, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. Falcon-40B: an open large language model with state-of-the-art performance. Technical report, Technology Innovation Institute, 2023.
- [AAMK22] Jesujoba O. Alabi, David Ifeoluwa Adelani, Marius Mosbach, and Dietrich Klakow. Adapting pre-trained language models to african languages via multilingual adaptive fine-tuning, 2022.
- [ABF⁺19] N. Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George F. Foster, Colin Cherry, Wolfgang Macherey, Z. Chen, and Yonghui Wu. Massively multilingual neural machine translation in the wild: Findings and challenges. *ArXiv*, abs/1907.05019, 2019.
- [ADH⁺23] Kabir Ahuja, Harshita Diddee, Rishav Hada, Millicent Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Mohamed Ahmed, Kalika Bali, and Sunayana Sitaram. MEGA: Multilingual evaluation of generative AI. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4232–4267, Singapore, December 2023. Association for Computational Linguistics.
- [ADT⁺24] Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat Venkitesh, Madeline Smith, Kelly Marchisio, Sebastian Ruder, Acyr F. Locatelli, Julia Kreutzer, Nick Frosst, Phil Blunsom, Marzieh Fadaee, A. Ustun, and Sara Hooker. Aya 23: Open weight releases to further multilingual progress. *ArXiv*, abs/2405.15032, 2024.
- [AFT⁺24] Mehdi Ali, Michael Fromm, Klaudia Thellmann, Richard Rutmann, Max Lübbering, Johannes Leveling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper Buschhoff, Charvi Jain, Alexander Weber, Lena Jurkschat, Hammam Abdelwahab, Chelsea John, Pedro Ortiz Suarez, Malte Ostendorff, Samuel Weinbach, Rafet Sifa, Stefan Kesselheim, and Nicolas Flores-Herr. Tokenizer choice for LLM training: Negligible or crucial? In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3907–3924, Mexico City, Mexico, June 2024. Association for Computational Linguistics.
- [ALTdJ⁺23] Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebr’on, and Sumit K. Sanghai. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. *ArXiv*, abs/2305.13245, 2023.
- [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.
- [BBC⁺23] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report, 2023.
- [Bha24] BhabhaAI. Gajendra-v0.1. <https://huggingface.co/BhabhaAI/Gajendra-v0.1>, 2024. Accessed: 2024-10-29.
- [BMR⁺20] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.

- [BSA⁺23] Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models across training and scaling, 2023.
- [CCE⁺18a] 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. *ArXiv*, abs/1803.05457, 2018.
- [CCE⁺18b] 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. *arXiv preprint arXiv:1803.05457*, 2018.
- [CCG⁺23] Hyung Won Chung, Noah Constant, Xavier Garcia, Adam Roberts, Yi Tay, Sharan Narang, and Orhan Firat. Unimax: Fairer and more effective language sampling for large-scale multilingual pretraining, 2023.
- [CCSH24] Justin Cui, Wei-Lin Chiang, Ion Stoica, and Cho-Jui Hsieh. Or-bench: An over-refusal benchmark for large language models, 2024.
- [CHL⁺22] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.
- [CKG⁺20] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL*, pages 8440–8451, Online, 2020.
- [CLL⁺23] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing GPT-4 with 90%* ChatGPT quality, March 2023.
- [CLLM20] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA: pre-training text encoders as discriminators rather than generators. In *Proceedings of the 8th International Conference on Learning Representations, ICLR*, Addis Ababa, Ethiopia, 2020.
- [CND⁺22] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [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. *arXiv preprint arXiv:2107.03374*, 2021.
- [DAR⁺23] Sumanth Doddapaneni, Rahul Aralikkatte, Gowtham Ramesh, Shreya Goyal, Mitesh M. Khapra, Anoop Kunchukuttan, and Pratyush Kumar. Towards leaving no indic language behind: Building monolingual corpora, benchmark and models for indic languages, 2023.
- [DCLT19a] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019.
- [DCLT19b] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, NAACL-HLT, pages 4171–4186, Minneapolis, MN, USA, 2019.
- [DJP⁺24] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Cantón Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab A. AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriele Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guanglong Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Laurens Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Ju-Qing Jia, Kalyan Vasuden Alwala, K. Upasani, Kate Plawiak, Keqian Li, Ken-591 neth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhatta, Lauren Rantala-Young, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline C. Muzzi, Mahesh Babu Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melissa Hall Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri S. Chatterji, Olivier Duchenne, Onur cCelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasić, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shao-liang Nie, Sharan Narang, Sharath Chandra Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenxin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiao-

qing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yiqian Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zhengxu Yan, Zhengxing Chen, Zoe Papakipos, Aaditya K. Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adi Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Ben Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Shang-Wen Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzm'an, Frank J. Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory G. Sizov, Guangyi Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Han Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kaixing(Kai) Wu, U KamHou, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, A Laverder, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollár, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sung-Bae Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Andrei Poenaru, Vlad T. Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xia Tang, Xiaofang Wang, Xiaojuan Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models. *ArXiv*, abs/2407.21783, 2024.

[DLVNN⁺23] Viet Dac Lai, Chien Van Nguyen, Nghia Trung Ngo, Thuat Nguyen, Franck Dernoncourt,

- Ryan A Rossi, and Thien Huu Nguyen. Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback. *arXiv e-prints*, pages arXiv-2307, 2023.
- [DSK⁺22] Raj Dabre, Himani Shrotriya, Anoop Kunchukuttan, Ratish Puduppully, Mitesh Khapra, and Pratyush Kumar. Indicbart: A pre-trained model for indic natural language generation. In *Findings of the Association for Computational Linguistics: ACL 2022*. Association for Computational Linguistics, 2022.
- [FFG⁺24] Manuel Faysse, Patrick Fernandes, Nuno M. Guerreiro, António Loison, Duarte M. Alves, Caio Corro, Nicolas Boizard, João Alves, Ricardo Rei, Pedro H. Martins, Antoni Bigata Casademunt, François Yvon, André F. T. Martins, Gautier Viaud, Céline Hudelot, and Pierre Colombo. Croissantllm: A truly bilingual french-english language model, 2024.
- [GAS⁺23] Deep Ganguli, Amanda Askell, Nicholas Schiefer, Thomas I. Liao, Kamilė Lukošiušė, Anna Chen, Anna Goldie, Azalia Mirhoseini, Catherine Olsson, Danny Hernandez, Dawn Drain, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jackson Kernion, Jamie Kerr, Jared Mueller, Joshua Landau, Kamal Ndousse, Karina Nguyen, Liane Lovitt, Michael Sellitto, Nelson Elhage, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robert Lasenby, Robin Larson, Sam Ringer, Sandipan Kundu, Saurav Kadavath, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, Christopher Olah, Jack Clark, Samuel R. Bowman, and Jared Kaplan. The capacity for moral self-correction in large language models. *arXiv preprint arXiv:2302.07459*, 2023.
- [GCA⁺23] Jay Gala, Pranjal A. Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M. Khapra, Raj Dabre, and Anoop Kunchukuttan. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages, 2023.
- [Gen24a] GenVRAdmin. Aryabhatta-gemmaorca-merged. <https://huggingface.co/GenVRAdmin/AryaBhatta-GemmaOrca-Merged>, 2024. Accessed: 2024-10-29.
- [Gen24b] GenVRAdmin. Aryabhatta-gemmaultra-merged. <https://huggingface.co/GenVRAdmin/AryaBhatta-GemmaUltra-Merged>, 2024. Accessed: 2024-10-29.
- [GFAEP⁺22] Asier Gutiérrez-Fandiño, Jordi Armengol-Estapé, Marc Pàmies, Joan Llop-Palao, Joaquin Silveira-Ocampo, Casimiro Pio Carrino, Carme Armentano-Oller, Carlos Rodriguez-Penagos, Aitor Gonzalez-Agirre, and Marta Villegas. Maria: Spanish language models. *Procesamiento del Lenguaje Natural*, page 39–60, 2022.
- [GFZ⁺24] Yiduo Guo, Jie Fu, Huishuai Zhang, Dongyan Zhao, and Yikang Shen. Efficient Continual Pre-training by Mitigating the Stability Gap. *arXiv preprint arXiv:2406.14833*, 2024.
- [GJH⁺24a] Jay Gala, Thanmay Jayakumar, Jaavid Aktar Husain, Aswanth Kumar M, Mohammed Safi Ur Rahman Khan, Diptesh Kanojia, Ratish Puduppully, Mitesh M. Khapra, Raj Dabre, Rudra Murthy, and Anoop Kunchukuttan. Airavata: Introducing hindi instruction-tuned llm. *arXiv preprint arXiv: 2401.15006*, 2024.
- [GJH⁺24b] Jay Gala, Thanmay Jayakumar, Jaavid Aktar Husain, Aswanth Kumar M, Mohammed Safi Ur Rahman Khan, Diptesh Kanojia, Ratish Puduppully, Mitesh M. Khapra, Raj Dabre, Rudra Murthy, and Anoop Kunchukuttan. Airavata: Introducing hindi instruction-tuned llm. *arXiv preprint arXiv: 2401.15006*, 2024.
- [GJH⁺24c] Jay Gala, Thanmay Jayakumar, Jaavid Aktar Husain, Aswanth Kumar M, Mohammed Safi Ur Rahman Khan, Diptesh Kanojia, Ratish Puduppully, Mitesh M. Khapra, Raj Dabre, Rudra Murthy, and Anoop Kunchukuttan. Airavata: Introducing hindi instruction-tuned llm, 2024.
- [GRS⁺23] Mingqi Gao, Jie Ruan, Renliang Sun, Xunjian Yin, Shiping Yang, and Xiaojun Wan. Human-like summarization evaluation with ChatGPT. *arXiv preprint arXiv:2304.02554*, 2023.

- [GTB⁺21] Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation v0.0.1. <https://doi.org/10.5281/zenodo.5371628>, September 2021.
- [GXR⁺24] Gurpreet Gosal, Yishi Xu, Gokul Ramakrishnan, Rituraj Joshi, Avraham Sheinin, Zhiming Chen, Biswajit Mishra, Natalia Vassilieva, Joel Hestness, Neha Sengupta, Sunil Kumar Sahu, Bokang Jia, Onkar Pandit, Satheesh Katipomu, Samta Kamboj, Samujjwal Ghosh, Rahul Pal, Parvez Mullah, Soundar Doraiswamy, Mohamed El Karim Chami, and Preslav Nakov. Bilingual adaptation of monolingual foundation models, 2024.
- [HBB⁺20] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Xiaodong Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *ArXiv*, abs/2009.03300, 2020.
- [HBB⁺22] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2022.
- [ILK⁺23] Ayyoob ImaniGooghari, Peiqin Lin, Amir Hossein Kargaran, Silvia Severini, Masoud Jalili Sabet, Nora Kassner, Chunlan Ma, Helmut Schmid, André Martins, François Yvon, and Hinrich Schütze. Glot500: Scaling multilingual corpora and language models to 500 languages. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, page 1082–1117. Association for Computational Linguistics, 2023.
- [JOOA⁺22] Odunayo Jude Ogundepo, Akintunde Oladipo, Mofetoluwa Adeyemi, Kelechi Ogueji, and Jimmy Lin. AfriTeVA: Extending ‘small data’ pretraining approaches to sequence-to-sequence models. In Colin Cherry, Angela Fan, George Foster, Gholamreza (Reza) Haffari, Shahram Khadivi, Nanyun (Violet) Peng, Xiang Ren, Ehsan Shareghi, and Swabha Swayamdipta, editors, *Proceedings of the Third Workshop on Deep Learning for Low-Resource Natural Language Processing*, pages 126–135, Hybrid, July 2022. Association for Computational Linguistics.
- [JSB⁺20] Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. The state and fate of linguistic diversity and inclusion in the NLP world. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online, July 2020. Association for Computational Linguistics.
- [JSK⁺24] Raviraj Joshi, Kanishk Singla, Anusha Kamath, Raunak Kalani, Rakesh Paul, Utkarsh Vaidya, Sanjay Singh Chauhan, Niranjana Wartikar, and Eileen Long. Adapting multilingual llms to low-resource languages using continued pre-training and synthetic corpus, 2024.
- [K24] Adithya S K. indic_eval. https://github.com/adithya-s-k/indic_eval, 2024. Accessed: 2024-10-29.
- [KBM⁺21] Simran Khanuja, Diksha Bansal, Sarvesh Mehtani, Savya Khosla, Atreyee Dey, Balaji Gopalan, Dilip Kumar Margam, Pooja Aggarwal, Rajiv Teja Nagipogu, Shachi Dave, Shruti Gupta, Subhash Chandra Bose Gali, Vish Subramanian, and Partha Talukdar. MuriL: Multilingual representations for indian languages, 2021.
- [KBR⁺24] Aditya Kusupati, Gantavya Bhatt, Aniket Rege, Matthew Wallingford, Aditya Sinha, Vivek Ramanujan, William Howard-Snyder, Kaifeng Chen, Sham Kakade, Prateek Jain, and Ali Farhadi. Matryoshka representation learning, 2024.
- [KKG⁺20] Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, Gokul N.C., Avik Bhattacharyya, Mitesh M. Khapra, and Pratyush Kumar. IndicNLP Suite: Monolingual corpora, evaluation benchmarks and pre-trained multilingual language models for Indian languages. In Trevor Cohn, Yulan He, and Yang Liu, editors, *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4948–4961, Online, November 2020. Association for Computational Linguistics.

- [KRLB20] Fajri Koto, Afshin Rahimi, Jey Han Lau, and Timothy Baldwin. Indolem and indobert: A benchmark dataset and pre-trained language model for indonesian nlp, 2020.
- [KSS⁺22] Aman Kumar, Himani Shrotriya, Prachi Sahu, Raj Dabre, Ratish Puduppully, Anoop Kunchukuttan, Amogh Mishra, Mitesh M. Khapra, and Pratyush Kumar. Indicnlg benchmark: Multilingual datasets for diverse nlg tasks in indic languages, 2022.
- [KYR⁺23] Hyunwoong Ko, Kichang Yang, Minho Ryu, Taekyoon Choi, Seungmu Yang, Jiwung Hyun, Sungho Park, and Kyubyong Park. A technical report for polyglot-ko: Open-source large-scale korean language models, 2023.
- [LH18] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Proceedings of the International Conference on Learning Representations*, ICLR, Vancouver, VC, Canada, 2018.
- [LHE21] Stephanie C. Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. In *Annual Meeting of the Association for Computational Linguistics*, 2021.
- [LJT⁺24] Peiqin Lin, Shaoxiong Ji, Jörg Tiedemann, André F. T. Martins, and Hinrich Schütze. Mala-500: Massive language adaptation of large language models, 2024.
- [LKL⁺23] Risto Luukkonen, Ville Komulainen, Jouni Luoma, Anni Eskelinen, Jenna Kanerva, Hanna-Mari Kupari, Filip Ginter, Veronika Laippala, Niklas Muennighoff, Aleksandra Piktus, Thomas Wang, Nouamane Tazi, Teven Le Scao, Thomas Wolf, Osmo Suominen, Samuli Sairanen, Mikko Merioksa, Jyrki Heinonen, Aija Vahtola, Samuel Antao, and Sampo Pyysalo. Fingpt: Large generative models for a small language, 2023.
- [LKW⁺23] Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. Bactrian-X: A multilingual replicable instruction-following model with low-rank adaptation. *arXiv preprint arXiv:2305.15011*, 2023.
- [LLG⁺19] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- [LMA⁺22] Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona Diab, Veselin Stoyanov, and Xian Li. Few-shot learning with multilingual language models, 2022.
- [LMZ⁺24] Lizhi Lin, Honglin Mu, Zenan Zhai, Minghan Wang, Yuxia Wang, Renxi Wang, Junjie Gao, Yixuan Zhang, Wanxiang Che, Timothy Baldwin, Xudong Han, and Haonan Li. Against the achilles’ heel: A survey on red teaming for generative models, 2024.
- [LOG⁺19] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT pre-training approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [LQN⁺23] Zhengzhong Liu, Aurick Qiao, Willie Neiswanger, Hongyi Wang, Bowen Tan, Tianhua Tao, Junbo Li, Yuqi Wang, Suqi Sun, Omkar Pangarkar, Richard Fan, Yi Gu, Victor Miller, Yonghao Zhuang, Guowei He, Haonan Li, Fajri Koto, Liping Tang, Nikhil Ranjan, Zhiqiang Shen, Xuguang Ren, Roberto Iriando, Cun Mu, Zhiting Hu, Mark Schulze, Preslav Nakov, Tim Baldwin, and Eric P. Xing. Llm360: Towards fully transparent open-source llms, 2023.
- [LTW⁺24] Zhengzhong Liu, Bowen Tan, Hongyi Wang, Willie Neiswanger, Tianhua Tao, Haonan Li, Fajri Koto, Yuqi Wang, Suqi Sun, Omkar Pangarkar, Richard Fan, Yi Gu, Victor Miller, Liqun Ma, Liping Tang, Nikhil Ranjan, Yonghao Zhuang, Guowei He, Renxi Wang, Mingkai Deng, Robin Algayres, Yuanzhi Li, Zhiqiang Shen, Preslav Nakov, and Eric Xing. Llm360 k2-65b: Scaling up fully transparent open-source llms. 2024.

- [LXA23] Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. ChatGPT as a factual inconsistency evaluator for abstractive text summarization. *arXiv preprint arXiv:2303.15621*, 2023.
- [LZK⁺23] Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. CMMLU: Measuring massive multitask language understanding in Chinese. *arXiv preprint arXiv: 2306.09212*, 2023.
- [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. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, EMNLP, pages 2381–2391, Brussels, Belgium, 2018.
- [MGU⁺22] Fatemehsadat Mireshghallah, Kartik Goyal, Archit Uniyal, Taylor Berg-Kirkpatrick, and Reza Shokri. Quantifying privacy risks of masked language models using membership inference attacks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, EMNLP, pages 8332–8347, Abu Dhabi, United Arab Emirates, 2022.
- [MHSB21] Stuart Mesham, Luc Hayward, Jared Shapiro, and Jan Buys. Low-resource language modelling of south african languages, 2021.
- [MPR22] Benjamin Minixhofer, Fabian Paischer, and Navid Rekabsaz. WECHSEL: Effective initialization of subword embeddings for cross-lingual transfer of monolingual language models. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3992–4006, Seattle, United States, July 2022. Association for Computational Linguistics.
- [MTG⁺23] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-Refine: Iterative refinement with self-feedback. *arXiv preprint arXiv:2303.17651*, 2023.
- [MWS⁺23] Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*, 2023.
- [NZL⁺24] Xuan-Phi Nguyen, Wenxuan Zhang, Xin Li, Mahani Aljunied, Zhiqiang Hu, Chenhui Shen, Yew Ken Chia, Xingxuan Li, Jianyu Wang, Qingyu Tan, Liying Cheng, Guanzheng Chen, Yue Deng, Sen Yang, Chaoqun Liu, Hang Zhang, and Lidong Bing. Seallms – large language models for southeast asia, 2024.
- [OAA⁺23] Akintunde Oladipo, Mofetoluwa Adeyemi, Orevaoghene Ahia, Abraham Owodunni, Odunayo Ogundepo, David Adelani, and Jimmy Lin. Better quality pre-training data and t5 models for African languages. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 158–168, Singapore, December 2023. Association for Computational Linguistics.
- [Ope23] OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [OWJ⁺22] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- [OZL21] Kelechi Ogueji, Yuxin Zhu, and Jimmy Lin. Small data? no problem! exploring the viability of pretrained multilingual language models for low-resourced languages. In Duygu Ataman, Alexandra Birch, Alexis Conneau, Orhan Firat, Sebastian Ruder, and Gozde Gul Sahin, editors, *Proceedings of the 1st Workshop on Multilingual Representation Learning*, pages 116–126, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.

- [PGL⁺22] Jonas Pfeiffer, Naman Goyal, Xi Victoria Lin, Xian Li, James Cross, Sebastian Riedel, and Mikel Artetxe. Lifting the curse of multilinguality by pre-training modular transformers. In *North American Chapter of the Association for Computational Linguistics*, 2022.
- [PLH⁺23] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with GPT-4. *arXiv preprint arXiv:2304.03277*, 2023.
- [PLMTB23] Aleksandar Petrov, Emanuele La Malfa, Philip HS Torr, and Adel Bibi. Language model tokenizers introduce unfairness between languages. *arXiv preprint arXiv:2305.15425*, 2023.
- [PTNT22] Long Phan, Hieu Tran, Hieu Nguyen, and Trieu H. Trinh. Vit5: Pretrained text-to-text transformer for vietnamese language generation, 2022.
- [RBC⁺22] Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training Gopher. *arXiv preprint arXiv:2112.11446*, 2022.
- [RPS⁺24] Gemma Team Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, L’eonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ram’e, Johan Ferret, Peter Liu, Pouya Dehghani Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stańczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Boxi Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Christopher A. Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, D. Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost R. van Amersfoort, Josh Gordon, Josh Lipschultz, Joshua Newlan, Junsong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, L. Sifre, L. Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Gorner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, S. Mc Carthy, Sarah Perrin, S’ebastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomás Kociský, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei,

- Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Brian Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeffrey Dean, Demis Hassabis, Koray Kavukcuoglu, Clément Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Keanealy, Robert Dadashi, and Alek Andreev. Gemma 2: Improving open language models at a practical size. *ArXiv*, abs/2408.00118, 2024.
- [RPV⁺21] Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych. How Good is Your Tokenizer? On the Monolingual Performance of Multilingual Language Models. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3118–3135, Online, August 2021. Association for Computational Linguistics.
- [RSR⁺20] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [RSR⁺23] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2023.
- [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.
- [Sar24] Sarvamai. sarvam-2b-v0.5. <https://huggingface.co/sarvamai/sarvam-2b-v0.5>, 2024. Accessed: 2024-10-29.
- [SBBC21] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. WinoGrande: An adversarial Winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- [SFA⁺23] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, Dragomir Radev, Eduardo González Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady Elsahar, Hamza Benyamina, Hieu Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jörg Frohberg, Joseph Tobing, Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon Weber, Long Phan, Loubna Ben allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario Šaško, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nurulaqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rhea Harliman, Rishi Bommasani, Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, Somaieh Nikpoor, Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heizerling, Chenglei Si, Davut Emre Taşar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza

Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Fevry, Trishala Neeraj, Urmish Thakker, Vikas Raut, Xiangru Tang, Zheng-Xin Yong, Zhiqing Sun, Shaked Brody, Yallow Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre François Lavallée, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Barua, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramanian, Aurélie Névél, Charles Lovering, Dan Garrette, Deepak Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Jordan Clive, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najoung Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, Shani Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan Belinkov, Zachary Bamberger, Zdeněk Kasner, Alice Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ana Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh Behroozi, Benjamin Ajibade, Bharat Saxena, Carlos Muñoz Ferrandis, Daniel McDuff, Danish Contractor, David Lansky, Davis David, Douwe Kiela, Duong A. Nguyen, Edward Tan, Emi Baylor, Ezinwanne Ozoani, Fatima Mirza, Frankline Ononiwu, Habib Rezanejad, Hessie Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jesse Passmore, Josh Seltzer, Julio Bonis Sanz, Livia Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akinlolu, Michael McKenna, Mike Qiu, Muhammed Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nour Fahmy, Olanrewaju Samuel, Ran An, Rasmus Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas Wang, Sourav Roy, Sylvain Viguier, Thanh Le, Tobi Oyeade, Trieu Le, Yoyo Yang, Zach Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, Antonio Miranda-Escalada, Ayush Singh, Benjamin Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourrier, Daniel León Perifán, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrmann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Imane Bello, Ishani Dash, Ji Hyun Kang, John Giorgi, Jonas Golde, Jose David Posada, Karthik Rangasai Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, Maria A Castillo, Marianna Nezhurina, Mario Sängler, Matthias Samwald, Michael Cullan, Michael Weinberg, Michiel De Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patrick Haller, Ramya Chandrasekhar, Renata Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, Sinee Sang-aaronsiri, Sriшти Kumar, Stefan Schweter, Sushil Bharati, Tanmay Laud, Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yash Shailesh Bajaj, Yash Venkatraman, Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and Thomas Wolf. BLOOM: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2023.

- [SFT⁺23] Oleh Shliazhko, Alena Fenogenova, Maria Tikhonova, Vladislav Mikhailov, Anastasia Kozlova, and Tatiana Shavrina. mgpt: Few-shot learners go multilingual, 2023.
- [SGB⁺24] Harman Singh, Nitish Gupta, Shikhar Bharadwaj, Dinesh Tewari, and Partha Talukdar. Indicgenbench: A multilingual benchmark to evaluate generation capabilities of llms on indic languages, 2024.
- [SHB16] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th An-*

nual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics.

- [SLP⁺21] Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *ArXiv*, abs/2104.09864, 2021.
- [SSJ⁺23] Neha Sengupta, Sunil Kumar Sahu, Bokang Jia, Satheesh Katipomu, Haonan Li, Fajri Koto, William Marshall, Gurpreet Gosal, Cynthia Liu, Zhiming Chen, Osama Mohammed Afzal, Samta Kamboj, Onkar Pandit, Rahul Pal, Lalit Pradhan, Zain Muhammad Mujahid, Massa Baali, Xudong Han, Sondos Mahmoud Bsharat, Alham Fikri Aji, Zhiqiang Shen, Zhengzhong Liu, Natalia Vassilieva, Joel Hestness, Andy Hock, Andrew Feldman, Jonathan Lee, Andrew Jackson, Hector Xuguang Ren, Preslav Nakov, Timothy Baldwin, and Eric Xing. Jais and jais-chat: Arabic-centric foundation and instruction-tuned open generative large language models, 2023.
- [SXD⁺22] Hao Sun, Guangxuan Xu, Jiawen Deng, Jiale Cheng, Chujie Zheng, Hao Zhou, Nanyun Peng, Xiaoyan Zhu, and Minlie Huang. On the safety of conversational models: Taxonomy, dataset, and benchmark. In *Findings of the Association for Computational Linguistics, ACL*, pages 3906–3923, Dublin, Ireland, 2022.
- [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.
- [TMS⁺23] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [Tör23] Petter Törnberg. ChatGPT-4 outperforms experts and crowd workers in annotating political twitter messages with zero-shot learning. *arXiv preprint arXiv:2304.06588*, 2023.
- [TRP⁺24] Liping Tang, Nikhil Ranjan, Omkar Pangarkar, Xuezhi Liang, Zhen Wang, Li An, Bhaskar Rao, Linghao Jin, Huijuan Wang, Zhoujun Cheng, Suqi Sun, Cun Mu, Victor Miller, Xuezhe Ma, Yue Peng, Zhengzhong Liu, and Eric P. Xing. Txt360: A top-quality llm pre-training dataset requires the perfect blend, 2024.
- [UAY⁺24] A. Ustun, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blunsom, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. Aya model: An instruction finetuned open-access multilingual language model. *ArXiv*, abs/2402.07827, 2024.
- [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems*, pages 5998–6008, Long Beach, CA, USA, 2017.
- [WA23] Minghao Wu and Alham Fikri Aji. Style Over Substance: Evaluation biases for large language models. *arXiv preprint arXiv:2307.03025*, 2023.

- [WGG⁺24] Chengyue Wu, Yukang Gan, Yixiao Ge, Zeyu Lu, Jiahao Wang, Ye Feng, Ying Shan, and Ping Luo. LLaMA pro: Progressive LLaMA with block expansion. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6518–6537, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [WGY⁺24] Ishaan Watts, Varun Gumma, Aditya Yadavalli, Vivek Seshadri, Manohar Swaminathan, and Sunayana Sitaram. Pariksha: A large-scale investigation of human-llm evaluator agreement on multilingual and multi-cultural data, 2024.
- [WLH⁺23] Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. Do-Not-Answer: A dataset for evaluating safeguards in LLMs. *arXiv preprint arXiv:2308.13387*, 2023.
- [WWS⁺22] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-Thought prompting elicits reasoning in large language models. NeurIPS, New Orleans, LA, USA, 2022.
- [XCR⁺20] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mT5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv:2010.11934*, 2020.
- [YRC23] Oleksandr Yermilov, Vipul Raheja, and Artem Chernodub. Privacy- and utility-preserving NLP with anonymized data: A case study of pseudonymization. In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing*, TrustNLP, pages 232–241, Toronto, ON, Canada, 2023.
- [YSM⁺23] Zheng-Xin Yong, Hailey Schoelkopf, Niklas Muennighoff, Alham Fikri Aji, David Ifeoluwa Adelani, Khalid Almubarak, M Saiful Bari, Lintang Sutawika, Jungo Kasai, Ahmed Barua, Genta Indra Winata, Stella Biderman, Edward Raff, Dragomir Radev, and Vassilina Nikoulina. Bloom+1: Adding language support to bloom for zero-shot prompting, 2023.
- [ZHB⁺19] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Annual Meeting of the Association for Computational Linguistics*, 2019.
- [ZLD⁺23] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. GLM-130B: an open bilingual pre-trained model. In *Proceedings of the Eleventh International Conference on Learning Representations*, ICLR, Kigali, Rwanda, 2023.
- [ZMH⁺23] Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. In *Proceedings of the Eleventh International Conference on Learning Representations*, ICLR, Kigali, Rwanda, 2023.
- [ZRS⁺21] Wei Zeng, Xiaozhe Ren, Teng Su, Hui Wang, Yi Liao, Zhiwei Wang, Xin Jiang, ZhenZhang Yang, Kaisheng Wang, Xiaoda Zhang, Chen Li, Ziyang Gong, Yifan Yao, Xinjing Huang, Jun Wang, Jianfeng Yu, Qi Guo, Yue Yu, Yan Zhang, Jin Wang, Hengtao Tao, Dasen Yan, Zexuan Yi, Fang Peng, Fangqing Jiang, Han Zhang, Lingfeng Deng, Yehong Zhang, Zhe Lin, Chao Zhang, Shaojie Zhang, Mingyue Guo, Shanzhi Gu, Gaojun Fan, Yaowei Wang, Xuefeng Jin, Qun Liu, and Yonghong Tian. Pangu- α : Large-scale autoregressive pretrained chinese language models with auto-parallel computation, 2021.