

PARAMANU: A FAMILY OF NOVEL EFFICIENT GENERATIVE FOUNDATION LANGUAGE MODELS FOR INDIAN LANGUAGES

Anonymous authors

Paper under double-blind review

ABSTRACT

We present PARAMANU (which means “atom” in multiple Indian languages), a *family* of novel language models for *Indian languages*. It is a collection of auto-regressive monolingual, bilingual, and multilingual Indian language models pretrained from scratch, currently covering 10 Indian languages (Assamese, Bangla, Hindi, Konkani, Maithili, Marathi, Odia, Sanskrit, Tamil, Telugu) across 5 scripts (Bangla, Devanagari, Odia, Tamil, Telugu). The models are pretrained with a context size of 1024 on a single GPU, and are of varying sizes ranging from 13.29 M to 367.5 M parameters. We proposed a RoPE embedding scaling method that enables us to pretrain language models from scratch at larger sequence length context size **on single GPU without increased GPU memory. ~~than the equivalent GPU memory.~~** We have also developed an efficient and advanced novel tokenizer **with least fertility score among existing LLMs** for Indian languages using a combination of BPE and Unigram that can also tokenize unseen languages written in the same script or the Roman script. We also proposed language specific tokenization for multilingual models and domain specific tokenization for monolingual language models. In order to avoid the “curse of multi-linguality” in our multilingual MPARAMANU model, we pretrained on comparable corpora by typological grouping using the same script. We proposed and performed pretraining for more than 1 epoch of training for most of our language models. From our results, we observed the language transfer phenomenon from low resource to high resource within languages of the same script and typology. We performed human evaluation of our pretrained models for open end text generation on grammar, coherence, creativity, and factuality metrics for several languages. Our Paramanu models outperformed standard **and multilingual** large language models (LLMs) by a large margin in performance despite being smaller in size by 64 to 20 times. We studied the impact of language specific tokenization versus language agnostic tokenization for bilingual language modeling. We also studied the impact of BPE versus Unigram tokenization for Devanagari script languages. We further created instruction-tuning datasets and instruction-tuned our pretrained models on 23,000 instructions in respective languages **except Hindi, for which we used 75,000 instructions.** Comparison with multilingual LLMs on various commonsense reasoning benchmarks for natural language understanding, natural language inference, and machine reading comprehension shows the advantage of our models. The performance of our Paramanu models leads to the conclusion that high quality generative language models are possible without high amount of compute power (FLOPS) and enormous number of parameters.

1 INTRODUCTION

World’s new age growth arenas are non-English speakers from countries such as India, China, Indonesia, Mexico, South Africa, etc., where more people will connect to the Internet for information need without language acting as boundaries. Multi-lingualism represents freedom of expression and diversity in a country like India. According to a Guardian article¹ by Andras Kornai, “95% of

¹<http://labs.theguardian.com/digital-language-divide/>

all languages in use today will never gain traction online”; this highlighting the digital language divide. Despite around 7,000 languages, current NLP technologies cover only 12%, leaving many non-English and non-European languages underserved. This, in our personal opinion, is an extremely incapacitating bias and language divide for the global digital world where the non-English, non-European language world has been left out. According to the Indian Census 2011, there are 22 official languages and more than hundred others with a sizeable number of speakers in India. Both Hindi and Bangla (Bengali), despite being the world’s 5th and 6th most spoken languages respectively according to Babel² are still underrepresented in today’s NLP technology. Multilingual NLP faces the challenges of having lack of quality benchmark datasets covering diverse languages from different language families and especially under-represented languages, and is typically referred as “low-resource” in the NLP community. Large language models (LLMs) such as GPTNeoX (Black et al., 2022), OPT (Zhang et al., 2022), LLaMa (Touvron et al., 2023), PaLM (Chowdhery et al., 2023), GPT-2-XL (Radford et al., 2019), GPT-J (Wang & Komatsuzaki, 2021), etc. have primarily focused on English and mostly European languages whereas other languages have not been given priority. Bloom (Workshop et al., 2023) is considered to be the biggest multilingual auto-regressive model that has been built till now; it has been pretrained on 45 languages including Indian languages. However, Indian languages are morphologically richer and typologically distinct than languages written in the Latin script [and, hence](#), grouping them together without considering linguistics nuances often leads to poor performance of LLMs for low resource languages.

This work is an attempt to make language technology more accessible for Indian languages. In this work, we focus on 10 Indian languages (Assamese, Bangla/Bengali, Hindi, Konkani, Maithili, Marathi, Odia, Sanskrit, Tamil, Telugu) written in 5 distinct scripts (Assamese-Bengali, Devanagari, Odia, Tamil, Telugu) comprising of more than billion speakers in the global world. [Our goal is to show that generative language models for low-resource Indian languages can be trained from scratch with limited compute and token budget, without using English corpora. We excluded English due to its linguistic differences with Indian languages in terms of typology, script, morphology and grammar, aiming to maintain language purity and typology grouping. Additionally, existing multilingual models often have an English-centric bias due to the large imbalance in data. To the best of our knowledge, the existing LLMs and multilingual LLMs struggle to generate grammatically correct and coherent sentences in Indian languages \(and reasoning ability is too far\) as shown in Appendix Tables 20, 26, 27, 28, 29, 30, 30, 32, 33, 35, 36, 37, 38, 39, 40. Many LLMs even generate text in even Arabic \(Table 29\) or Japanese \(Table 39\) scripts when prompted with prompts in Bengali, Hindi, and Sanskrit. We also found that LLMs including ChatGPT \(Oct’23\) were not able to distinguish languages of the same script such as Bengali versus Assamese. Human evaluation by annotators confirm this \(Section A.5 in Appendix\). Our models, based on Transformer decoders \(Vaswani et al., 2017\), are enhanced with improvements in the architecture. These enhancements, improvements and novelties make our models efficient, small but strong. They have been pretrained from scratch to support context size of 1024 without requiring higher physical memory on a single NVIDIA A100-PCIE-40GB GPU. Our model architecture has the ability to capture a much higher context size without requiring equivalent physical memory; ~~in contrast, the LLaMa-1 model was pretrained from scratch using 2048 GPUs for a context size of 2048.~~](#) Our models are of varying sizes and of three types: (1) monolingual language models, (2) bilingual language models, and (3) multilingual language models. Both bilingual and multilingual language models are pretrained from scratch on comparable corpora and with typological grouping of languages to avoid curse of multilinguality.

We summarize our contributions as follows.

1. We propose a RoPE (Su et al., 2022) embedding scaling method that enables us to pretrain language models from scratch at larger sequence length context size than the equivalent GPU memory. We scaled the RoPE embedding through a shrinking factor by dividing the target context length y by `max_permissible_context_size_length` on single GPU keeping all other hyperparameters fixed such as batch size, vocabulary size.
2. We propose a novel tokenization method using a combination of both Byte-Pair Encoding (BPE) (Sennrich et al., 2016) and Unigram (Kudo, 2018) algorithms.
3. We perform language specific tokenization for multilingual model and domain specific tokenization for monolingual language modeling.

²<https://www.babbel.com/en/magazine/the-10-most-spoken-languages-in-the-world>

4. We perform pretraining on comparable corpora for multilingual/bilingual generative language model to handle data imbalance and curse-of-multi-linguality in multilingual language model.
5. We perform more than 1 epoch of pretraining for most our language models. From our results, we observed the language transfer from low resource to high resource within language of the same script and typology. We found that small models trained for more than 1 epoch on high quality data is better than bigger model trained for 1 epoch on not so good quality bigger corpora. We studied the impact of language specific tokenization versus language agnostic tokenization for bilingual language modeling. We also studied the impact of BPE vs Unigram tokenization for Devanagari script languages as a proof-of-concept.
6. We curated books for pretraining dataset and created an instruction-tuning dataset of 23000 instructions and instruction-tune our pretrained Bangla, ~~Hindi~~, Marathi, Tamil, and Telugu models using 23000 instructions each in respective languages ~~but for Hindi, we used additional 52,000 Alpaca translated Hindi (Taori et al., 2023), total of 75,000 instructions.~~

2 RELATED WORK

The rise of large language models (LLMs) has led to numerous new models competing across various benchmarks, but most remain English-centric. Bloom, the largest multilingual LLM with 170 billion parameters, utilizes 40% English data and 60% from 44 other languages. Despite extensive multilingual training, many LLMs exhibit a significant English bias. Research has shown that a substantial portion of neurons in these models remains inactive; for instance, in the 66 billion parameter OPT model, over 70% of the feed-forward network (FFN) neurons in certain layers are “dead” meaning they do not activate even on diverse datasets. This sparsity in neuron activation limits the model’s ability to learn and generalize across languages, particularly those that differ significantly from English (Voita et al., 2024). Consequently, even when trained on multilingual data, these models struggle to effectively process and generate text in languages beyond English, revealing a critical limitation in their design. Recent efforts to adapt English-centric models like Llama for Indian languages (Airavata (Gala et al., 2024), OpenHathi for Hindi, and Tamil) have involved extending vocabulary and fine-tuning techniques such as LoRA (Hu et al., 2022) and QLoRA (Dettmers et al., 2023). However, these models still exhibit a strong English bias and struggle to generate high-quality text in Indian languages. Massively multilingual models (MMTs) (Devlin et al., 2019), (Conneau et al., 2020), (Xue et al., 2021) are pretrained on large corpora but often lack alignment between languages, leading to poor transfer performance for distant languages (Lauscher et al., 2020). The “curse of multilinguality” (CoM) indicates that adding more languages can degrade per-language performance, necessitating larger model capacities and corpora. While language-specific adapters like MAD-X (Pfeiffer et al., 2020) improve performance, they do not generalize well to unseen languages. Recent advancements, such as MAD-G (Ansell et al., 2021) and BAD-X (Parović et al., 2022), focus on bilingual adapters to enhance language transfer. Our work aims to create dedicated models for low-resource Indian languages by developing language-specific generative models from scratch, emphasizing linguistic features, typology, and tailored tokenization.

3 METHODOLOGY

3.1 DATASET FOR PRETRAINING

Pretrained data was split into 95%-5% training and validation sets so that we do not lose much data for pretraining as the purpose of this work is a step towards developing pretrained generative effective language models from scratch using our novel architecture for Indian languages. Pretraining data covers web scrapped news, blogs from IndicCorp v2 (Doddapaneni et al., 2023), Bangla literature from Vacaspati (Bhattacharyya et al., 2023), Wikipedia articles, curated books of various genres, subjects, education books, magazines in respective Indian languages representing each distinct language community to cover Indian culture, rich history, and knowledge. Our pretraining corpora have no source code, scientific journals/articles, medical and engineering education books, research papers as these are generally written in English in India. Dataset details can be found in Table 19 in Appendix and [data distribution can be found in the Figure 3 in Appendix](#). We apply the following data cleaning and preprocessing techniques as mentioned in Appendix A.4.

| Models | Size | #Tokens pretrained | # Training A100 hours | Script & Family | # Speakers |
|--------------------------|------|--------------------|-----------------------|-------------------------------------|------------|
| Bloom | 7.1B | 340 B | 1.08M | Multilingual | 1B+ |
| OpenHathi (Llama) | 7B | 1 T | N/A | Hindi (Indo-European) | 692M |
| Sarvam | 2B | 4 T | N/A | Multilingual (Devanagri, Dravidian) | 1B+ |
| Paramanu-Bangla (ours) | 108M | 26.21 B | 19.45 | Bengali (Indo-European) | 300M |
| Paramanu-Hindi (ours) | 367M | 66 B | 239 | Devanagari (Indo-European) | 692M |
| Paramanu-Marathi (ours) | 208M | 28.83 B | 88 | Devanagari (Indo-European) | 99M |
| Paramanu-Odia (ours) | 87M | 52.42 B | 84.5 | Odia (Indo-European) | 43M |
| Paramanu-Sanskrit (ours) | 139M | 45 B | 110 | Devanagari (Indo-European) | 0.025M |
| Paramanu-Tamil (ours) | 208M | 26.2 B | 208 | Tamil (Indo-Dravidian) | 77M |
| Paramanu-Telugu (ours) | 208M | 39.32 B | 112.5 | Telugu (Indo-Dravidian) | 95M |
| mParamanu (ours) | 162M | 26.2 B | 118 | Multilingual (Devanagari) | 1B+ |

Table 1: Pretrained LLM tokens, training hours, [script](#), [language family](#) and [#speakers](#).

3.2 DATASET FOR INSTRUCTION TUNING

We created a dataset of 5,000 human-annotated instructions covering tasks like poem and novel writing, article summarization, grammar correction, and Q&A on topics such as climate change and healthcare in Bangla, reflecting Bengali culture and linguistics. We then used the Google Translate API (goo, 2023) to translate these instructions into Hindi, Marathi, Tamil, and Telugu. Additionally, we translated 15,000 instructions from Dolly (Conover et al., 2023) to these languages and generated 3,000 instructions using the self-instruct technique (Wang et al., 2023) in five Indian languages. Finally, we fine-tuned our pretrained models [except Hindi](#) on a total of 23,000 instructions, including human-annotated, machine-generated, and translated instructions. [Since creating multilingual datasets using automatic machine translation from Bengali to other Indian languages introduce translation errors but we also performed human checks and corrected them. We found around 8% word errors on average for Hindi and Marathi using Google Translate.](#)

3.3 TOKENIZATION

Figure 1 shows the flowchart of our novel tokenization technique both for monolingual and multilingual settings. We performed domain adaptive tokenization for monolingual models using a combination of Byte-Pair encoding (BPE) (Sennrich et al., 2016) and Unigram (Kudo, 2018). We trained separate Byte-Pair encoding (BPE) (Sennrich et al., 2016) and Unigram tokenizers using Sentencepiece (Kudo & Richardson, 2018) module on the high quality part of the pretraining data from scratch. Then, we merge both the [independent tokenizers by intersection respectively of size \$V'\$](#) . This is similar to merging two different data structures of same size (considering each tokenizer as a list of tokens) by intersection to remove overlapping elements. We used the merged tokenizer to tokenize the pretraining data. During pre-tokenization, NFC normalization was performed on the processed data; digits are split into individual tokens and unknown UTF-8 characters were reduced to byte granularity. For monolingual tokenization, we trained individual BPE and Unigram tokenizers on domain specific data for respective language with same vocabulary size and merged the tokenizers via merge by intersection to remove overlapping tokens and make our specialised tokenizer compact, optimized, effective and highly effective for monolingual data and also performed the same approach to tokenize the multilingual pretraining corpora by performing language specific tokenization in multi-task way where every task is a language and merged the distinct language specific tokenizers by union via intersection, i.e, removing the overlapping tokens. In this way, our mBharat tokenizer was able to learn language specific tokens based on typology and can generalize to unseen languages of the same script. mBharat tokenizer was also exposed to little amount of English high quality corpus to learn its ability to tokenize languages of the Roman script such as English. From 2, we observe that mBharat tokenizer has the least fertility score of 1.66 for languages in Assamese-Bengali, 1.25 for Devanagari (Hindi) and 1.75 for Odia script among tokenizers of LLMs like Sarvam 2B (Sarvam2B, 2024), LLama-3.1 (Dubey et al., 2024), Gemma-2 (Team et al., 2024), and GPT-4o.

3.4 MODEL ARCHITECTURE

Multilingual mParamanu, monolingual Paramanu models (Assamese, Bangla, Hindi, Konkani, Marathi, Odia, Tamil, Telugu, and Paramanu-Sanskrit) and bilingual Konkani-Maithili models are based on transformer (Vaswani et al., 2017) based causal decoder architecture (Radford et al.,

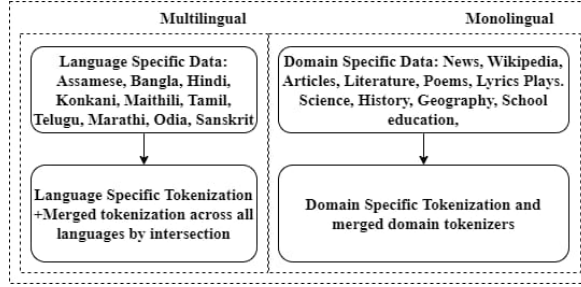


Figure 1: Tokenization technique for monolingual and multilingual setting.

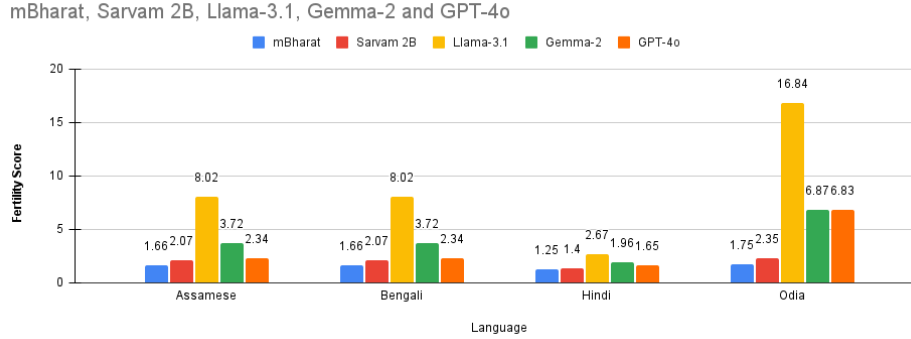


Figure 2: Fertility score of mBharat tokenizer v/s LLMs across languages of 3 scripts (Bengali, Devanagari, and Odia). LLMs score are reported from (Sarvam2B, 2024)

2019) with modifications. The model architecture uses RMSNorm as pre-normalization layer with $\text{norm_epsilon} = 1e-5$, approximate GeGLU Activations (Shazeer, 2020) as non-linearity by replacing the standard ReLU non-linearity activation function. The model architecture uses a scaled version of RoPE embedding (Su et al., 2022) with $\theta=10,000$. We scaled the RoPE embedding through a shrinking factor by dividing the target context length y by *permissible_context_size_length* on single GPU, keeping all other hyperparameters fixed such as batch size, vocabulary size, etc. This allows every *position_ids* to be divided by the shrinking ratio in the RoPE embedding methodology. For instance, if the *permissible_context_size_length* for a given physical memory hardware is 256, then we apply shrinking factor of 16 for target context size of 4096 on Single A100 40G chip during pretraining. Then, a token with *position_ids* = 4000 becomes $4000/16 = 250$, and the neighbouring token 4001 becomes $4001/16 = 250.06$, to be within 0 to 256. This is how we can capture higher context size during pretraining on limited physical memory required to pretrain model at higher context size outside the *permissible_context_size_length*. This modification allows us to pretrain models from scratch at much higher context size than the physical memory required for pretraining. Hence, with limited physical memory and limited GPUs, we can pretrain generative language models from scratch at much higher desired context size. Following (Chowdhery et al., 2023), we remove all biases from dense layers to improve the training stability. Table 17 lists down the various configuration to develop our 13.29M, 26.59M, 87M, 87.25M, 92.63M, 108.5M, 139.3M, 162M, 208M, 350M models. Table 17 lists the different model configuration but we used weight tying (Press & Wolf, 2017) to improve the performance of language models by sharing the weights of the embedding and softmax layers. Therefore, the total number of parameters in our models are typically smaller due to weight tying and different tokenizer sizes of each language specific model than the sizes mentioned in Table 17.

3.5 TRAINING

We performed hyperparameter tuning on 15M models to find the optimal tokenizer size for each language specific monolingual corpus, learning rate, learning rate scheduler, and weight decay. We used the concept of μP transfer (Yang et al., 2021), and transferred the learned hyperparameters to our bigger models. For further training and hyperparameter tuning details, we mentioned in

Appendix A.3. We find the optimal size of tokenizers by training range of tokenizers of various sizes and optimizing the validation perplexity of the language models and fertility score of the tokenizers.

3.6 IMPACT OF LANGUAGE SPECIFIC MERGED-TOKENIZER V/S LANGUAGE AGNOSTIC TOKENIZER FOR BILINGUAL MODEL (KONKANI-MAITHILI 27M GPT) OF TYPOLOGICALLY SAME SCRIPT (DEVANAGARI)

For the language agnostic approach, we merged both the Konkani and Maithili corpus which are of comparable size after our preprocessing in terms of number of sentences (lines) then we trained SentencePiece’s BPE algorithm to a vocabulary size of 1750 on the whole corpus and trained our 27M Konkani-Maithili bilingual GPT model and we report the validation perplexity score to be 12.43393. For the language specific approach, we trained individual SentencePiece (Kudo & Richardson, 2018) BPE tokenizers independently on Konkani (vocabulary size of 1k) and Maithili (vocabulary size of 750) corpus respectively. After that, we merged both the language specific tokenizers and also removing the overlapped tokens as both Konkani and Maithili are written in Devanagari script. We used this merged tokenizer to train another 27M bilingual Maithili-Konkani GPT and found the validation perplexity of the bilingual modeling drops to 8.53827 from 12.43393. From our experiments, we infer that language specific tokenization is very important to preserve the language specific linguistics even for three typologically similar languages (Devanagari script) and the merging operation of the tokenization also helps to omit the overlapping tokens which results in optimized vocabulary size for language modeling resulting in better open-end text generation, lower perplexity score due to optimal size of the tokenizer for multilingual language modeling.

3.7 LANGUAGE MODELING FOR PARAMANU-KONKANI (DEVANAGARI SCRIPT) 15M WITH BPE TOKENIZER V/S UNIGRAM TOKENIZER

From our experiments where we tokenized the Konkani corpus with Sentencepiece’s Unigram model with a vocabulary size of 1000 and trained a 15 Model and similarly we also trained another 15M model with SentencePiece BPE tokenizer on the same setting to see whether there is an impact of these two different tokenizers on the perplexity score of monolingual language modeling. We found the validation perplexity for Unigram 15M model to be 11.88412 whereas for BPE 15M model is 11.74314 which is slightly lower but not significantly different. However, Unigram tokenizer takes longer time than BPE to get trained on the corpus.

3.8 QUANTITATIVE RESULTS BASED ON PERPLEXITY METRIC

Table 2 lists the validation perplexity and MFU metrics of our various pretrained models. In terms of quantitative evaluation of language modeling, the lower the perplexity, the better is the language model. From Table 2, Paramanu-Sanskrit exhibited the lowest validation perplexity among our models, followed by Odia, Bangla, Telugu, mParamanu (162M), Tamil, Marathi, and Hindi, which had the highest perplexity. It is important to note that most models were trained for the same number of steps, regardless of their pretraining dataset size, meaning some may be under-trained. Continued pretraining, such as for Paramanu-Hindi, could further improve perplexity scores.

3.9 BASED ON QUANTITATIVE BENCHMARK

We evaluated our models on key Natural Language Understanding (NLU), Natural Language Inference, and Commonsense Reasoning tasks, including their multilingual variants, while benchmarking against multilingual LLMs like Sarvam-2B, Bloom, and Bloomz in two groups, one with LLMs of size $\leq 2B$ and another group of LLMs of size $> 2B$. We utilized the translated ARC, HellaSwag, and MMLU datasets from (Lai et al., 2023), employing Eleuther AI’s LM Evaluation Harness (Sutawika et al., 2024) for evaluation. Key assessments included HellaSwag (Zellers et al., 2019), which tests common sense reasoning by predicting scenario endings; MMLU (Hendrycks et al., 2021), measuring broad knowledge across diverse subjects; and ARC-Challenge (Clark et al., 2018), which examines complex reasoning with scientific questions. We encountered misalignment issues with the LM Evaluation Harness datasets and our models, preventing evaluations on HellaSwag except for Hindi. Other evaluations included XCOPA (Ponti et al., 2020), assessing cross-lingual common-sense reasoning; XNLI (Conneau et al., 2018), which assesses cross-lingual sentence classification

| Model | Perplexity |
|---|------------|
| Paramanu-Assamese 26.59M | 6.620 |
| Paramanu-Bangla 87.25M | 5.069 |
| Paramanu-Bangla 108.5M | 4.102 |
| Paramanu-Hindi 162M | 16.992 |
| Paramanu-Hindi 367.5M | 11.052 |
| Paramanu-Konkani-Maithili 13.29M (merged language specific tokenizer) | 8.538 |
| Paramanu-Konkani-Maithili 13.29M (language agnostic tokenizer) | 12.433 |
| Paramanu-Odia 87M | 3.068 |
| Paramanu-Sanskrit 139.33M | 1.748 |
| mParamanu 92.63M | 8.443 |
| mParamanu 162M | 6.924 |
| Paramanu-Marathi 207.73M | 8.943 |
| Paramanu-Telugu 208.25M | 5.400 |
| Paramanu-Tamil 207.84M | 7.618 |

Table 2: Perplexity of models

| N-shot | XNLI-Hindi | XStoryCloze-Hindi | XStoryCloze-Telugu | XCOPA-Tamil |
|--------|--------------|-------------------|--------------------|--------------|
| 0 | 33.49 | 52.42 | 56.06 | 54.00 |
| 5 | 34.04 | 51.49 | 54.67 | 52.40 |
| 25 | 33.23 | 52.02 | 55.92 | 49.80 |

Table 3: N-shot evaluation of pretrained Paramanu models across various benchmarks.

across 15 languages; and XStoryCloze (Lin et al., 2022), evaluating story understanding by selecting the correct ending to a four-sentence story. Together, these benchmarks comprehensively assess model performance and reasoning capabilities.

Table 4, Table 5, Table 6, Table 7 and Table 8 evaluate model performance in a zero-shot setting using accuracy metrics across translated benchmarks (ARC, MMLU, HellaSwag) for Bangla, Hindi, Tamil, and Telugu. Table 3 presents n-shot evaluations for XNLI in Hindi, XStoryCloze in Hindi and Telugu, and XCOPA in Tamil. Additionally, Table 5 and Table 6 assess various pretrained monolingual and multilingual models for Devanagari script across MMLU, HellaSwag, ARC, XStoryCloze, XNLI, and Belebele, highlighting cross-lingual language transfer among Devanagari languages (Hindi, Marathi). Notably, mParamanu, pretrained on low-resource Devanagari languages (Sanskrit, Konkani, Maithili), achieved scores of 25.86 for MMLU-Marathi, 24.84 for Hindi, 28 for Belebele-Marathi, and 25.44 for Belebele-Hindi, indicating effective language transfer to medium (Marathi) and high-resource (Hindi) languages using the mBharat tokenizer. Interestingly, neither Paramanu-Sanskrit nor mParamanu were pretrained on Hindi or Marathi but still performed well on their benchmarks. This is possibly due to the same script.

Table 5, Table 13 in Appendix and Table 6, Table 14 in Appendix show that Paramanu models exhibit superior performance across various benchmarks in Devanagari languages, despite their smaller sizes and being pretrained on fewer tokens than larger multilingual LLMs. Specifically, Paramanu-Marathi (208M) outperformed Sarvam (2B), OpenHathi (7B), and Bloom (560M) on the Marathi benchmark. Similarly, mParamanu (162M) outperformed Paramanu-Hindi (367M) and demonstrated competitiveness against larger models. Notably, Paramanu-Hindi-instruct (356M) surpassed all larger multilingual LLMs, except Bloomz (7B), by a significant margin, benefiting from instruction tuning on a dataset of 27,000 Hindi instructions and additional 52,000 Alpaca machine-translated instructions. In contrast, Bloomz (7B) was trained on hundreds of thousands of instructions. If the 367M Hindi model underwent more training steps, it could potentially achieve even better performance, as many models were pretrained for the same duration regardless of dataset size. The stronger performance of mParamanu in Hindi illustrates effective language transfer within the same script and typology. Notably, Paramanu-Sanskrit (139M), pretrained on 45 billion tokens, achieved an average score of 31.05, surpassing both Hindi models and closely approaching Bloom (560M) and Bloomz (560M). Its lower perplexity (1.75) compared to Paramanu-Hindi (11.05) further supports the notion that additional pretraining for the Hindi model could significantly enhance its downstream performance.

From Table 3, we observe that the performance of our models drop from zero-shot setting to 25 shot setting on XNLI-Hindi, XStoryCloze for Hindi and Telugu, and XCOPA for Tamil. This type of phenomenon has also been observed in PlanningBench (Valmeekam et al., 2024) where GPT-3.5-Turbo, GPT-4, GPT4-o performance on Blocksworld dropped from 0 shot to 1-shot significantly.

| Models | MMLU-Bangla | ARC-Bangla | Belebele-Bangla | Average (Bangla) | Belebele-Assamese |
|--------------------------------------|--------------|--------------|-----------------|------------------|-------------------|
| Paramanu-Bangla 108M (ours) | 23.82 | 25.75 | 25.11 | 24.89 | 25.33 |
| Paramanu-Bangla-instruct 108M (ours) | 27.60 | 28.50 | 32.45 | 29.52 | 30.54 |
| mParamanu 162M (ours) | 25.29 | 20.19 | 27.44 | 24.31 | 29.00 |
| Bloom 560M | 22.61 | 26.00 | 22.89 | 23.83 | 22.78 |
| Bloomz 560M (instruction-tuned) | 25.82 | 23.43 | 22.77 | 24.01 | 25.11 |
| Bloom 1.1B | 23.90 | 24.37 | 26.00 | 24.75 | 26.89 |
| Sarvam 2B | 24.05 | 28.40 | 23.22 | 25.22 | 27.78 |

Table 4: Zero-shot evaluation of LLMs ($\leq 2B$) across translated benchmarks of MMLU, HellaSwag, ARC datasets, and Belebele in Bengali script. All benchmarks report Accuracy except for ARC which reports Normalized Accuracy. **Max scores are in bold.**

| Models | MMLU-Marathi | ARC-Marathi | Belebele-Marathi | Average (Marathi) |
|-------------------------------------|--------------|--------------|------------------|-------------------|
| mParamanu 162M (ours) | 25.68 | 22.16 | 28.00 | 25.28 |
| Paramanu-Hindi 367M (ours) | 23.78 | 24.16 | 24.66 | 24.20 |
| Paramanu-Hindi-instruct 367M (ours) | 28.72 | 27.85 | 32.00 | 29.52 |
| Paramanu-Marathi 208M (ours) | 25.39 | 26.49 | 27.33 | 26.40 |
| Paramanu-Sanskrit 139M (ours) | 24.96 | 26.49 | 24.33 | 25.26 |
| Bloom 560M | 22.78 | 24.50 | 27.00 | 24.76 |
| Bloomz 560M (instruction-tuned) | 26.20 | 24.24 | 25.44 | 25.29 |
| Bloom 1B | 23.93 | 25.10 | 28.33 | 25.78 |
| Sarvam 2B | 23.96 | 27.53 | 26.77 | 26.08 |

Table 5: Zero-shot evaluation of LLMs ($\leq 2B$) for cross-lingual language transfer in Marathi. All benchmarks report Accuracy except for ARC (Normalized Accuracy). **Max scores are in bold.**

Perhaps n-shot examples become additional soft constraints on the generation which might be the reason of degradation of performance from the original training dataset. From Table 9, we see that Paramanu-Bangla 108M outperformed Bloom 560M by 1.06% points, Bloomz 560M by 1.05% points, Bloom 1.1B on average score across MMLU, ARC, and Belebele benchmarks, and by 1.21% points on MMLU over Bloom 560M despite being smaller by 10.2 times compared to Bloom 1.1B and being pretrained on 26.21 billion of tokens. However, Paramanu-Bangla 108M is extensively trained only on Bangla literature corpus. With further instruction-tuning on 27k Bangla instructions, Paramanu-Bangla-instruct 108M outperformed Bloom 560M, Bloomz 560M, Bloom 1.1B, Sarvam 2B, Bloom 7B on average score of MMLU, ARC, and Belebele benchmarks for Bangla respectively.

Table 7 and Table 8 compare pretrained multilingual LLMs ($< 2B$) and instruction-tuned models on Tamil and Telugu benchmarks. [Table 15 and Table 16 in Appendix compares our model with LLMs \(\$> 2B\$ \)](#) Our model, Paramanu-Tamil (208M), outperformed larger multilingual LLMs like Bloom (560M), Bloomz (560M), and Bloom (1.1B) across four benchmarks (Belebele, XCOPA, MMLU, and ARC) in Tamil, coming close to Sarvam (2B) despite being much smaller and trained on fewer tokens. On MMLU-Tamil, both Paramanu-Tamil and Paramanu-Tamil-instruct outperformed Sarvam (2B) by 2.89 percentage points, with Paramanu-Tamil pretrained on 26.2 billion tokens. Paramanu-Tamil-instruct surpassed Bloom (7B), despite being 34 times smaller, and outperformed Bloom (1.1B) by 2.28 points, Bloom (560M) by 3 points, and Bloomz (560M) by 1.17 points. For Telugu, Paramanu-Telugu-instruct (208M) outperformed Bloom (560M) by 2.72 points, Bloomz by 2.25 points, Bloom (1.1B) by 1.46 points, and Bloom (7B) by 1.22 points, with Paramanu-Telugu pretrained on 39.32 billion tokens. The improvements in metric scores for Tamil and Telugu instruction-tuned models were modest, likely due to lower-quality machine translations from Bangla compared to Hindi. Nonetheless, these results show strong performance of our models on various NLP tasks despite their smaller size and fewer tokens, challenging the notion that larger models are always better. Our findings suggest that smaller pretrained models can excel when trained on high-quality, preprocessed data over multiple epochs, outperforming larger models trained on lower-quality data for an epoch.

3.9.1 BASED ON HUMAN EVALUATION

We hard-prompted various LLMs (GPT-2 XL, GPT Neo 1.3B, LLaMa 2 7B, OPT 6.7B, and the multilingual Bloom series) alongside our pretrained models (Paramanu-Bangla, Paramanu-Hindi, and mParamanu for Sanskrit) without fine-tuning. The prompts reflected the local, cultural, and literary contexts of Assamese, Bangla, Hindi, Konkani, Maithili, Odia, and Sanskrit. Due to resource constraints, extensive evaluations focused on Paramanu-Bangla, Paramanu-Hindi, and mParamanu. The top three predictions from each model were generated with temperature = 1.0 and $\text{top}_p = 0.9$.

| Models | MMLU-Hindi | HellaSwag-Hindi | ARC-Hindi | XStoryCloze-Hindi | XNLI-Hindi | Belebele-Hindi | Average (Hindi) |
|-------------------------------------|--------------|-----------------|--------------|-------------------|--------------|----------------|-----------------|
| mParamanu 162M (ours) | 24.84 | 24.87 | 22.35 | 49.24 | 33.70 | 25.44 | 30.07 |
| Paramanu-Hindi 367M (ours) | 24.38 | 24.83 | 27.05 | 47.92 | 32.00 | 23.33 | 29.92 |
| Paramanu-Hindi-instruct 367M (ours) | 30.25 | 29.42 | 30.23 | 58.00 | 40.25 | 42.78 | 40.14 |
| Paramanu-Marathi 208M (ours) | 25.49 | 26.59 | 23.97 | 48.71 | 33.73 | 27.33 | 30.97 |
| Paramanu-Sanskrit 139M (ours) | 25.16 | 25.64 | 25.17 | 50.23 | 34.46 | 25.66 | 31.05 |
| Bloom 560M | 23.67 | 27.50 | 23.88 | 54.79 | 40.80 | 26.44 | 32.84 |
| Bloomz 560M (instruction-tuned) | 25.87 | 26.48 | 24.40 | 55.53 | 35.58 | 26.00 | 32.31 |
| Bloom 1B | 23.86 | 28.28 | 24.74 | 55.59 | 42.77 | 28.00 | 33.87 |
| Sarvam 2B | 24.54 | 33.66 | 28.00 | 60.29 | 46.74 | 24.44 | 36.27 |

Table 6: Zero-shot evaluation of LLMs($\leq 2B$) for cross-lingual language transfer in Hindi. All benchmarks report Accuracy except for ARC (Normalized Accuracy). **Max scores are in bold.**

| Models | Belebele-Tamil | XCOPA-Tamil | MMLU-Tamil | ARC-Tamil | Average (Tamil) |
|-------------------------------------|----------------|--------------|--------------|--------------|-----------------|
| Paramanu-Tamil 208M (ours) | 26.88 | 57.60 | 24.37 | 24.51 | 33.34 |
| Paramanu-Tamil-instruct 208M (ours) | 30.22 | 56.00 | 26.95 | 26.04 | 34.80 |
| Bloom 560M | 27.22 | 55.80 | 23.95 | 25.57 | 33.13 |
| Bloomz 560M (instruction-tuned) | 23.55 | 58.60 | 25.78 | 25.30 | 33.30 |
| Bloom 1.1B | 25.77 | 57.00 | 24.67 | 24.34 | 32.94 |
| Sarvam 2B | 27.44 | 63.00 | 24.06 | 26.53 | 35.25 |

Table 7: Zero-shot evaluation of LLMs ($\leq 2B$) in Tamil script models. All benchmarks report Accuracy except for ARC which reports Normalized Accuracy. **Max scores are in bold.**

~~Native speakers assessed the outputs on Grammar, Coherence, Creativity, and Factuality, each scored from 0 (worst) to 5 (best).~~ For human evaluation, we asked 10 annotators to evaluate top-3 responses from models for each prompt on a scale of 0 (worst) to 5 (best). We report the average score of all ratings. We also have reported normalised scores of ratings in Table 19 in appendix to handle inconsistencies among annotators. We reached inter-annotator kappa score of 0.85 for Bengali, 0.79 for Hindi, and 0.72 for Sanskrit. Figure 4 in the Appendix shows the bar chart for inter-annotator agreement’s Fleiss Kappa score. For Factuality, higher scores indicated better alignment with real events, with some evaluators assigning a score of 0 when premises could not be verified.

Table 24 displays the human evaluation of Paramanu-Bangla 87.25M model for the mentioned Bangla prompts. Paramanu-Bangla 87.25M model scored an average score of 3.5/5 on grammar, 3.325/5 on coherence, 3.225/5 on creativity, and 3.2/5 on factuality metrics across top 3 generations for each Bangla prompt. Table 10 compares the performance of Paramanu-Hindi 162M and other LLMs including multilingual Bloom which was pretrained on Indian languages. In this table, we can complete see that none of the open source LLMs have the ability to generate grammatically, coherent sentences in Hindi except the Bloom series. Our monolingual model, Paramanu-Hindi 162M has performed better by 17.25% on grammar, by 46.05% on coherence, by 62.5% on creativity, and by 238.5% on factuality compared to Bloom 3B model despite being 19 times smaller in size. Table 11 in Appendix compares the performance of mParamanu-162M and other LLMs including multilingual Bloom which was pretrained on Indian languages. We can see that none of the LLMs have the ability to generate grammatically, coherent sentences in Sanskrit keeping aside the factuality. Our multilingual model, mParamanu-162M has scored the highest among all on grammar (3.75/5), coherence (3.166/5), creativity (2.166/5), and factuality (1.75/5) whereas Bloom 3B scored 0.166/5 on grammar, 0.0833/5 on coherence, and 0/5 for both creativity and factuality metrics respectively. GPT-3.5-Turbo (ChatGPT) has scored very poorly 0.25/5 on grammar & coherence metrics, 0.1818/5 in creativity and 0.33/5 on factuality metrics respectively for Sanskrit text generation. mParamanu-162M is smaller by 44.25 times compared to 7B LLaMa-2 model and yet it has shown its high quality text generation in Sanskrit as compared to ChatGPT, LLaMa, and Bloom series of models.

4 CONCLUSIONS AND FUTURE WORK

In this paper, we present a series of novel efficient small monolingual, bilingual, and multilingual pretrained auto regressive models the range of 13.5M to 367.5M for 10 Indian languages across 5 scripts **excluding English**, thereby enabling GenAI NLP technology accessible to billion speakers in the world and working towards including underrepresented languages to reduce the language bias and divide in today’s NLP technology. **We proposed a RoPE embedding scaling method that enables us to pretrain language models from scratch at larger sequence length context without increased GPU memory.** We proposed novel tokenization technique of combining both BPE and Unigram tokenizers **into a single tokenizer**. We also proposed and performed language specific tokenization for

| Models | Belebele-Telugu | XStoryCloze-Telugu | MMLU-Telugu | ARC-Telugu | Average (Telugu) |
|--------------------------------------|-----------------|--------------------|--------------|--------------|------------------|
| Paramanu-Telugu 208M (ours) | 26.00 | 51.42 | 25.12 | 26.32 | 32.22 |
| Paramanu-Telugu-instruct 208M (ours) | 27.50 | 58.00 | 26.75 | 25.75 | 34.50 |
| Bloom 560M | 23.55 | 55.65 | 24.10 | 23.85 | 31.78 |
| Bloomz 560M (instruction-tuned) | 22.44 | 54.86 | 26.82 | 24.91 | 32.25 |
| Bloom 1.1B | 26.88 | 56.38 | 24.53 | 24.38 | 33.04 |
| Sarvam 2B | 27.66 | 60.09 | 24.67 | 25.78 | 34.55 |

Table 8: Zero-shot evaluation of LLMs ($\leq 2B$) in Telugu script models. All benchmarks report Accuracy except for ARC which reports Normalized Accuracy. **Max scores are in bold.**

| Model | Grammar | Coherence | Creativity | Factuality |
|-------------------------------|----------------|----------------|----------------|----------------|
| GPT2-XL | 0.45833 | 0.37500 | 0.37500 | 0.37500 |
| GPT-Neo 1.3B | 0.91666 | 0.91666 | 0.91666 | 0.91666 |
| OPT 6.7B | 0.70833 | 0.70833 | 0.70833 | 0.70833 |
| GPT-J 6B | 1.12500 | 0.95833 | 0.95833 | 0.95833 |
| LLaMa 2 7B | 0.70833 | 0.70833 | 0.70833 | 0.70833 |
| Bloom 560M | 1.70833 | 1.41666 | 1.41666 | 1.37500 |
| Bloom 1.1B | 1.33333 | 1.29166 | 1.29166 | 1.29166 |
| Bloom 3B | 1.54166 | 1.29166 | 1.33333 | 1.33333 |
| Bloom 7.1B | 1.75000 | 1.16666 | 1.16666 | 1.08333 |
| GPT-3.5-Turbo | 0.58330 | 0.75000 | 0.58330 | 0.58330 |
| Paramanu-Bangla 108.5M (ours) | 4.66666 | 4.58333 | 3.76280 | 3.45833 |

Table 9: Human evaluation results (average scores of top-3 generations per prompt) of open-end text generation of Paramanu-Bangla v/s LLMs for 4 Bangla prompts on Grammar, Coherency, Creativity, and Factuality metrics. Scale is 0 (worst) to 5 (best). GPT-3.5-Turbo was accessed in October 2023.

multilingual language modeling and domain specific tokenization for monolingual language modeling for strong performance. Our multilingual tokenizer, mBharat shows the best fertility scores among Indian language tokenizers. We considered typological grouping and pretraining on comparable size of each monolingual language corpus for our multilingual mParamanu 162M to handle data imbalance and curse of multilinguality. We evaluated our models for open-end text generation with human evaluators on grammar, coherency, creativity, and factual metrics. **We reached inter-annotator kappa score of 0.85 for Bengali, 0.79 for Hindi, and 0.72 for Sanskrit.** In our evaluation, we found that none of the popular existing LLMs can generate grammatically correct and coherent sentences in 10 Indian languages despite being pretrained on Indian language corpora. Our efficient generative language models have performed better than Bloom 7B, LLaMa-2 7B, OPT 6.7B, GPT-J 6B, GPTNeo 1.3B, GPT2-XL models for open-end text generation in Assamese, Bangla, Hindi, Odia, and Sanskrit despite being 66 times to 20 times smaller in size. We also evaluated our models and compared with several multilingual LLMs across various NLU, NLI, and commonsense reasoning benchmarks. Our models **outperformed most multilingual LLMs of size 2B and performed very competitive or even better than LLMs of size 7B on various LLM benchmarks despite being smaller in size by multiple order of magnitude compared with LLMs whose size is bigger by multiple order of magnitude.** We observed language transfer phenomena from low-resource to high resource languages of same script and typology. We also instruction-tuned our pretrained models for Bangla, Hindi, Marathi, Tamil, and Telugu and show their task handling capabilities.

In future, we would like to extend our multilingual model to 22 official Indian languages and align our generative language models with multimodal encoders to develop multimodal generative language models for Indian languages.

| Model | Grammar | Coherence | Creativity | Factuality |
|----------------------------|---------------|---------------|---------------|---------------|
| GPT2-XL | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| GPT-Neo 1.3B | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| OPT 6.7B | 0.5833 | 0.1667 | 0.1667 | 0.0000 |
| GPT-J 6B | 0.5833 | 0.2500 | 0.0000 | 0.0000 |
| LLaMa 2 7B | 1.3333 | 0.3333 | 0.5000 | 0.2083 |
| Bloom 560M | 2.7917 | 2.4583 | 1.0000 | 1.1667 |
| Bloom 1.1B | 3.2917 | 2.7917 | 1.6250 | 1.3333 |
| Bloom 3B | 4.0833 | 3.1666 | 2.0000 | 1.0833 |
| Bloom 7.1B | 3.2917 | 2.7917 | 1.6250 | 1.3333 |
| Paramanu-Hindi 162M (ours) | 4.7917 | 4.6250 | 3.2500 | 3.6666 |

Table 10: Human Evaluation results (avg scores of top 3 generations per prompt) of open-end text generation of Paramanu-Hindi 162M v/s LLMs for 4 Hindi prompts on Grammar, Coherency, Creativity, and Factuality metrics. Scale is 0 (worst) to 5 (best)

REFERENCES

- Google Translate Cloud REST API. <https://cloud.google.com/translate/docs/reference/rest>, 2023. Accessed: 2023-10-30.
- Alan Ansell, Edoardo Maria Ponti, Jonas Pfeiffer, Sebastian Ruder, Goran Glavaš, Ivan Vulić, and Anna Korhonen. MAD-G: Multilingual adapter generation for efficient cross-lingual transfer. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 4762–4781, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-emnlp.410. URL <https://aclanthology.org/2021.findings-emnlp.410>.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016.
- Pramit Bhattacharyya, Joydeep Mondal, Subhadip Maji, and Arnab Bhattacharya. VACASPATI: A diverse corpus of Bangla literature. In Jong C. Park, Yuki Arase, Baotian Hu, Wei Lu, Derry Wijaya, Ayu Purwarianti, and Adila Alfa Krisnadhi (eds.), *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1118–1130, Nusa Dua, Bali, November 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.ijcnlp-main.72. URL <https://aclanthology.org/2023.ijcnlp-main.72>.
- Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanishu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. Gpt-neox-20b: An open-source autoregressive language model, 2022.
- 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. *Journal of Machine Learning Research*, 24(240):1–113, 2023. URL <http://jmlr.org/papers/v24/22-1144.html>.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018. URL <https://arxiv.org/abs/1803.05457>.
- Alexis Conneau, Rutu Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. XNLI: Evaluating cross-lingual sentence representations. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2475–2485, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1269. URL <https://aclanthology.org/D18-1269>.
- 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 Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetraault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 8440–8451, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.747. URL <https://aclanthology.org/2020.acl-main.747>.
- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. Free dolly: Introducing the world’s first

truly open instruction-tuned llm, 2023. URL <https://www.databricks.com/blog/2023/04/12/dolly-first-open-commercially-viable-instruction-tuned-llm>.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning of quantized LLMs. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=OUIFPHEgJU>.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.

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. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12402–12426, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.693. URL <https://aclanthology.org/2023.acl-long.693>.

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 Roziere, 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 Canton 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 AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan 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, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, 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 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, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Del-pierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha

Jain, Adam Kelsey, Adam Shajnfeld, Adithya 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, Benjamin 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, Daniel 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 Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damla, 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, Kai Wu, Kam Hou U, 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, Lavender A, 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 Keaneally, 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 Dollar, 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, Rutu Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin 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, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaoqian 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, 2024. URL <https://arxiv.org/abs/2407.21783>.

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. URL <https://arxiv.org/abs/2401.15006>.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference*

- on *Learning Representations*, 2021. URL <https://openreview.net/forum?id=d7KBjmI3GmQ>.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Thomas Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karén Simonyan, Erich Elsen, Oriol Vinyals, Jack Rae, and Laurent Sifre. An empirical analysis of compute-optimal large language model training. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 30016–30030. Curran Associates, Inc., 2022a. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022b.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- Taku Kudo. Subword regularization: Improving neural network translation models with multiple subword candidates, 2018.
- Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In Eduardo Blanco and Wei Lu (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 66–71, Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-2012. URL <https://aclanthology.org/D18-2012>.
- 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, 2023.
- Anne Lauscher, Vinit Ravishankar, Ivan Vulić, and Goran Glavaš. From zero to hero: On the limitations of zero-shot language transfer with multilingual Transformers. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 4483–4499, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.363. URL <https://aclanthology.org/2020.emnlp-main.363>.
- 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 generative language models. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 9019–9052, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.616. URL <https://aclanthology.org/2022.emnlp-main.616>.
- Sharan Narang, Hyung Won Chung, Yi Tay, Liam Fedus, Thibault Fevry, Michael Matena, Karishma Malkan, Noah Fiedel, Noam Shazeer, Zhenzhong Lan, Yanqi Zhou, Wei Li, Nan Ding, Jake Marcus, Adam Roberts, and Colin Raffel. Do transformer modifications transfer across implementations and applications? In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5758–5773, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.465. URL <https://aclanthology.org/2021.emnlp-main.465>.
- Marinela Parović, Goran Glavaš, Ivan Vulić, and Anna Korhonen. BAD-X: Bilingual adapters improve zero-shot cross-lingual transfer. In Marine Carpuat, Marie-Catherine de Marneffe, and

- Ivan Vladimir Meza Ruiz (eds.), *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1791–1799, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.130. URL <https://aclanthology.org/2022.naacl-main.130>.
- Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. MAD-X: An Adapter-Based Framework for Multi-Task Cross-Lingual Transfer. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 7654–7673, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.617. URL <https://aclanthology.org/2020.emnlp-main.617>.
- Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. XCOPA: A multilingual dataset for causal commonsense reasoning. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 2362–2376, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.185. URL <https://aclanthology.org/2020.emnlp-main.185>.
- Ofir Press and Lior Wolf. Using the output embedding to improve language models, 2017.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019. URL <https://api.semanticscholar.org/CorpusID:160025533>.
- Sarvam2B. sarvamai/sarvam-2b-v0.5 @ Hugging Face — huggingface.co. <https://huggingface.co/sarvamai/sarvam-2b-v0.5>, 2024. [Accessed 15-09-2024].
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In Katrin Erk and Noah A. Smith (eds.), *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1162. URL <https://aclanthology.org/P16-1162>.
- Noam Shazeer. Glu variants improve transformer, 2020.
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2022.
- Lintang Sutawika, Hailey Schoelkopf, Leo Gao, Stella Biderman, Baber Abbasi, Jonathan Tow, ben fattori, Charles Lovering, farzanehnakhaee70, Jason Phang, Anish Thite, Fazz, Aflah, Niklas Muennighoff, Thomas Wang, sdtbck, gakada, nopperl, researcher2, tttyuntian, Chris, Julen Etxaniz, Zdenk Kasner, Khalid, Jeffrey Hsu, Hanwool Albert Lee, Anjor Kanekar, AndyZwei, Pawan Sasanka Ammanamanchi, and Dirk Groeneveld. Eleutherai/lm-evaluation-harness: v0.4.1, January 2024. URL <https://doi.org/10.5281/zenodo.10600400>.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, 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, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogoziska, Dustin 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-Pluciska, Harleen Batra,

- Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, 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 Ardeshtir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, 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 Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Faret, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. Gemma 2: Improving open language models at a practical size, 2024. URL <https://arxiv.org/abs/2408.00118>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. Planbench: an extensible benchmark for evaluating large language models on planning and reasoning about change. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA, 2024. Curran Associates Inc.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17*, pp. 60006010, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- Elena Voita, Javier Ferrando, and Christoforos Nalmpantis. Neurons in large language models: Dead, n-gram, positional. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 1288–1301, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.75. URL <https://aclanthology.org/2024.findings-acl.75>.
- Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions, 2023.
- BigScience Workshop, :, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ili, 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 Froberg, 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 ako, 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, Rheza Harliman, Rishi Bommasani, Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, So-maieh 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 Heinzerling, Chenglei Si, Davut Emre Taar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesh 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 Raunak, 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 Baruwa, Amanpreet Singh, Anasztasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, 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, Zdenk Kasner, Alice Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ana Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh Haji-Hosseini, 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 Burynek, 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 Oyebade, 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 Perrián, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrmann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Imane Bello, Ishani Dash, Jihyun 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 Sängner, 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, Srishti Kumar, Stefan Schweter, Sushil Bharati, Tanmay Laud, Théo Gi-

gant, 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, 2023. URL <https://arxiv.org/abs/2211.05100>.

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. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 483–498, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.41. URL <https://aclanthology.org/2021.naacl-main.41>.

Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tuning large neural networks via zero-shot hyperparameter transfer. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=Bx6qKuBM2AD>.

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. URL <https://api.semanticscholar.org/CorpusID:159041722>.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt: Open pre-trained transformer language models, 2022.

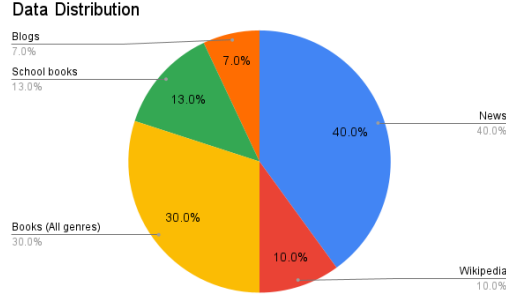


Figure 3: Pretraining Data distribution.

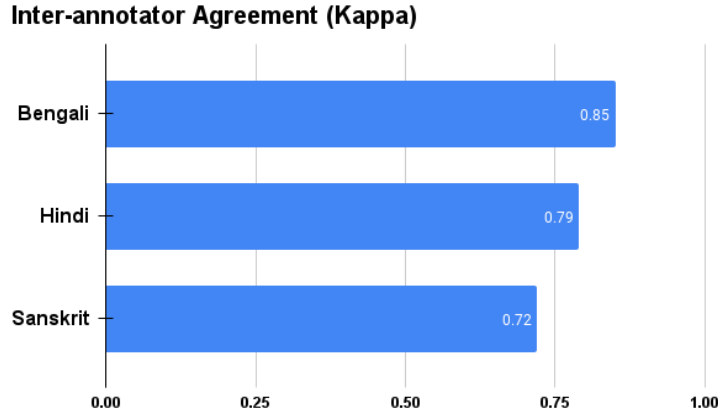


Figure 4: Inter-Annotator Agreement (Kappa)

A APPENDIX

| Model | Grammar | Coherence | Creativity | Factuality |
|-----------------------|-------------|-------------|-------------|-------------|
| GPT2-XL | 0.17 | 0.08 | 0.00 | 0.00 |
| GPT-Neo 1.3B | 0.25 | 0.00 | 0.00 | 0.00 |
| OPT 6.7B | 0.00 | 0.00 | 0.00 | 0.00 |
| GPT-J 6B | 0.33 | 0.33 | 0.00 | 0.00 |
| LLaMa 2 7B | 0.42 | 0.33 | 0.42 | 0.00 |
| Bloom 560M | 0.08 | 0.00 | 0.08 | 0.00 |
| Bloom 1.1B | 0.08 | 0.00 | 0.00 | 0.00 |
| Bloom 3B | 0.17 | 0.08 | 0.00 | 0.00 |
| GPT-3.5-Turbo | 0.25 | 0.25 | 0.18 | 0.33 |
| mParamanu 162M (ours) | 3.75 | 3.17 | 2.17 | 1.75 |

Table 11: Human evaluation results (average scores for top-3 generations per prompt) of open-end text generation of mParamanu v/s LLMs for 4 Sanskrit prompts on various metrics. Scale is from 0 (worst) to 5 (best). GPT-3.5-Turbo was accessed in October 2023.

A.1 MODEL PARAMETERS

A.2 LANGUAGE DEMOGRAPHICS

A.3 TRAINING

Following (Hoffmann et al., 2022b), we set lr decay steps to max_steps and the minimum lr is set nearly to $0.1 \cdot lr$. The lr schedule starts with a linear warm-up from 0 to the maximum lr at 1000

| Models | MMLU-Bangla | ARC-Bangla | Belebele-Bangla | Average (Bangla) | Belebele-Assamese |
|--------------------------------------|-------------|------------|-----------------|------------------|-------------------|
| Paramanu-Bangla 108M (ours) | 23.82 | 25.75 | 25.11 | 24.89 | 25.33 |
| Paramanu-Bangla-instruct 108M (ours) | 27.60 | 28.50 | 32.45 | 29.52 | 30.54 |
| mParamanu 162M (ours) | 25.29 | 20.19 | 27.44 | 24.31 | 29.00 |
| Bloom 7B | 27.10 | 26.09 | 23.22 | 25.47 | 23.11 |
| Bloomz 7B (instruction-tuned) | 32.46 | 27.20 | 53.67 | 37.77 | 48.00 |

Table 12: Zero-shot evaluation of LLMs (> 2B) across translated benchmarks of MMLU, HellaSwag, ARC datasets, and Belebele in Bengali script. All benchmarks report Accuracy except for ARC which reports Normalized Accuracy.

| Models | MMLU-Marathi | ARC-Marathi | Belebele-Marathi | Average (Marathi) |
|-------------------------------------|--------------|--------------|------------------|-------------------|
| mParamanu 162M (ours) | 25.68 | 22.16 | 28.00 | 25.28 |
| Paramanu-Hindi 367M (ours) | 23.78 | 24.16 | 24.66 | 24.20 |
| Paramanu-Hindi-instruct 367M (ours) | 28.72 | 27.85 | 32.00 | 29.52 |
| Paramanu-Marathi 208M (ours) | 25.39 | 26.49 | 27.33 | 26.40 |
| Paramanu-Sanskrit 139M (ours) | 24.96 | 26.49 | 24.33 | 25.26 |
| Bloom 7B | 27.30 | 25.54 | 24.00 | 25.61 |
| Bloomz 7B (instruction-tuned) | 32.62 | 27.44 | 53.00 | 37.68 |
| OpenHathi 7B | 26.09 | 24.24 | 25.88 | 25.40 |
| Airavata 7B (instruction-tuned) | 26.15 | 23.90 | 29.89 | 26.64 |

Table 13: Zero-shot evaluation of LLMs (>2B) for cross-lingual language transfer in Marathi. All benchmarks report Accuracy except for ARC which reports Normalized Accuracy.

steps, followed by a cosine decay to the minimum lr until the end of an epoch of training. We used the following equation for lr decay ratio.

$$lr_{decay_ratio} = \frac{t - warmup_{steps}}{lr_{decay_steps} - warmup_{steps}}$$

where t is the current training step.

A.3.1 26.58M PARAMANU-ASSAMESE MODEL

We used the same training procedure as mentioned in 3.5 but with a batch size of 64, gradient accumulation steps of 4, and the maximum sequence length set to 1024, i.e., 262,144 tokens per iteration and transferred the learned hyperparameters from 15M model to 42M model using (μP) transfer. We set maximum learning rate (lr) to $3e-3$ (max), weight decay to $1e-1$. We trained our bigger models with fused *AdamW* optimizer for an epoch of training with $\beta_1=0.9$, $\beta_2=0.95$, dropout of 0.0, and scaled the gradient norms using a maximum norm clipping value of 1.0, and weight decay of 0.1. For our experiments and modeling, we implemented our code using Pytorch 2.0, in-house optimized CUDA kernels and used `torch.compile` feature for every model. To further speedup training, we used BF16 mixed precision training.

A.3.2 87.25M PARAMANU-BANGLA MODEL

We used the same training procedure as mentioned in 3.5 but with a batch size of 32, gradient accumulation steps of 8, and the maximum sequence length set to 1024, i.e., 262,144 tokens per iteration and transferred the learned hyperparameters from 15M model to 110M model using (μP) transfer.

A.3.3 108.5M PARAMANU-BANGLA MODEL

We used the same training procedure as mentioned in 3.5 but with a batch size of 32, gradient accumulation steps of 8, and the maximum sequence length set to 1024, i.e., 262,144 tokens per iteration and transferred the learned hyperparameters from 15M model to 140M model using (μP) transfer.

A.3.4 162M PARAMANU-HINDI MODEL

We used the same training procedure as mentioned in 3.5 but with a batch size of 32, gradient accumulation steps of 8, the maximum sequence length set to 1024, i.e., 262,144 tokens per iteration and transferred the learned hyperparameters from 15M model to 162M model using (μP) transfer.

| Models | MMLU-Hindi | HellaSwag-Hindi | ARC-Hindi | XStoryCloze-Hindi | XNLI-Hindi | Belebele-Hindi | Average (Hindi) |
|-------------------------------------|------------|-----------------|-----------|-------------------|------------|----------------|-----------------|
| mParamanu 162M (ours) | 24.84 | 24.87 | 22.35 | 49.24 | 33.70 | 25.44 | 30.07 |
| Paramanu-Hindi 367M (ours) | 24.38 | 24.83 | 27.05 | 47.92 | 32.00 | 23.33 | 29.92 |
| Paramanu-Hindi-instruct 367M (ours) | 30.25 | 29.42 | 30.23 | 58.00 | 40.25 | 42.78 | 40.14 |
| Paramanu-Marathi 208M (ours) | 25.49 | 26.59 | 23.97 | 48.71 | 33.73 | 27.33 | 30.97 |
| Paramanu-Sanskrit 139M (ours) | 25.16 | 25.64 | 25.17 | 50.23 | 34.46 | 25.66 | 31.05 |
| Bloom 7B | 27.04 | 31.39 | 26.36 | 60.55 | 47.18 | 23.00 | 35.92 |
| Bloomz 7B (instruction-tuned) | 35.55 | 28.57 | 29.36 | 57.71 | 40.52 | 53.11 | 40.80 |
| OpenHathi 7B | 27.69 | 30.54 | 25.51 | 57.04 | 39.03 | 32.66 | 35.41 |
| Airavata 7B (instruction-tuned) | 30.43 | 29.53 | 25.60 | 55.59 | 39.04 | 41.44 | 36.93 |

Table 14: Zero-shot evaluation of LLMs (>2B) for cross-lingual language transfer in Hindi. All benchmarks report Accuracy except for ARC which reports Normalized Accuracy.

| Models | Belebele-Tamil | XCOPA-Tamil | MMLU-Tamil | ARC-Tamil | Average (Tamil) |
|-------------------------------------|----------------|-------------|------------|-----------|-----------------|
| Paramanu-Tamil 208M (ours) | 26.88 | 57.60 | 24.37 | 24.51 | 33.34 |
| Paramanu-Tamil-instruct 208M (ours) | 30.22 | 56.00 | 26.95 | 26.04 | 34.80 |
| Bloom 7B | 25.55 | 59.20 | 26.39 | 24.69 | 33.95 |
| Bloomz 7B (instruction-tuned) | 50.66 | 57.40 | 29.48 | 28.10 | 41.41 |

Table 15: Zero-shot evaluation of LLMs (>2B) in Tamil script models. All benchmarks report Accuracy except for ARC which reports Normalized Accuracy.

A.3.5 367.5M PARAMANU-HINDI MODEL

We used the same training procedure as mentioned in 3.5 but with a batch size of 32, gradient accumulation steps of 16, the maximum sequence length set to 1024, i.e., 524,288 tokens per iteration and transferred the learned hyperparameters from 15M model to 367.5M model using (μ P) transfer. After 1 epoch of training, the average validation perplexity is 11.05240 whereas the average training perplexity is 10.99616.

A.3.6 87M PARAMANU-ODIA MODEL

We used the same training procedure as mentioned in 3.5 but with a batch size of 64, gradient accumulation steps of 8, and the maximum sequence length set to 1024, i.e., 524,288 tokens per iteration and transferred the learned hyperparameters from 15M model to 110M model using (μ P) transfer.

A.3.7 139.3M PARAMANU-SANSKRIT MODEL

We used the same training procedure as mentioned in 3.5 but with a batch size of 64, gradient accumulation steps of 8, and the maximum sequence length set to 1024, i.e., 524,288 tokens per iteration and transferred the learned hyperparameters from 15M model to 175M model using (μ P) transfer.

A.3.8 13.5M BILINGUAL KONKANI-MAITHILI GPT MODEL

Both Konkani and Maithili languages are typologically similar (Devanagari script). For bilingual model, we trained two 27M models with language agnostic tokenization and language specific tokenization to study the impact of language specific tokenization against language agnostic tokenization for multilingual language modeling. For training with the language specific tokenization, where we basically trained independent BPE tokenizer on Konkani and Maithili corpora with a tokenizer size of 1000 and 750 respectively and merged them together. For language agnostic tokenization, we trained BPE tokenizer on the merged corpora of Konkani and Maithili with a tokenizer size of 1750.

We used the same training procedure as mentioned in 3.5 but with a batch size of 128, gradient accumulation steps of 2, and the maximum sequence length set to 1024, i.e., 262,144 tokens per iteration. We transferred the learned hyperparameters from 15M model to 27M model using (μ P) transfer.

A.3.9 92.63M MULTILINGUAL MPARAMANU MODEL

For our multilingual mParamanu model, we selected the languages on the basis of typological grouping and having comparable corpora as shown in the Table 19, to avoid pretraining our multilingual model on disproportionate corpora of multiple languages. Thus, we avoid adding the Hindi corpus

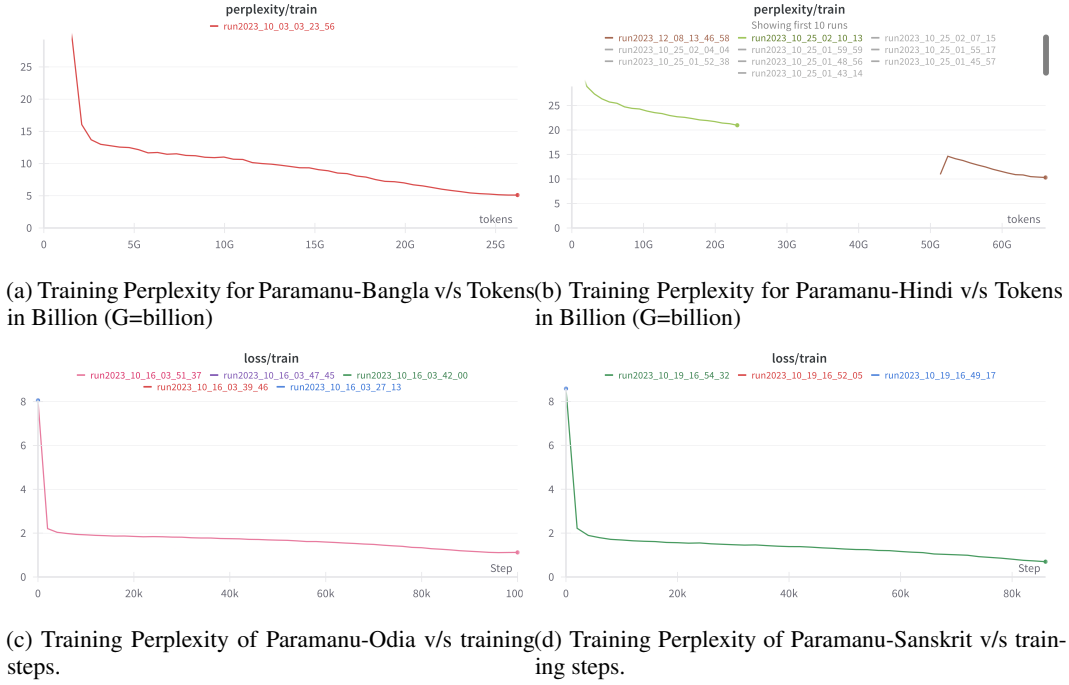


Figure 6: Training Perplexity v/s Tokens and Training Perplexity v/s training steps for Paramanu pretrained models.

| n_params | d_model | n_layers | n_heads | dim_head | max_seq_length |
|----------|---------|----------|---------|----------|----------------|
| 15M | 288 | 6 | 6 | 48 | 512 |
| 27M | 512 | 4 | 8 | 64 | 1024 |
| 42M | 512 | 8 | 8 | 64 | 1024 |
| 110M | 768 | 12 | 12 | 64 | 1024 |
| 140M | 768 | 15 | 12 | 64 | 1024 |
| 350M | 1024 | 12 | 16 | 64 | 1024 |
| 425M | 1280 | 18 | 10 | 128 | 1024 |

Table 17: Model architectures considering tokenizer size of 32000

- Removal of English literals, Roman digits, French, German, Italian, Russian, Chinese literals and punctuation following Unicode representation using regular expressions
- Removal of emoticons, symbols, pictographs, transport & map symbols, and iOS flags following Unicode representation using regular expression

| Language | Family | Script | #Speakers |
|----------|----------------|------------------|-----------|
| Assamese | Indo-European | Assamese-Bengali | 24 M |
| Bangla | Indo-European | Bengali | 300 M |
| Hindi | Indo-European | Devanagari | 692 M |
| Konkani | Indo-European | Devanagari | 2 M |
| Maithili | Indo-European | Devanagari | 14 M |
| Marathi | Indo-European | Devanagari | 99 M |
| Odia | Indo-European | Odia | 43 M |
| Sanskrit | Indo-European | Devanagari | 0.025 M |
| Tamil | Indo-Dravidian | Tamil | 77 M |
| Telugu | Indo-Dravidian | Telugu | 95 M |

Table 18: Speaker estimates according to the Indian Census 2011

6. Removal of links, emails, HTML/XML tags, emojis, language specific punctuation, personal information like phone number, address, ID number using regular expression. We also deduplicated web scrapped pretraining corpora in respective languages

| Language | Corpus Source | Corpus Size | #Sentences |
|----------|---|-------------|------------|
| Assamese | Indic Corp v2 + Wikipedia + Curated books (ours) | 3.2 GB | 5,734,166 |
| Bangla | Vacasapati + Wikipedia + Curated books (ours) | 3.6 GB | 22,533,608 |
| Hindi | IITB monolingual + Wikipedia + Curated books (ours) | 15.8 GB | 52,124,643 |
| Konkani | Indic Corp v2 | 516.5 MB | 1,337,693 |
| Maithili | Indic Corp v2 | 191.3 MB | 947,386 |
| Marathi | Indic Corp v2 + Wikipedia + Curated books (ours) | 12.5 GB | 34,567,839 |
| Odia | Indic Corp v2 + Wikipedia | 6.2 GB | 14,657,392 |
| Sanskrit | Indic Corp v2 | 6.7 GB | 17,034,631 |
| Tamil | Indic Corp v2 + Wikipedia + Curated books (ours) | 10.7GB | 27,872,768 |
| Telugu | Indic Corp v2 + Wikipedia + Curated books (ours) | 13.5 GB | 40,241,847 |

Table 19: Pretraining Data details after data cleaning.

A.5 HUMAN EVALUATION

A.5.1 MPARAMANU-162M VS BLOOM 1.1B EVALUATION

From Table 20 we see that Bloom 1.1B model could not distinguish languages of the same script so when we prompt Bloom with Konkani, Maithili, and Sanskrit (Devanagari script), Bloom only generated incoherent, grammatically incorrect text in Hindi whereas our multilingual model, mParamanu-162M was able to recognise prompt in respective distinguished languages Konkani, Maithili, and Sanskrit (Devanagari) to generate grammatically sound and coherent text in respective languages. Similarly, Bloom 1.1B could not identify Assamese (Assamese-Bengali Script) and generated out of context Bangla words including series of Bangla alphabets for the Assamese prompt as shown in the table 20 whereas mParamanu-162M model generated coherent sentences in Assamese when prompted with Assamese prompt. However, mParamanu-162M model is 7 times smaller in size than Bloom 1.1B and has better text generation capabilities in terms of grammar, coherence metrics for languages under consideration. This makes us to ascertain that despite Bloom being pretrained on 45 languages including Indian languages yet its text generation capabilities is very limited in terms of grammatically correct and coherent text generation for Indian languages. Therefore, its not the right way to pretrain on many languages together varying script at the same time without considering linguistics typology, features, and disproportionate pretraining low quality corpus in multilingual setting.

A.5.2 BANGLA EVALUATION

We hard prompted the pretrained models with the following Bangla prompts, prompt 1: ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! (phēludā tōpasēkē balala, yā lālamōhanabābukē khabara dē) – Feluda told Tapshe, Go and call Lalmohanbabu (from a popular detective novel series), prompt 2: এমন সময় হঠাৎ বাঘের ডাক! (ēmana samāya haṭhāt bāghēra ḍāka!) – suddenly there was a tiger roar (astonishing fearful expression), prompt 3: অপু এসে ডাকল, মা! (apu ēsē ḍākala, mā!) – Apu came and called, Mother (Apu is a character in a popular novel, also adapted into a popular movie), prompt 4: সে খুব খুশী হলো। (sē khuba khuśī halō) – she became very happy (generic statement). All these prompts are related either to common expression in Bangla or covers diverse Bangla literature and cultural context in Bangla speaking community in West Bengal, India.

Table 24 displays the human evaluation of Paramanu-Bangla 87.25M model for the mentioned Bangla prompts. We see Paramanu-Bangla 87.25M model scored an average score of 3.5/5 on grammar, 3.325/5 on coherence, 3.225/5 on creativity, and 3.2/5 on factuality metrics across top 3 generations for each Bangla prompt. Table 9 compares the performance of Paramanu-Bangla 108.5M and other LLMs including multilingual Bloom which was pretrained 45 languages including Indian languages. In this table, we can see that none of the open source LLMs have the ability

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

1344
1345
1346
1347
1348
1349

| Input Prompt | Text Generation |
|-----------------------------------|---|
| কামাখ্যা মন্দিৰ (Kamakhya Temple) | কামাখ্যা মন্দিৰত উপস্থিত হ'ব গৃহস্থানী অমিত যাহ। (kāmakhyā' mandirata upasthita ha'bā gṛhamantrī amita śbāha.) কামাখ্যা মন্দিৰ দৰ্শনৰ পাছতে দুখামত্ৰী হিমত বিৰি সৰ্মাই ক্লেগ স্ব' মুকলি কৰিব। (kāmakhyā' mandira darśanara pāchatē mukhyamantrī himanta biśba sarmāi phlēga śba' mukali kariba.) কামাখ্যা মন্দিৰলৈ যোৱা পথত যোৱা কোনো এজন লোককে ক'ভিড টেষ্ট কৰিবলৈ অন্তৰ্বেধ আৰু আহ্বান। (kāmakhyā' mandiralai yōvā pathata yōvā kōnō ējana lōkākē ka'bhida tēsta karibalai anurōdhā āru āhṛāṇā.) |
| অসম চাহ (Assam Tea) | অসম চাহ মজদুৰ সংঘৰ মৰিগণ শাখাৰ ধৰ্মঘট। (asama chāha majadura saṁghara dharmā ghaṁḍa) (asama chāha majadura saṁghara margigāṇō śākhāra dharmaghata.) অসম চাহ মজদুৰ সংঘৰ নিৰ্দেশ মৰ্মে প্ৰমিকৰ মৰুবি ৩৫০তকালৈ বৃদ্ধি কৰাৰ দাবীত ১ নংঘৰ পৰা ৩ নংঘৰলৈ তিনিদিনীয়া কাৰ্যসূচীৰে বাগিছা সমূহৰ কাৰ্যবাহী পৰিচালক মজলিলত চাহ শ্ৰমিকৰ বিভিন্ন দাবীত চতুৰ্থ বৰ্ষৰ কৰ্মচাৰীৰ এক বিশাল গণ ধৰ্মঘট কাৰ্যসূচী ৰূপায়ণ কৰা হৈছে। (asama chāha majadura saṁghara nirdēśa marmē śramikara majuri 350takalāi bṛddhi karāra dābita 1 navēmbara parā 3 navēmbaralāi tinidiniyā kāryasūciṛē bāgichā samūhara kāryabāhī paricālakā majallilā chāha śramikara bibhinna dābita catuttha bargara karmacārira ēka biśālā gaṇa dharmaghata kāryasūci rūpāyāṇa karā hāya.) অসম চাহ মজদুৰ সংঘৰ যোহাট শাখাৰ সভাপতি ৰাজেন গোঁহাই এক সম্পাদক সৌৰভ গগনেৰু কৰা বৈ পৰৱৰ্তী সময়ত উপযুক্ত ন্যায় বাৰ্তাৰ ক্ষেত্ৰত অসম চৰকাৰৰ মিনিস্ট্ৰি অব বেকেল লট অব অসম আৰু অসম আৰাম গৃহৰ ব্যৱস্থাপক বাবে এক লোকক অসমীয়া দৰিদ্ৰ ৰজাঘৰীয়া আৰ্থিক সাহায্য প্ৰদান কৰি আহিছে। (asama chāha majadura saṁghara yōrahātā śākhāra sabhāpātī rājēna gōmhāi āru sampādaka saurabha gāyanē kayā yē paravartī samayāt upayukta nyāya byavasthāra kśētrata asama carakārāra minisṭrī aba bēmka laṭa ava asama āru asama ārama gṛhara byavasthāra bābē bahu lōkaka asama carakārā prēmī daridra rājāgharīyā ārthika sāhāyā pradāna karāi āhicē.) |
| ভূপেন হাজৰিকা (Bhupen Hazarika) | ভূপেন হাজৰিকা ন্যাস প্ৰৱাহটোয়ে এই ন্যাস বাবে প্ৰদান কৰিছে জনজাতীয় পুৰস্কাৰ। (bhūpēna hājārikā nyāsa guvāhātīyē ēi nyāsara bābē pradāna karicē janajātiyā puraskāra.) ভূপেন হাজৰিকা হেমচন্দ্ৰ বৰুৱা অসমীয়া বিদোদন শাস্ত্ৰৰ পিতৃ স্বৰ্গীয় পিতাপত্ৰৰ স্বধাৱিকাৰী হেমচন্দ্ৰ গোস্বামী আছিল অসমীয়া ভাষা সাহিত্যৰ এগৰাকী বিশিষ্ট প্ৰতিভাশালী ব্যক্তি। (bhūpēna hājārikā hēmacandra baruvā asamiyā binōdāna sāstrara pitṛ sbargīyā pitāpatra sratbādhikārī hēmacandra gōśbānī āchila asamiyā bhāśā sāhityara ēgarākī biśiṭṭhā pratibhāśālī byakti.) ভূপেন হাজৰিকা সেতুঘাট ৰ ওপৰত উজনি অসমত চলোৱা এক বহুৎ অভিযানত বৃহৎ পৰিমাণৰ ভ্ৰূপাঙ্গ জন্ম কৰা হৈছে। (bhūpēna hājārikā sētughāṭa ra ōparata ujani asamata calōvā ēka bahat abhiyānata brhat paramāṇara brāgacha jadbā karā haicē.) |
| বিহু (Bihu) | বিহু অসমৰ এবিধ লোকনাট্যনৃত্যৰ অন্যতম এক লোকনাট্যনৃত্য। (biḥu asamara ēbidha lōkanāṭyānuṣṭhānara anyatama ēka lōkanāṭyānuṣṭhāna.) বিহু উদযাপনৰ জৰিয়তেও মাক দেউতাকৰ উপস্থিতিত শিশুৰ বেলিকাও স্বচ্ছলতা লাভ কৰিব পৰা যায়। (biḥu udayāpanara jariyātē māka dēutākara upasthitiyā śiśura bēlikāō sbacacalātā lābha kariba parā yāya.) বিহু কেৰে মাজুলীৰ আকাশ বতাহ মথখিত। (biḥu cērē mājulīra ākāśa batāḥa mukharita.) |

Table 21: Paramanu-Assamese open end text generation for Assamese prompts capturing local context and culture.

| Input Prompt | Text Generation |
|--|---|
| ଜଗନ୍ନାଥ ମନ୍ଦିର (Jagannātha Mandir) | <p>ଜଗନ୍ନାଥ ମନ୍ଦିରକୁ ଦର୍ଶନ ପାଇଁ ଆସୁଥିବା ପିଲାମାନେ ମହାପ୍ରଭୁଙ୍କ ଆଦ୍ୟ ସେବାକ ଶାସନ ଗାଦିରେ ବସିଛନ୍ତି । (jagannātha mandiraku darśana pāim̃ āsūthibā pilāmānē mahāprabhūṅka ādya sēbaka śāśana gadīrē basichanti ।)</p> <p>ଜଗନ୍ନାଥ ମନ୍ଦିରକୁ ସଂଯୋଗ କରୁଥିବା ଜଗନ୍ନାଥ ମନ୍ଦିରର ବିଦ୍ୟା ସ୍ଥାନର ବିଶେଷ ପାଇଁ ଶେକେଟି ଟଙ୍କା ଅନୁଦାନ ଦିଆଯାଇଥିବାବେଳେ ବସୁ ଶ୍ଯାସ୍ତ୍ରର ନବୀନକଣ୍ଠ ନିର୍ମାଣ କରାଯାଇ । (jagannātha mandiraku saṁyōga karuthibā jagannātha mandirara tiḥya sthānara bikāśa pāim̃ 5kōṭi [ṅka] anudāna diāyāitibhābēlē baṣ ṣṭāṇḍara nabikaraṇa nirmāṇa karāyāi ।)</p> <p>ଜଗନ୍ନାଥ ମନ୍ଦିର ପରିଚାଳନା କରିଥିବା ପୂଜକର ପଦସ୍ଥ ହେଉଛନ୍ତି ଅନନ୍ତ କୁମାର ଜେନା । (jagannātha mandira paricālanā kamiṭira pūrbatana sadasya hēuchanti ananta kumāra jēnā ।)</p> |
| ପୁରୀ ସାଗର ବେଳାଭୂମି (Puri Sea Beach) | <p>ପୁରୀ ସାଗର ବେଳାଭୂମିରେ ଶୁଭାରମ୍ଭ ହେଉଛି ଆଜ୍ଞାତମାୟ ଖ୍ୟାତି ସମ୍ପନ୍ନ ବାଲୁକା ଶିଳ୍ପୀ ସୁଦର୍ଶନ ପଟ୍ଟନାୟକଙ୍କ ଦ୍ଵାରା ବିଶ୍ଵ ବାଲୁକା କୃତୀ । (purī sāgara bēlābhūmīrē śubhārambha hōichi antarjāṭīya khyāti samparṇṇa bālūkā śilpī sūdarśana paṭṭanāyakaṅka dvārā biśva bālūkā kṛtī ।)</p> <p>ପୁରୀ ସାଗର ବେଳାଭୂମିର ମହୋତ୍ସବର ପ୍ରଥମ ଦିନରେ ପୁରୀ ଶ୍ରୀଜଗନ୍ନାଥ ମନ୍ଦିରରେ ଆକର୍ଷଣୀୟ ଦୃଶ୍ୟକାଣ୍ଡ ଦୃଶ୍ୟ । (purī sāgara bēlābhūmī mahōtsabara prathamā dinārē purī śrījagannātha mandirārē ākarsāṇīya dr̥śyakaṅḍa dr̥śya ।)</p> <p>ଜଗନ୍ନାଥ ମନ୍ଦିର ବେଳାଭୂମିରେ ୩୨ ଫୁଟର ପାନ୍ଥେରା । (purī sāgara bēlābhūmīrē 32 phutara pānthēra ।)</p> |
| ଚିଲିକା ହ୍ରଦ (Chilika Lake) | <p>ଚିଲିକା ହ୍ରଦ ଭିତରକନିକା ପ୍ରଭୃତିରେ ବିଭିନ୍ନ ପ୍ରଜାତିର ଚଳଚଫଳ ମିଳୁଥିଲା । (ciliikā hrada bhitarkanikā prabhṛtīrē bibhinna prajāṭira ḍalaphina miḷuthilā ।)</p> <p>ଚିଲିକା ହ୍ରଦ ଭୂଇଁରେ ଛୁଇଁବ ସମୁଦ୍ର ଆସିଲା ସାମାନ୍ୟ ବହଳିଆ ଡୁଆର । (ciliikā hrada bhūim̃irē chuiṁba samudra āsilā sāmānya bahaliā ḍuāra ।)</p> <p>ଚିଲିକା ହ୍ରଦ ଗୁଡୁଣୀ ବୈଦିକାଙ୍କୁ ଏ ଅନ୍ୟ ଜଳପ୍ରାପ୍ତ ଦେଇ ଏହି ନଦୀ ପୁରୀକୁ ପ୍ରଜନନ କରେ ଏ ଗୁଡୁଣୀ ନଦୀର ଜଳ ଯାଇପୁର ବିଭାୟନନପୁରରେ ପ୍ରସ୍ତୁତିତ ରାଘବପତି ମୁନୁରାବଳୀରେ ବସି ଲିଖିଲା ପ୍ରଭୁପାଦ ଶ୍ରୀ । (ciliikā hrada brāhmanī baitarāṇi ānāya jaladvipa dēhi ēhi nadi purāṇa prajānana karē ṇ brāhmanī nadira jala yājapura ṭikāyatanandapurārē prastāḍita rāsapathara punaruddhāra kari ē jillā pragatira sutrapāṭa huē ।)</p> |
| କୋଣାର୍କ ସୂର୍ଯ୍ୟ ମନ୍ଦିର (Konark Sun Temple) | <p>କୋଣାର୍କ ସୂର୍ଯ୍ୟ ମନ୍ଦିରକୁ ବିଶ୍ଵ ବିଦ୍ୟାର ମାନ୍ୟତା ଦାବିରେ ମୁଖ୍ୟମନ୍ତ୍ରୀଙ୍କୁ ଚିଠି ଲେଖିଲେ କୋଣାର୍କ ବିଧାୟକ । (kōṇārka sūryya mandiraku biśva tiḥyara mānyāta dābirē mukhyāmantrin̄ku ciṭhi lēkhilē kōṇārka bidhāyaka ।)</p> <p>କୋଣାର୍କ ସୂର୍ଯ୍ୟ ମନ୍ଦିର ପାଦଦେଶର ଉପର ଘାଟି । (kōṇārka sūryya mandira pādādēśārē kōṇārka mahōtsaba ud ghāṭitā ।)</p> <p>କୋଣାର୍କ ସୂର୍ଯ୍ୟ ମନ୍ଦିରର ବିଭିନ୍ନ ଅବଶ୍ୟାସ୍ତ୍ର ନେଇ କୋଣାର୍କ ଯାତ୍ରା ପ୍ରତିବାଦ କରିବା ପରେ ବିଦେଶୀ ଯାତ୍ରୀଙ୍କ ଆଗମନ ଉପରେ ପ୍ରଶାସନ ୧୪୪ ଧାରା ଲାଗୁ କରିଥିଲା । (kōṇārka sūryya mandirara bibhinna abhyābasthā nēi kōṇārka yātrī pratibāda karibā parē bidēśī yātrin̄ka āgamana unārē prāśāśana 144 dhārā lāgu karithilā ।)</p> |

Table 22: Paramanu-Odia’s text generation for 4 Odia prompts capturing local context and culture.

| Language | Input Prompt | Text Generation |
|----------|-------------------------|--|
| Assamese | ভাৰত মোৰ মাতৃভূমি | মোৰ দেশ ভাৰত বিশ্বৰ বিভিন্ন ঠাইৰ সংস্কৃতিকালৈয়ে গঢ় লৈ উঠিছে। মোৰ দেশ ভাৰত এখন সাংস্কৃতিক কেন্দ্ৰ। |
| Bangla | ভাৰত আমাৰ মাতৃভূমি | আমাৰ দেশ ভাৰতবৰ্ষ আমাদেৱৰ দেশ। আমাৰ দেশ ভাৰতবৰ্ষৰ বিৰাট অংশ এক সময় বঙ্গদেশৰ মত ছিল। |
| Dogri | সাদা দেশ भारत है | সাদা দেশ भारत है ही आज के समय से आने वाला उत्तर भारतीय है जहाँ धीरूभाई अम्बानी फॉर युवा का मुजुफ्फुरपुर का लोग है वहाँ के बावजूद बिग बाँस नाम की तरफ से भी कुछ लकी। |
| Konkani | भारत हो म्हजो देश | भारत हो म्हजो देश ह्या राष्ट्रची वाट लायल्या भारत हो म्हजो देश भारत म्हज्या इतिहासांत देव म्हज्या इतिहासांत आस्थावता आदलें पैकूच आसा त्या इतिहासाचो आदलो काळ आशिल्लो हाका इतिहासाचो चवथो भाग म्हण्टा। |
| Maithili | भारत हमर देश अछि | भारत हमर देश अछि और जो देशक सर्वाधिक पदरें ओतप्रोत अछि भारत हो म्हजो देश हो मोलांत। |
| Odia | ଭାରତ ହେଉଛି ମୋର ମାତୃଭୂମି | ଭାରତ ହେଉଛି ମୋର ମାତୃଭୂମି ଭାରତ ଏକତାକରଣ ହୋଇଛି ଭାରତ ହେଉଛି ମୋର ମାତୃଭୂମି ଓ ମୋ ନିଜଭୂମି। |
| Sanskrit | अस्माकं देशः भारतवर्षः | अस्माकं देशः भारतवर्षः सर्वातिथायी। अस्माकं देशः भारतवर्षः इति ख्यातः। |

Table 23: mParamanu 162 M pretrained model text generation

| Prompt | Text Generation (top-3) | Grammar | Consistency | Coherency | Factuality |
|--|---|----------|-------------|-----------|------------|
| ফেনুমা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে। | ফেনুমা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে। আমরা আলাম ঘর থেকে বেরোতেই ধূজাটিবাবুর আর্দনাম গুনলাম। তারপর ফেনুনার লগল পেললাম। ধূজাটিবাবুর ঘরের সামনে যে লোক বসেছিল, সে ফেনুনার দিকে চেয়ে জিজ্ঞেস করল, আপনাদের নাম কী? ফেনুমা আর্দনাম বলে ফেনুমা বাইরে চলে গেল। আমি আরও মিনিটিবাকের বাইরে হটলাম। ফেনুমা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে। উনি আসছেন। ফেনুমা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে। ভদ্রলোক জিজ্ঞেস করবেন, আর সজি বলতে কী, আমার সন্দেহ আছে। বাইতাকালনাঃ-এও অবিশি উনি ছাড়া। | 3/0.66 | 2.8/0.53 | 2.6/0.4 | 2.5/0.33 |
| এমন সময় হঠাৎ বাঘের ডাক! | এমন সময় হঠাৎ বাঘের ডাক! বাঘের সে ডাক বাঘের স্বর নয়! দূরে শিরিগঞ্জের ধানখেতের পিছন হইতে যেন কীপিতে কীপিতে ছুটিয়া আসিতেছে। এমন সময় হঠাৎ বাঘের ডাক! সম্ভ্র ভয়ে চী- চী করিয়া উঠিল। মনে হইল যেন সেই বনের হরিণগুলো আবার জঙ্গলের নিক হইতে আসিয়া উপস্থিত হইল। এমন সময় হঠাৎ বাঘের ডাক! আমার ভয় পাবার হয়তো কিছু ছিল। নইলে দুবে সকলের চোখের ওপরে একটা শজাকর মতো বড়ো চোখা শি- অলকলিন আসেই কোরেছি কি-না? না, না, আমি নিশ্চয় ভয় করি না। সে গলার আগুয়াজ এবার কাছেই কোথায় মিলিয়ে গেল। | 3.3/0.60 | 2.9/0.40 | 2.8/0.30 | 2.8/0.3 |
| অপু এসে ডাকল, মা! | অপু এসে ডাকল, মা! বেতে মা! অপু এসে ডাকল, মা, ও মা, কোথা যাও? অপু এসে ডাকল, মা, শিশুর এসে, গ্রামের ওদিক থেকে ডাকছে। | 4 | 4 | 4 | 4 |
| সে খুব খুশী হলো। | সে খুব খুশী হলো। সবাই তাকে এত সম্মানের চোখে দেখে। প্রথম দিন থেকে দেখে এবং অভ্যর্থনা পাওয়ার পর তার প্রতি প্রদ্বা জানায়। সে খুব খুশী হলো। ভালো হয়েছে, মা। সে খুব খুশী হলো। সবকিছু তার খোলসা হয়ে গেছে। | 3.7/0.70 | 3.6/0.6 | 3.5/0.50 | 3.5/0.50 |

Table 24: Human evaluation results of open-end text generation of Paramanu-Bangla 87.25M for 4 Bangla prompts on Grammar, Coherency, Creativity, and Factuality metrics. We use the format (a/b) where a is the average score of top 5 generations on scale of 0 (worst) to 5 (best) and b is the normalized score according to a ; $b = a_{min} / (a_{max} - a_{min})$

to generate grammatically, coherent sentences in Bangla except the Bloom series keeping aside the factuality. Our monolingual model, Paramanu-Bangla 108.5M model has scored 202.7% better and 166.6% better on grammar metric, 254.84 % and 292.86% on coherence metric, 182.22% and 222.53% on creativity metric, 159.4% and 219.23% on factuality metric than Bloom 3B, and Bloom 7B models respectively, despite Paramanu-Bangla is 28.3 times and 67 times smaller in terms of total number of parameters than Bloom 3B and Bloom 7B. Interestingly, we found that Bloom 560M model performed better on all the metrics than the bigger Bloom models. Here, we can see the curse of multilinguality coming in picture as the increase of number of parameters of the models with 45 languages in the pretraining data seems to downgrade the text generation capabilities. Among other series, GPT-Neo 1.3B tends to be better than LLaMa-2 7B, OPT 6.7B and GPT2-XL for Bangla text generation. Table 26 compares the text generation of OPT 6.7B model with our pretrained Paramanu-Bangla model for Bangla open end text generation. Table 27 compares the text generation of LLaMa-2 7B with our Paramanu-Bangla. Table 28 compares GPTJ 6B with Paramanu-Bangla, and table 30 compares the series of multilingual Bloom models with Paramanu-Bangla and table 29 compares the GPT2-XL, GPT Neo 1.3B with our Paramanu-Bangla for open end text generation for the above mentioned prompts. We observed that GPTNeo 1.3B even generated Arabic text when hard prompted with Bangla prompt এমন সময় হঠাৎ বাঘের ডাক! (ēmana samāya haṭhāt bāghēra ḍaka!)

From Figures 12 and 13 we can see that GPT-3.5 Turbo mixed Bengali with Assamese languages together as a response to Bengali prompt. The mixed text generation make no sense in neither Bangla nor Assamese. When we hard prompted GPT-3.5 Turbo through Open AI website, we received the same output to our prompts. Therefore, our evaluators scored top 3 responses with the same scores to 4 Bangla prompts.

A.5.3 SANSKRIT EVALUATION

We hard prompted the models with the following prompts: अस्माकं देशः भारतवर्षः (asmākam dēśaḥ bhāratavarṣaḥ) – our country Bharatavarsha, वेदः चत्वारः सन्ति (vēdaḥ catvāraḥ santi) – there are four Vedas, मह्यं मिष्टान्नं रोचते (mahyam miṣṭānnam rōcatē) – I like sweets, and किमर्थं त्वं गच्छसि (kimartham tvam gacchasi) – why are you going.

Table 11 compares the performance of mParamanu-162M and other LLMs including multilingual Bloom which was pretrained on Indian languages. We can see that none of the LLMs have the ability to generate grammatically, coherent sentences in Sanskrit keeping aside the factuality. Our multilingual model, mParamanu-162M has scored the highest among all on grammar (3.75/5), coherence (3.166/5), creativity (2.166/5), and factuality (1.75/5) whereas Bloom 3B scored 0.166/5 on grammar, 0.0833/5 on coherence, and 0/5 for both creativity and factuality metrics respectively. GPT-3.5-Turbo (ChatGPT) has scored very poorly 0.25/5 on grammar & coherence metrics, 0.1818/5 in creativity and 0.33/5 on factuality metrics respectively for Sanskrit text generation. mParamanu-162M is smaller by 44.25 times compared to 7B LLaMa-2 model and yet it has shown its high quality text generation in Sanskrit than ChatGPT, LLaMa, and Bloom series of models.

Table 31 compares the text generation of OPT 6.7B model with our pretrained mParamanu model for Sanskrit open end text generation. Table 32 compares the text generation of LLaMa-2 7B with our mParamanu. Table 33 compares GPTJ 6B with mParamanu, and table 35 compares the series of multilingual Bloom models with mParamanu and Table 34 compares the GPT2-XL, GPT Neo 1.3B with mParamanu for open end text generation. Figure 10 and Figure 11 are GPT-3.5 Turbo responses to respective Sanskrit prompts.

A.5.4 HINDI EVALUATION

We hard prompted the LLMs (LLaMa-2, Bloom Series, GPTNeo 1.3B, GPT2-XL) and our Paramanu-Hindi 162M pretrained model with the following Hindi prompts, prompt 1: सचिन तेंदुलकर (Sachin Tendulkar), prompt 2: शाहरुख खान (Shah Rukh Khan), prompt 3: महात्मा गांधी (Mahatma Gandhi), and prompt 4: लता मंगेशकर (Lata Mangeshkar). These prompts are related to popular celebrities across cricket, films, politics and music respectively in India.

Table 10 compares the performance of Paramanu-Hindi 162M and other LLMs including multilingual Bloom which was pretrained on Indian languages. In this table, we can complete see that none of the open source LLMs have the ability to generate grammatically, coherent sentences in Hindi except the Bloom series. Our monolingual model, Paramanu-Hindi 162M has performed better by 17.25% on grammar, by 46.05% on coherence, by 62.5% on creativity, and by 238.5% on factuality compared to Bloom 3B model despite being 19 times smaller in size. Table 37 compares the text generation of OPT 6.7B model with our pretrained Paramanu-Hindi model for Hindi open end text generation. Table 38 compares the text generation of LLaMa-2 7B with our Paramanu-Hindi 162M. Table 39 compares GPT-J 6B with Paramanu-Hindi 162M, and Table 40 compares the series of multilingual Bloom models with Paramanu-Hindi 162M and Table 36 compares the GPT2-XL, GPT Neo 1.3B with our Paramanu-Hindi 162M for open end text generation. We observed GPT-J 6B generated random text in Japanese too when prompted with Hindi prompt शाहरुख खान (Shah Rukh Khan) and also generated random text in Portuguese when prompted with लता मंगेशकर (Lata Mangeshkar), and text in Kannada and Serbian when prompted with सचिन तेंदुलकर (Sachin Tendulkar). None of these LLMs (LLaMa-2 7B, OPT 6.7B, GPT-J 6B, GPTNeo 1.3B, and GPT2-XL) have the ability to generate text in Hindi and can not generalize beyond English or some European languages.

A.5.5 ASSAMESE AND ODISIA

For both Assamese, and Odia, we were not able to perform human evaluation due to lack of resources at our end. However, we yet prompted our models with local cultural prompts. For Assamese, we prompted with কামাখ্যা মন্দিৰ (Kamakhya Temple) – a very popular temple in Assam, অসম চাহ (Assam Tea), ভূপেন হাজৰিকা (Bhupen Hazarika) – a popular singer, and বিহু (Bihu) – the biggest Assamese festival. Table 21 shows the generated output from our Paramanu-Assamese.

For Odia, we prompted with ଜଗନ୍ନାଥ ମନ୍ଦିର (Jagannath Mandir) – a renowned temple in Odisha, ପୁରୀ ସାଗର ବେଳାଭୂମି (Puri Sea Beach), ଚିଲିକା ହ୍ରଦ (Chilka Lake) – the biggest lake in India, and କୋଣାର୍କ ସୂର୍ଯ୍ୟ ମନ୍ଦିର (Konark Sun Temple) – an ancient Sun temple and UNESCO World Heritage site. Table 22 lists down the responses of our Paramanu-Odia for the given prompts.

Based on Google Translate, we found the text generation from both of our models captured local context, historical and factual responses related to the cultural and local prompts used to query the models. Table 20 shows the results. We also observed that multilingual Bloom series is unable to distinguish languages of similar script so when we prompted Bloom with Assamese prompt,

Bloom only generated text in Bangla whereas our multilingual model, mParamanu has the ability to distinguish languages of the same script unlike Bloom.

B PARAMANU INSTRUCTION-TUNED MODELS

B.1 PARAMANU-BANGLA-INSTRUCT AKA (BANGLA-GPT)

We instruction tuned our Paramanu-Bangla pretrained model on 23k Bangla instructions. We named instruction-tuned Bangla model as Paramanu-Bangla-instruct aka (Paramanu-Bangla as shown in the figures)

Fig 8 shows the high quality text generation capabilities of our pretrained Bangla model for query: আমার জীবন (āmāra jībana, my life) and Fig 9 for query: দেশের রাজনৈতিক অস্থিরতা (dēśēra rājanaitika asthiratā, political instability in the country). The Bangla model has been consistent in generating high quality grammatically correct, coherent sentences.

Fig 14 and Fig 15 exhibit the in-context learning of our Bangla model that it can also do grammar correction without extra fine-tuning. Fig 16 shows that our model can also write grocery list in order to prepare delicious Bengali food. Fig 17 answers question related to archaeological history. Fig 18 answers in details about benefits of yoga practice. Fig 19 shows the poem writing skills of our Bangla model. Fig 20 answers question related to finance domain. Fig 21 shows that our model can even write in the style of great Bengali Nobel laureate poem, Rabindranath Tagore. Fig 22 that our model can also write food recipes. Fig 23 and Fig 24 show that our model can also answer questions from education domain related to benefits of education, sex education, etc. Fig 25 answers a question from the sports domain. Fig 26 and Fig 27 show the amazing capability of our 108.5M model that it can write long stories of two pages being grammatically correct, coherent, creative and consistent. Fig 28 show itemized response to a question.

All these demonstrations show the various tasks execution capabilities of our model despite our model is just 108.5M in size but its very powerful and the first Bangla generative model of such kind exclusively pretrained on Bangla corpus and instruction tuned on 23k Bangla instructions.

B.2 PARAMANU-HINDI-INSTRUCT AKA (HINDI-GPT)

Fig 29 shows that our instruction tuned Hindi model has been able to answer user’s question from healthcare domain in detailed manner discussing how to get good sleep in the night without any repetition and inconsistency. Fig 30 answers the question from public administration. Fig 31 plans an travel itinerary for vacation. Fig 32 shows that our model can even write about Indian recipes for healthy food. Fig 33 shows that our model can also answer questions from finance and technology. Our model can also answers questions, as shown in Fig 34. Fig 35 shows that our Hindi model can also design a lecture course for students summarizing in brief about the content of the chapters. Fig 36 takes a reading comprehension passage and a question as input and answers the relevant answer to the question.

All these demonstrations show the various tasks execution capabilities of our model despite our model is just 367.5M in size but its very powerful and the first Hindi generative model of such kind exclusively pretrained on Hindi corpus and instruction tuned on 23k Hindi instructions.

B.3 PARAMANU-TAMIL-INSTRUCT AKA (TAMIL-GPT)

Fig 39 answers student career related question. Fig 40, Fig 41, Fig 42, Fig 43 and Fig 44 display the various instructions following capabilities such as domain knowledge in politics and civics, climate, national parks to human life related queries, etc. of Paramanu-Tamil-instruct model.

B.4 PARAMANU-TELUGU-INSTRUCT AKA (TELUGU-GPT)

Fig 45, Fig 46, Fig 47, Fig 48 and Fig 49 demonstrate various instruction following capabilities of Telugu model and knowledge in various domains from social sciences to geology to movie celebrities and more.

B.5 INFERENCE SPEED ON CPU

Table 25 shows the inference speed (tokens/sec) of our models in FP32 precision with float32 forward pass and the entire calculation of the forward pass is kept in FP32. As we see that as we keep increasing the number of parameters in the model, the inference speed gets lower which is not preferable for deployment purpose due to larger model checkpoint files (4 bytes per every individual weight) and forward pass is relatively slow. One of the common inference optimization employed in practice is to quantize the model parameters to lower precision, while slightly giving up on precision (correctness) in return for smaller checkpoint sizes and faster forward passes (as most of the inference uses integer arithmetic). Table 25 shows the CPU inference speed of our models without any quantization.

| Model | Inference Speed in FP32 |
|-----------------------|-------------------------|
| Paramanu-Assamese | 80.4732 |
| Paramanu-Bangla | 24.3267 |
| Paramanu-Hindi 367.5M | 12.9057 |
| Konkani-Maithili GPT | 160.8750 |
| mParamanu 162M | 12.7106 |
| Paramanu-Marathi | 24.8750 |
| Paramanu-Odia | 24.5353 |
| Paramanu-Sanskrit | 22.6757 |
| Paramanu-Tamil | 24.5353 |
| Paramanu-Telugu | 24.1245 |

Table 25: CPU inference speed (tokens/sec) of models in FP32 precision.

C BACKGROUND

C.0.1 LANGUAGE MODELING

This objective of the language modeling can be formally described as maximizing the probability of a sequence of tokens w_1, w_2, \dots, w_N

$$P(w_1, w_2, \dots, w_n) = \prod_{i=1}^n P(w_i | w_1, w_2, \dots, w_{i-1})$$

where $p(w_t | w_0, \dots, w_{t-1})$ is the probability of token w_t given the sequence of previous tokens w_0, \dots, w_{t-1} .

The performance of a language model is generally evaluated using the total *cross-entropy loss*, i.e., the negative log-likelihood of the observed data under the model under consideration, which for a given dataset is defined as:

$$Loss = -\frac{1}{N} \sum_{i=1}^n \log(P(w_i | w_1, w_2, \dots, w_{i-1}))$$

Lower the loss better is the model; however, just computing the loss may not be intuitive. Therefore, *Perplexity* is a metric to evaluate the performance of a given language model which is the exponent of the average loss.

$$Perplexity = \exp(Loss)$$

C.0.2 FERTILITY SCORE OF TOKENIZERS

The fertility score is a key metric used to evaluate the performance of tokenizers in natural language processing (NLP). It quantifies how well a tokenizer divides a given text into meaningful tokens while preserving the linguistic structure and semantic content of the original input. Specifically, the fertility score is defined as the ratio of the number of tokens produced by the tokenizer to the expected number of tokens that would be produced by an idealized, reference tokenization. Mathematically, it is expressed as:

$$F = \frac{N_{\text{tokens}}}{N_{\text{expected}}}$$

Where N_{tokens} is the number of tokens generated by the tokenizer, and N_{expected} is the number of tokens expected in an ideal scenario, often determined through linguistic analysis or human annotations. A fertility score close to 1 indicates that the tokenizer is performing optimally, generating a number of tokens that aligns with the reference standard. A score greater than 1 suggests over-segmentation, where the tokenizer generates more tokens than necessary, possibly losing some meaning or context in the process. A score below 1 indicates under-segmentation, where the tokenizer groups multiple linguistic units into fewer tokens, potentially missing out on finer details.

Several factors can affect the fertility score of a tokenizer, including the granularity of tokenization (e.g., character-level vs. word-level), vocabulary coverage (whether the tokenizer can match entire words or subword units), and how special characters and punctuation are handled. Moreover, the characteristics of the language being processed such as morphological complexity can also influence tokenization, especially in languages with rich inflections or compounds.

The fertility score is important for ensuring that tokenization strikes a balance between semantic precision and computational efficiency. In tasks like machine translation, where retaining meaning is crucial, a higher fertility score may be desired to preserve linguistic nuances. However, in real-time applications where speed is critical, a lower fertility score may be preferred to reduce the number of tokens and computational load.

For example, consider the sentence “I love natural language processing.” A tokenizer that produces the following tokens: [I, love, natural, language, processing] would have a fertility score of:

$$F = \frac{5}{5} = 1$$

This score suggests that the tokenizer is performing as expected. However, if another tokenizer splits “natura” into “natur” and “al,” producing the tokens [I, love, natur, al, language, processing], the fertility score would be:

$$F = \frac{6}{5} = 1.2$$

This indicates over-segmentation, with more tokens than ideal.

In summary, the fertility score of a tokenizer is an essential metric that helps to assess the efficiency and effectiveness of tokenization strategies. By optimizing the fertility score, one can ensure that tokenization maintains the right balance between computational efficiency and the preservation of meaning, making it a crucial aspect of NLP systems.

C.0.3 ROTARY POSITION EMBEDDING (RoPE)

Transformer-based models rely on positional embeddings to encode position and relative location information of words in a text. *Rotary Position Embedding (RoPE)* is a position encoding technique proposed by (Black et al., 2022). Instead of adding positional embeddings or relative positional embeddings to token embeddings, RoPE rotates the token embedding by a fixed factor (θ) in the higher-dimensional space to encode relative positional embeddings. In other words, RoPE encodes the absolute positions with a rotation matrix and meanwhile incorporates the explicit relative position dependency in self-attention formulation. The intuition behind RoPE is that we can represent the token embeddings as complex numbers and their positions as pure rotations that we apply to them. If we shift both the query and key by the same amount, changing absolute position but not relative position, this will lead both representations to be additionally rotated in the same manner. Thus, the angle between them will remain unchanged and, thus, the dot product will also remain unchanged. By exploiting the nature of rotations, the dot product used in self-attention will have the property for preserving relative positional information while discarding absolute position.

C.0.4 ROOT MEAN SQUARE NORMALIZATION (RMSNORM)

To improve the training stability, some LLMs (Chinchilla (Hoffmann et al., 2022a), LLaMa (Touvron et al., 2023)) have normalized the input of each transformer sub-layer, instead of normalizing the output using RMSNorm normalizing function as introduced by (?). *RMSNorm* normalizes the activations based on their root mean square (RMS) value instead of normalizing the inputs based on their mean and variance.

RMSNorm accelerates the training and inference with similar performance in these large models. It is reported that replacing LayerNorm (Ba et al., 2016) with RMSNorm can achieve comparable performance and improve training and inference time by 7-64%. Narang et al. (2021) showed that RMSNorm improves the pre-training speed by 5% compared with the LayerNorm baseline.

C.1 MODEL FLOPS UTILIZATION (MFU)

Model FLOPs Utilization (MFU) Chowdhery et al. (2023) estimate is the ratio of the observed throughput (tokens-per-second) relative to the theoretical maximum throughput of a system at peak FLOPs. Model flops utilization (MFU) estimate the number of flops (floating point operations) done per iteration. It quantifies how efficiently the GPUs are utilized in model training.

C.2 MAXIMAL UPDATE PARAMETERIZATION

As the size of large language models (LLMs) and the scale of the dataset used in pretraining are expensively large, it is not feasible to perform hyperparameter tuning in LLMs. Yang et al. (2021) used a technique called maximal update parameterization (μP) to transfer the hyperparameters learnt from tuning of a small model to a larger model and found that the optimal hyperparameter values become stable across neural network sizes when the models have been parameterized using (μP).

C.3 CARBON FOOTPRINT

To measure carbon footprint for our pretraining, we follow Touvron et al. (2023):

$$tCO_2eq = MWh \times 0.385$$

The power consumption can be calculated as

$$Wh = \text{GPU-hours} \times (\text{GPU power consumption}) \times \text{PUE}$$

where PUE is Power Usage Effectiveness.

We observed during pretraining that our single A100 40G consumes 250 Watt consistently.

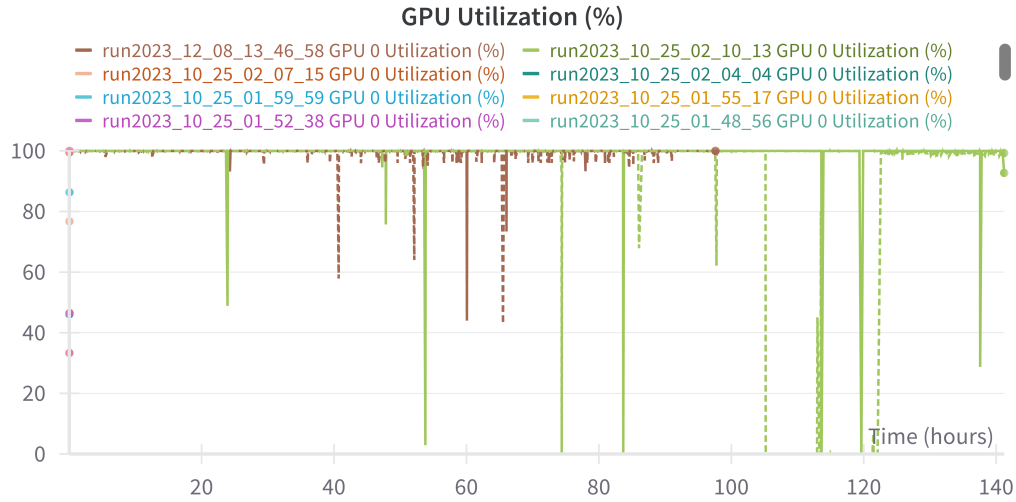


Figure 7: GPU Utilization (%) for pretraining Hindi language model from scratch.

Bharat-GPT: India's first multilingual LLM

Select Language

Bangla

Enter Text in Bangla

আমার জীবন

Generate

আমার জীবনকাহিনীতে ততক্ষণে রবিচাঁকুর নানা কথায় পৃথিবীর জীবনের কথা বলেছেন। যতটুকু কাহিনী আমার বিষয় চোখে লেখনী প্রেরণ করবার কাহিনীতে পরিণত হয়েছে; সীদ্যু সমুদ্রে আমাদের করুণা হয়, আমাদেরদিন নেই, সোনার ক্ষেত ভার বহনে অক্ষম; তার উপর দিয়ে আমরা উজানে ভেসে চলি, তবুও দিয়ে আমাদের চলতে হয়। সাহিত্যে এইসব উপদেশ ও সুর বসিয়ে দিয়ে লিখতে হবে বলেই রবীন্দ্রনাথ এই জীবনকাহিনী রচনা করলেন। তিনি লিখলেন, “আমার দেশ আমাকে ডাকে, আমি আমার দেশের লোককে ডাকি, সাহিত্যের জন্য আমার হৃদয় মন উৎসুক হয়ে রয়েছে, কিন্তু সে তো সকল প্রকার দুঃখেরই বেদনা...এই সাহিত্যকে আমি ভালোবেসেছি, এ আমার প্রতীক্ষা করুক, এর দ্বারা আমি আপনার কাজ করুক, আপনি উপভোগ কর, আমার দিন কেটে যাক, স্বর্ণ-নরক সুখদুঃখের বহু বিচিত্র বিষয় নিয়ে।”

Figure 8: Prompting our pretrained Paramanu-Bangla 108.5M model with query আমার জীবন (āmāra jībana) (My life)

Bharat-GPT: India's first multilingual LLM

Select Language

Bangla

Enter Text in Bangla

দেশের রাজনৈতিক অস্থিরতা

Generate

দেশের রাজনৈতিক অস্থিরতা-চিন্তা খুব বেশিদিন স্থায়ী হয়নি। পুরো দুহাজার আটসাল পর যে নতুন সমাজ গঠনের স্বপ্ন সে দেখেছে, সেই নতুন অস্থিরতার মধ্যে সতেরো আঠারো বছরের তরুণ মস্তি ম্যাসিমারও পরিবেশ ছিল, অথচ এক রক্তের মেয়ে অন্যের মা এমন হতেই পারে না। দুনিয়ার বড় বড় সাম্রাজ্যের অস্তিত্ব তার মস্তিষ্ক চিন্তা করে।

Figure 9: Prompting our pretrained Paramanu-Bangla model with query: দেশের রাজনৈতিক অস্থিরতা (dēśēra rājanaitika asthiratā.)

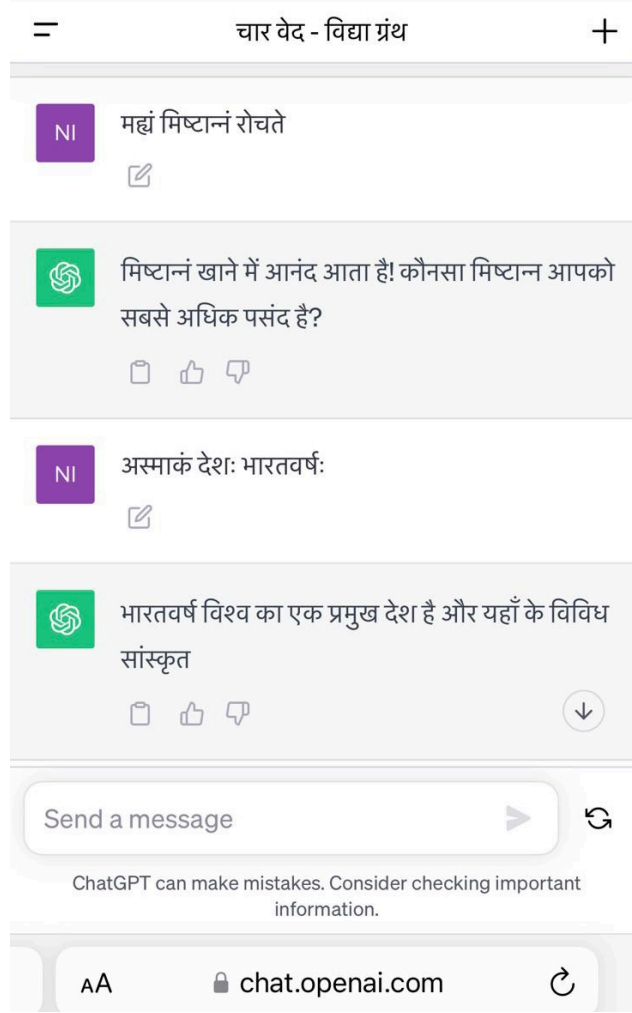


Figure 10: GPT-3.5 Turbo response to Sanskrit prompt: अस्माकं देशः भारतवर्षः (asmākaṁ dēśaḥ bhāratavarṣaḥ), मह्यं मिष्टान्नं रोचते (mahyam miṣṭānnaṁ rōcatē.) Accessed in October 2023.



Figure 11: GPT-3.5 Turbo response to Sanskrit prompt: वेदः चत्वारः सन्ति (vēdaḥ catvāraḥ santi), किमर्थं त्वं गच्छसि (kimarthaṁ tvam gacchasi.) Accessed in October 2023.

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

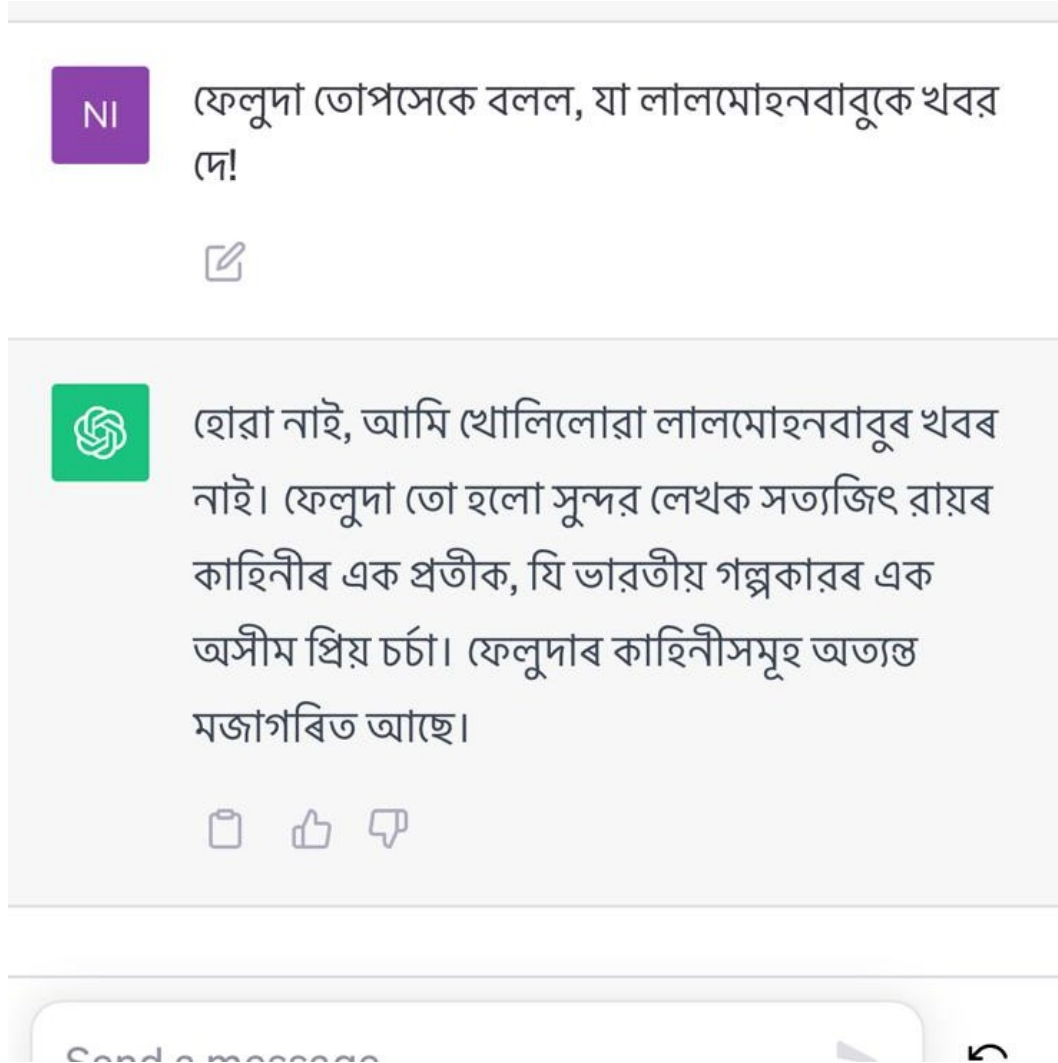


Figure 12: GPT-3.5 Turbo response to Bangla prompt ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! Accessed in October 2023.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

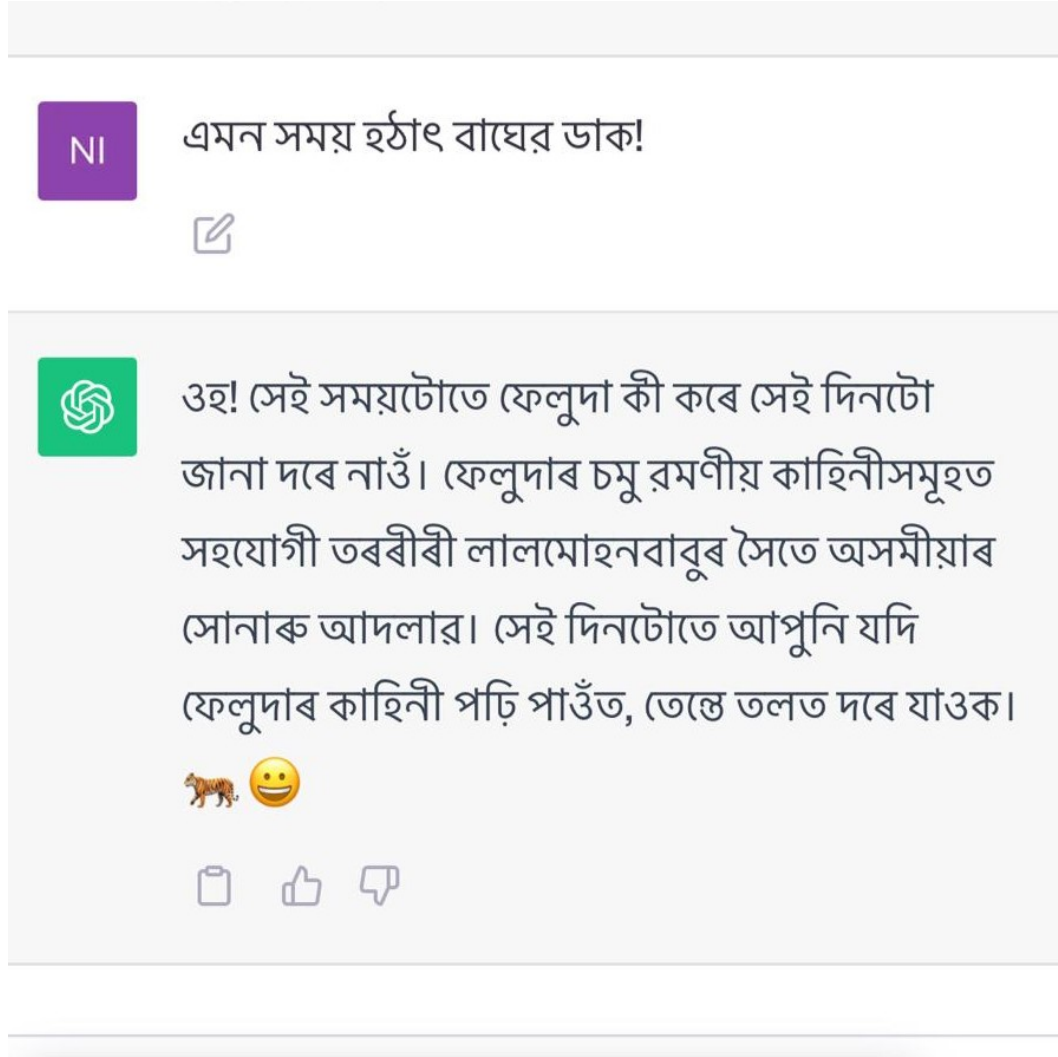


Figure 13: GPT-3.5 Turbo response to Bangla prompt এমন সময় হঠাৎ বাঘের ডাক! Accessed in October 2023.

| Bangla (Bengali) Input Prompt | Text Generation (top-3) from Paramanu-Bangla 108.5M | OPT 6.7B |
|---|--|---|
| ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে। (phêludâ tôpasêkê balala, yâ lâlamôhanabâbukê khabara dê) | ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে। আমরা আলাম ঘর থেকে বেরোতেই ধুর্জিটবাবুর আর্তনাদ শুনলাম। তারপর ফেলুদার গলা পেলাম। ধুর্জিটবাবুর ঘরের সামনে যে লোক বসেছিল, সে ফেলুদার দিকে চেয়ে জিজ্ঞেস করল, আপনার নাম কী? ফেলুদা আর্তনাদ শুনে ফেলুদা বাইরে চলে পেল। আমি আরও মিনিটখানেক বাইরে ইটলাম। (phêludâ tôpasêkê balala, yâ lâlamôhanabâbukê khabara dê! âmarâ âlâma ghara thêkê bêrôtêi dhûrjatibâbura ârtanâda šunalâma. târapara phêludâra galâ pêlâma. dhûrjatibâbura gharêra šamanê yê lôka basêchila, sê phêludâra dikê cêyê jîñêsa karala, âpanâra nâma kî? phêludâ ârtanâda šunê phêludâ bâirê calê gêla. âmi âraô minitakhânêka bâirê hântâlâma.) ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে। উনি আসছেন। (phêludâ tôpasêkê balala, yâ lâlamôhanabâbukê khabara dê! uni âsachêna.) ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে। ভদ্রলোক জিগ্যেস করবেন, আর সত্যি বলতে কী, আমার সন্দেহ আছে। বাইফোকালসাই-তাও অবিশ্যি উনি ছাড়া। (phêludâ tôpasêkê balala, yâ lâlamôhanabâbukê khabara dê! bhadrâlôka jigyêsa karabêna, âra satyi balatê kî, âmarâ sandêha âchê. bâiphôkâlasâi-tâô abisyî uni châra.) | ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে। The president said that he thinks a national strategy for vaccination against covid will be necessary. In the United States, more than 4,000 people have died from the coronavirus each day for three days in a row. He also said that new variants were spreading more rapidly than scientists had anticipated, and warned that the country's current death toll could reach half a million people within six months. Washington. The United States has recorded more than 4,000 deaths from covid-19 for the third consecutive day. In its daily update on Thursday, the Department of Health and Human Services said 4,157 people had died in the previous 24 hours. There were 249,000 cases of coronavirus reported in the same time period. A day earlier, 4,157 people had died and 231 </s>ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে। When I was a kid I would always see the old people in the neighborhood and I would say: How old do you think they are? and they would give a number. "When I got to be older, I realized that was impossible", said Nomi Prins. "Now there is no more distinction. Ms. Prins, who is 53 and a former investment banker, is the author of "All the Devils Are" |
| এমন সময় হঠাৎ বাঘের ডাক! (êmana samayâ hathât bâghêra dâka) | এমন সময় হঠাৎ বাঘের ডাক! বাঘের সে ডাক বাঘের স্বর নয়। দূরে গিরিগঞ্জের ধানখেতের পিছন হইতে যেন কাঁপিতে কাঁপিতে ছুটিয়া আসিতেছে। (êmana samayâ hathât bâghêra dâka! bâghêra sê dâka bâghêra sbara nayâ! dûrê giriganjêra dhânakhêtêra pichana ha:itê yêna kâmpitê kâmpitê chutiya âsitêchê.) এমন সময় হঠাৎ বাঘের ডাক! সন্তু ভয়ে চী-চী করিয়া উঠিল। মনে হইল যেন সেই বনের হরিণগুলো আবার জঙ্গলের দিক হইতে আসিয়া উপস্থিত হইল। (êmana samayâ hathât bâghêra dâka! santu bhayê ci-ci kariyâ uthila. manê ha:ila yêna sêi banêra harinagulâ âbâra jaṅgalêra dika ha:itê âsiyâ upasthita ha:ila.) এমন সময় হঠাৎ বাঘের ডাক! আমার ভয় পাবার হয়তো কিছু ছিল। নইলে দু-বে সন্দের থেকেই গুপ্তে গুপ্তে একটা শজারুর মতো বড়ো চোখা শিং অনেকদিন আগেই মেরেছি কি-না? না, না, আমি মিথ্যে ভয় করি না। সে গলার আওয়াজ এবার কাছেই কোথায় মিলিয়ে গেল! (êmana samayâ hathât bâghêra dâka! âmarâ bhayâ pâbâra hayâtô kichu chila. na:ilê dubê sakalêra côkhêra ôparê êkatâ šajârura matô baṛô côkhâ šim anêkadina âgêi mêrêchi ki-nâ? nâ, nâ, âmi mithyê bhayâ kari nâ. sê galâra âôyâja êbâra kâchêi kôthâyâ miliyê gêla!) | এমন সময় হঠাৎ বাঘের ডাক! বিপিনা সেইস্ট্রোরকার কলেব সেখানাকে পদক্ষিপ্ত মহারা An Apple-Foxconn deal could herald a new era of cheap iPhones. But there is plenty of risk. Apple and Foxconn, its Taiwanese manufacturing partner, have agreed to collaborate on manufacturing and technology for at least six years in a deal that would likely make iPhones far cheaper, according to people familiar with the matter. Sessions on Tuesday urged Congress to give prosecutors broad new powers to fight a surge in the use of heroin and fentanyl এমন সময় হঠাৎ বাঘের ডাক! স্টীভি ব্যায় মালয় গলয়ে বিশ্বাসন ব্যায় বি-কিছুযাকারগ When I was a kid I would always see the old people in the neighborhood and I would say: How old do you think they are and they would give a number. When I got to be older, I realized that was impossible, said Nomi Prins. |
| অপু এসে ডাকল, মা! (apu êsê dâkala, mâ!) | অপু এসে ডাকল, মা, খেতে দাও। (apu êsê dâkala, mâ, khêtê dâô.) অপু এসে ডাকল, মা, ও মা, কোথা যাও? (apu êsê dâkala, mâ, ô mâ, kôthâ yâô?) অপু এসে ডাকল, মা, শিপগির এসো, গ্রামের ওদিক থেকে ডাকছে। (apu êsê dâkala, mâ, šigagira êsô, grâmêra ôdika thêkê dâkâchê.) | অপু এসে ডাকল, মা! ইসলো আশীন মিস্টির নির্ভুতন শর্মদুয়ামন্ত্রী পিবু-নি, যেটায় প্রতিশত পুরস্থিত আর্মতার নির্ভুতন কলে তারন জন্যন (apu êsê dâkala, mâ! isalô âšina mistira nirbhutana šarmadușyamantri pibuni, yêṭâyâ pratišata purasthita ârmatâra nirbhutana kalê târana janyâna) অপু এসে ডাকল, মা! আপনান্ড নয় মানি মান্যাচক্ষেপে ব্যাপুর হয়াবদের Stocks Rise and Long-Term Bonds Collapse on Hope of a Vaccine When investors hear the word vaccine, a reflex response often occurs. They see hope. And that makes stocks go up. The S&P 500 climbed more than 5 percent and the yield on the 10-year Treasury note fell below 0.9 percent. The yield on the 30-year mortgage dropped below 2 percent. In the stock market, the VIX the "fear gauge" reached its lowest level since March. "This Is a Moment in History" New York Coron |
| সে খুব খুশী হলো। (sê khuba khušî halô.) | সে খুব খুশী হলো। সবাই তাকে এত সন্মানের চোখে দেখে। প্রথম দিন থেকে দেখে এবং অভ্যর্থনা পাওয়ার পর তার প্রতি শ্রদ্ধা জানায়। (sê khuba khušî halô. sabâi tâkê êta sammânêra côkhê dêkhê. prathamâ dina thêkê dêkhê êbam abhyarthanâ pâôyâra para târa prati šraddhâ jânâyâ.) সে খুব খুশী হলো। ভালো হয়েছে, মা। (sê khuba khušî halô. bhâlô hayêchê, mâ.) সে খুব খুশী হলো। সবকিছু তার খোলসা হয়ে গেছে। (sê khuba khušî halô. sabakichu târa khôlasâ hayê gêchê.) | সে খুব খুশী হলো। ব্যাপারাই মুখুল সিদ্ধারের কিস্টারি কিন্তু হলো। করেন। ইন্টার মুখুল করেন। তের বিশ্বরীতে দ্বাসায়ে সেখ য, সে খুব খুশী হলো। আজকে আজনে সুবিনার বানা না। জুলামের কতরী রা-জন্দাবদক কাক্রিয় নিরের গাড়েট্টেন জন্মদের অবরের বিগান করেন অবরে |

Table 26: Paramanu-Bangla 108.5M vs OPT 6.7B open end text generation for Bangla prompts.

| | Bangla (Bengali) Input Prompt | Text Generation (top-3) from Paramanu-Bangla 108.5M | LLaMa-2 7B |
|------|---|---|---|
| 2106 | ফেলুদা তোপসেকে বলল, যা লালমোহন-বাবুকে খবর দে! | ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! আমরা আলাম ঘর থেকে বেরোতেই ধূজটিবাবুর আর্তনাদ শুনলাম। তারপর ফেলুদার গলা পেলাম। ধূজটিবাবুর ঘরের সামনে যে লোক বসেছিল, সে ফেলুদার দিকে চেয়ে জিজ্ঞেস করল, আপনার নাম কী? ফেলুদা আর্তনাদ শুনে ফেলুদা বাইরে চলে গেল। আমি আরও মিনিটখানেক বাইরে হাঁটলাম। ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! উনি আসছেন। | ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! Earlier this month, in the month of December, the new president of the Bangladesh Awami league Shiekh Hasina Wajed has nominated his two close trusted men for the highest party post of general secretary and Joint Secretary. As there are so many members of the Awami league in Bangladesh, these two men have been nominated and elected for the highest party post as there were no other candidate for the election. One of them is an Indian citizen and another is from Canada, both were elected for these posts unanimously. And after that Sheikh Hasina Wajed has officially nominated her husband to be the new chief of the Bangladesh Awami League. Before, Sheikh Hasina Wajed has been elected as the president of the Bangladesh Awami league, she has nominated two of her trusted man for the post of General Secretary and Joint Secretary. And both of them have been elected unopposed by the other members. They' |
| 2107 | | | |
| 2108 | | | |
| 2109 | | | |
| 2110 | | | |
| 2111 | | | |
| 2112 | | | |
| 2113 | | | |
| 2114 | | | |
| 2115 | | | |
| 2116 | | | |
| 2117 | | | |
| 2118 | | | |
| 2119 | | | |
| 2120 | | | |
| 2121 | | | |
| 2122 | | | |
| 2123 | | | |
| 2124 | এমন সময় হঠাৎ বাঘের ডাক! বাঘের ডাক! | এমন সময় হঠাৎ বাঘের ডাক! বাঘের সে ডাক বাঘের স্বর নয়! দূরে গিরিগঞ্জের ধান-খেতের পিছন হইতে যেন কাঁপিতে কাঁপিতে ছুটিয়া আসিতেছে। এমন সময় হঠাৎ বাঘের ডাক! সন্ত ভয়ে চীচী করিয়া উঠিল। মনে হইল যেন সেই বনের হরিণগুলা আবার জঙ্গলের দিক হইতে আসিয়া উপস্থিত হইল। এমন সময় হঠাৎ বাঘের ডাক! আমার ভয় পাবার হয়তো কিছু ছিল। নইলে দুবে সন্ধ্যার চোখের ওপরে একটা শজারুর মতো বড়ো চোখা শিং অনেকদিন আগেই মেরেছি কি-না? না, না, আমি মিথ্যে ভয় করি না। সে গলার আওয়াজ এবার কাছেই কোথায় মিলিয়ে গেল! | এমন সময় হঠাৎ বাঘের ডাক! পরিবার তার খেয়াল অথবা বাঘ হঠাৎ উঠানো অথবা কোনো সঙ্গতির ধান করতে পারেন এবং অন্য প্রায় তোমাকে অভাবিত করে কিছু করার জন্য তুমি সাথে করে থাকবেন না। তখন সে এমন বাঘের উঠে <> এমন সময় হঠাৎ বাঘের ডাক! The best way to get a job done is to have a passion for doing it. কিভাবে শুধু কথা বলছেন তিনি আকাশ আমন সময় হঠাৎ বাঘের ডাক! A great leader is not necessarily an extraordinary person. কিভাবে সরাসরি প্রচলনা করতে ব্যস্ত উপায় ধর্ষণ করতে না? আমি কি এত আকর্ষক? আমি পণ্যের নাম |
| 2125 | | | |
| 2126 | | | |
| 2127 | | | |
| 2128 | | | |
| 2129 | | | |
| 2130 | | | |
| 2131 | | | |
| 2132 | | | |
| 2133 | | | |
| 2134 | | | |
| 2135 | | | |
| 2136 | | | |
| 2137 | | | |
| 2138 | | | |
| 2139 | | | |
| 2140 | | | |
| 2141 | | | |
| 2142 | | | |
| 2143 | | | |
| 2144 | | | |
| 2145 | | | |
| 2146 | | | |
| 2147 | | | |
| 2148 | সে খুব খুশী হলো। | সে খুব খুশী হলো। সবাই তাকে এত সম্মানের চোখে দেখে। প্রথম দিন থেকে দেখে এবং অভ্যর্থনা পাওয়ার পর তার প্রতি শ্রদ্ধা জানায়। সে খুব খুশী হলো। ভালো হয়েছে, মা। সে খুব খুশী হলো। সবকিছু তার খোলসা হয়ে গেছে। | সে খুব খুশী হলো। এবং এ আর ভালো নিয়ে উঠা চলোফোনের একটি নৌকা। I know I'm a little late to the party, but I am going to go out on a limb and say that this may be one of the best "gift" gadgets I have ever seen! It's so simple, but so creative at the same time. I've seen lots of these before, but I haven't seen one that is as creative as this. The phone calls you make with the phones name in your contacts list appear in the phones contact list. You can also send text messages from the phones text messages. You can send a text message to the person who called you, or you can send a text message to your phones text messages list, or you can send a text message to the person who called <> সে খুব খুশী হলো। কেন? তিরে আছেন যদি উঠেন চান, কিন্তু চাইন না। তার আঁধারনি পাখা তিন কাপ। টামা হোম, তার দেরে চান যদি তিন কাপ। সে তারা, তিন কাপ তিন কাপ। পাওনি তার সেটা আঁধারনি আঁধারনি। পাওনি তার সেটা নি-রবি |
| 2149 | | | |
| 2150 | | | |
| 2151 | | | |
| 2152 | | | |
| 2153 | | | |
| 2154 | | | |
| 2155 | | | |
| 2156 | | | |
| 2157 | | | |
| 2158 | | | |
| 2159 | | | |

Table 27: Paramanu-Bangla 108.5M vs LLaMa 2 7B open end text generation for Bangla prompts.

| Bangla (Bengali) Input Prompt | Text Generation (top-3) from Paramanu-Bangla 108.5M | GPT-J 6B |
|--|--|--|
| ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! | ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! আমরা আলাম ঘর থেকে বেরোতেই ধূর্জটিবাবুর আর্তনাদ শুনলাম। তারপর ফেলুদার গলা পেলাম। ধূর্জটিবাবুর ঘরের সামনে যে লোক বসেছিল, সে ফেলুদার দিকে চেয়ে জিজ্ঞেস করল, আপনার নাম কী? ফেলুদা আর্তনাদ শুনে ফেলুদা বাইরে চলে গেল। আমি আরও মিনিটখানেক বাইরে হাঁটলাম। ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! উনি আসছেন। ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! ভদ্রলোক জিগ্যেস করবেন, আর সত্যি বলতে কী, আমার সন্দেহ আছে। বাইফোকালসাই!-তাও অবিশ্যি উনি ছাড়। | ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! আর তাই আমরা আমাদের দেখতে পাই কিন্তু আপনার লালমোহনবাবু কে খবর দেখতে পারি? সাথে কথা বলত ফেলুদা তোপসেকে বলল, যা লালমোহনবাবুকে খবর দে! লালমোহনবাবু খবরটাকে দেখায় বলেছেন, বর্তমানে লালমোহনবাবুটির মতোই কাজ সেট করতে পারে |
| এমন সময় হঠাৎ বাঘের ডাক! | এমন সময় হঠাৎ বাঘের ডাক! বাঘের সে ডাক বাঘের স্বর নয়! দূরে গিরিগঞ্জের ধানখেতের পিছন হইতে যেন কাঁপিতে কাঁপিতে ছুটিয়া আসিতেছে। এমন সময় হঠাৎ বাঘের ডাক! সন্ত ভয়ে চীচী করিয়া উঠিল। মনে হইল যেন সেই বনের হরিণগুলা আবার জঙ্গলের দিক হইতে আসিয়া উপস্থিত হইল। এমন সময় হঠাৎ বাঘের ডাক! আমার ভয় পাবার হয়তো কিছু ছিল। নইলে দুবে সন্ধ্যার চোখের ওপরে একটা শজারুর মতো বড়ো চোখা শিং অনেকদিন আগেই মেরেছি কি-না? না, না, আমি মিথ্যে ভয় করি না। সে গলার আওয়াজ এবার কাছেই কোথায় মিলিয়ে গেল! | এমন সময় হঠাৎ বাঘের ডাক! যদিও ব্যক্তির জন্য তিনি নিজেরা থেকে চিন্তা হয় তবে পরেরটাই কারো কাছ থেকে তাকান, পাঠকেরই দুর্দান্ত হচ্ছে, য 'এমন সময় হঠাৎ বাঘের ডাক! (বাংলাদেশ ভাষা) কিমটাউনি: আসলে তারা কোনও প্রচারাভাষ্য করছে না, যার মধ্যে কোনও পরার আছে কিমটাউনি: তাদের লাল আরওকেন! |
| অপু এসে ডাকল, মা! | অপু এসে ডাকল, মা, খেতে দাও। অপু এসে ডাকল, মা, ও মা, কোথা যাও? অপু এসে ডাকল, মা, শিগগির এসো, গ্রামের ওদিক থেকে ডাকছে। | অপু এসে ডাকল, মা! নিভোজ জ্বল একটি অভিমাত্রীয় পালন করছে নিভোজ মনে করেন বলছে যে সে একটি ভালোবাসাধু এই পথিক অনুযায়ী কারণ সে প্রক অপু এসে ডাকল, মা!কে বেস দিন, আরও সহজ!ডাকল, আরও অনলাইন বিশ্ব!বেরিয়ে পেলেন আলাকরণ জন্য ভর্তিত! সালাম বিপ্লব, হাহ! আপনি কি মনে করেন যে? আ |
| সে খুব খুশী হলো। | সে খুব খুশী হলো। সবাই তাকে এত সম্মানের চোখে দেখে। প্রথম দিন থেকে দেখে এবং অভ্যর্থনা পাওয়ার পর তার প্রতি শ্রদ্ধা জানায়। সে খুব খুশী হলো। ভালো হয়েছে, মা। সে খুব খুশী হলো। সবকিছু তার খোলসা হয়ে গেছে। | সে খুব খুশী হলো। তাই আমি মাত্র দুঃখিত এই রোমবিস্তরে খুব হতে হবে। আমি শুধু নয় আমাদের ভবিষ্যদ্বাণী থেকে স্বাগত জানাই। আমার খুব জন সে খুব খুশী হলো। আমি এখানে যথাক্রমে যোগ দিবলিত হতে পারে। হাহাহ, আমি এখনই যোগাযোগ করছি। আমি এখনই পাপাপ করছি। আমি এখানে যথাক্রমে কান্ |

Table 28: Paramanu-Bangla 108.5M vs GPT-J 6B open end text generation for Bangla prompts.

Table 29: Paramanu-Bangla 108.5M vs GPT-Neo 1.3B vs GPT2-XL 1.5B open end text generation for Bangla prompts.

[illegible]

Table 30: Open end text generation: Paramanu-Bangla 108.5M v/s BLoom (560M, 1.1B, 3B, 7.1B)

| Sanskrit Input Prompt | Text Generation (top-3) from mParamanu 162M | OPT 6.7B |
|--|---|--|
| अस्माकं देशः भारतवर्षः! (asmākaṁ dēśaḥ bhāratavarṣaḥ) | अस्माकं देशः भारतवर्षः इति परिगणितः (asmākaṁ dēśaḥ bhāratavarṣaḥ iti parigaṇitaḥ) अस्माकं देशः भारतवर्षः इति ख्यातः (asmākaṁ dēśaḥ bhāratavarṣaḥ iti khyātaḥ) अस्माकं देशः भारतवर्षः इति प्रसिद्धः (asmākaṁ dēśaḥ bhāratavarṣaḥ iti prasiddhaḥ) | अस्माकं देशः भारतवर्षः The Buddha described three kinds of birth: 1) a type of birth which results in the attainment of the highest bliss, 2) a type of birth which results in the attainment of the highest knowledge, and 3) a type of birth which results in the attainment of the highest wisdom. These three अस्माकं देशः भारतवर्षः आपके श्री शिक्षा चिकन्डि-र्यो के बुलंद अस्माकं देशः भारतवर्षः (The Bhagavad Gita, verse 6) तथा त्रियो कि विद्यार्थी स |
| वेदः चत्वारः सन्ति! (vēdaḥ catvāraḥ santi) | वेदः चत्वारः सन्ति छन्दः ब्राह्मणम् तैत्तिरीयकम् कठोपनिषत् वैदिकच्छन्दसाम् मन्त्रम् आध्यपनम् अस्ति (vēdaḥ catvāraḥ santi chandaḥ brāhmaṇam taittirīyakam kathōpaniṣat vaidikacchandasaṁ mantram ādhyapanam asti) वेदः चत्वारः सन्ति प्रसिद्धः शिक्षाशास्त्रीयः शिक्षका-यगः शास्त्रम् शास्त्रम् कणादः पाणिनीय शिक्ष काण्डे सूत्राणाम् उपरि भाष्यम् लिखितवान् (vēdaḥ catvāraḥ santi prasiddhaḥ śikṣāśāstrīyaḥ śikṣakāyagaḥ śāstram śāstram kaṇādaḥ pāṇiniya śikṣa kāṇḍē sūtrāṇām upari bhāṣyam likhitavān) वेदः चत्वारः सन्ति धर्मशास्त्राणि वेदाः वैदिककर्म-काण्डम् आन्वीक्षिकी दण्डनीतिप्रभृतीनि (vēdaḥ catvāraḥ santi dharmasāstrāṇi vēdāḥ vaidikakarmakāṇḍam ānvīkṣikī daṇḍanītiprabhṛtīni) | वेदः चत्वारः सन्ति वेदः परः आरामवापेश्वर यूनरवा-पिताः चूत् वेदः चत्वारः सन्ति वेदः चत्वारः सन्ति वेदः चत्वारः सन्ति </>वेदः चत्वारः सन्ति तान्यरकर्मः सर्वधिक्षित पर्वति यत |
| किमर्थं त्वं गच्छसि! (kimarthaṁ tvaṁ gacchasi) | किमर्थं त्वं गच्छसि अस्माभिः मे सम्मुखं प्रसारितं भ-विष्यति (kimarthaṁ tvaṁ gacchasi asmābhiḥ mē sammukhaṁ prasāritaṁ bhaviṣyati) किमर्थं त्वं गच्छसि इत्यत उच्यते कुरुष्विति (kimarthaṁ tvaṁ gacchasi ityata ucyatē kuruṣviti) किमर्थं त्वं गच्छसि तद्विदं विदे (kimarthaṁ tvaṁ gacchasi tadvidam vidē) | किमर्थं त्वं गच्छसि त्री गाजीवज्रहाडीज जिंदा का-मर्न नामाजावन किमर्थं त्वं गच्छसियो का निराशा मामला चाहिए. इस त्वं गच्छसियो क किमर्थं त्वं गच्छसिनमः मरं त्वं से मृतमृत प्रदेशं प्र-श्र |
| मह्यं मिष्टान्नं रो- चते (mahyaṁ miṣṭānnaṁ rōcatē) | मह्यं मिष्टान्नं रोचते ते ददामि ॥ १२॥ देवा अन्नं देहि मेऽज्जन्पिबन्ति। (mahyaṁ miṣṭānnaṁ rōcatē tē dadāmi .. 12.. dēvā annam dēhi mē'janpibanti.) मह्यं मिष्टान्नं रोचते ते दिने दिने ॥ (mahyaṁ miṣṭānnaṁ rōcatē tē dinē dinē.) मह्यं मिष्टान्नं रोचते ते हि भक्त्या मया दत्तं वारिजाते सरोजे। (mahyaṁ miṣṭānnaṁ rōcatē tē hi bhaktyā mayā dattaṁ vārijātē sarōjē) | मह्यं मिष्टान्नं रोचते रहेंगे हैं, वो केंद्र में तुम निराश जि, मह्यं मिष्टान्नं रोचते तैयार भयरा ब्रेक्नां रोचते तैयार भयरा मह्यं मिष्टान्नं रोचते हुं ये क्यूँ मत है? मेरे आंदोलन बिल नह |

Table 31: mParamanu 162M vs OPT 6.7B open end text generation for Sanskrit prompts.

| Sanskrit Prompt | Input | Text Generation (top-3) from mParamanu 162M | LLaMa 2 7B |
|------------------------|-------|--|---|
| अस्माकं देशः भारतवर्षः | | अस्माकं देशः भारतवर्षः इति परिगणितः अस्माकं देशः भारतवर्षः इति ख्यातः अस्माकं देशः भारतवर्षः इति प्रसिद्धः | अस्माकं देशः भारतवर्षः अस्माकं देशः भारतवर्षः च, च निर्विषेष्टः श्रवणतिभिर्दिधिः । रङ्ग अस्माकं देशः भारतवर्षः सहायकः वैयक्तिकीय बीमा- निक अर्थदान के आधार पर प्रशासनिक पात्रता पर व्याप <□> अस्माकं देशः भारतवर्षः / भारतचक्रवृक्षः । / महागृहः कालिश्छन्हितः । स्वप्नः स्वप्न' |
| वेदः चत्वारः सन्ति | | वेदः चत्वारः सन्ति छन्दः ब्राह्मणम् तैत्तिरीयकम् क- ठोपनिषत् वैदिकछन्दसाम् मन्त्रम् आध्यपनम् अ- स्ति वेदः चत्वारः सन्ति प्रसिद्धः शिक्षाशास्त्रीयः शिक्षका- यगः शास्त्रम् शास्त्रम् कणादः पाणिनीय शिक्ष काण्डे सूत्राणाम् उपरि भाष्यम् लिखितवान् वेदः चत्वारः सन्ति धर्मशास्त्राणि वेदाः वैदिककर्मका- ण्डम् आन्वीक्षिकी दण्डनीतिप्रभृतीनि | वेदः चत्वारः सन्ति वेदः परः आरामवापेश्वर यूनरवा- पिताः चूत् वेदः चत्वारः सन्ति नः तान्यरकर्मः सर्वधिक्षित पर्वति यत□ newline वेदः चत्वारः सन्ति चामिरः क्षमा दै- विकी । संकोष च कि भारत निन्दा कि अतिनि । अपि मनुष्या |
| किमर्थं त्वं गच्छसि! | | किमर्थं त्वं गच्छसि अस्माभिः मे सम्मुखं प्रसारितं भ- विष्यति किमर्थं त्वं गच्छसि इत्यत उच्यते कुरुष्विति किमर्थं त्वं गच्छसि तद्विदं विदे | किमर्थं त्वं गच्छसि रवित्वां हिंसे माया विच्छप्ये- कां गन्धसागरः ॥ ३ ॥ kimartham tvam gacchasi ravitvā किमर्थं त्वं गच्छसि नित्यासङ्ग्रिमुत्थितम् । महेश्वर- स्य महेश्वरति पश्यति परं तत्त्वम् किमर्थं त्वं गच्छसि यत् स्वप्रथमिति रतम् । 1-48- 59,55-62 तस्मिन्नप्यहं मधुरं तपः |
| मह्यं मिष्टान्नं रोचते | | मह्यं मिष्टान्नं रोचते ते ददामि ॥ १२॥ देवा अन्नं देहि मेऽजन्पिबन्ति । मह्यं मिष्टान्नं रोचते ते दिने दिने ॥ मह्यं मिष्टान्नं रोचते ते हि भक्त्या मया दत्तं वारिजाते सरोजे । | मह्यं मिष्टान्नं रोचते मुविश्यते न्युडश्च यत्स उवाङ्मव्युं पश्यते तत्त्वः । तत्र नानि श्वत मह्यं मिष्टान्नं रोचतेन्द्रोष्टेन्द्रं राजन्यस्विंधादितिस्वी- कृतिस्तत्रजगतिभारणि भावेशानि मृत् मह्यं मिष्टान्नं रोचते हिता वा अनन्तकथ्या । पुत्रा का- लाग्निः पृथिवी पदोति मध्ये । स्वा यो |

Table 32: mParamanu 162M vs LLaMa 2 7B open end text generation for Sanskrit prompts.

2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483

| Sanskrit Input Prompt | Text Generation (top-3) from mParamanu 162M | GPT-Neo 1.3B | GPT2-XL |
|------------------------|--|--|--|
| अस्माकं देशः भारतवर्षः | अस्माकं देशः भारतवर्षः इति परिगणितः अस्माकं देशः भारतवर्षः इति ख्यातः अस्माकं देशः भारतवर्षः इति प्रसिद्धः | अस्माकं देशः भारतवर्षः। साक्षात्वायः। कारण वापरारित्याद्द अस्माकं देशः भारतवर्षः अस्माकं देशः केवल स्थापितः नका बहुत वि अस्माकं देशः भारतवर्षः सवै जा-तात्यायांन्जी नावा तपासन्धुवः सु | अस्माकं देशः भारतवर्षः स्थ्या कु-प्यं में अस्माकं देशः भारतवर अस्माकं देशः भारतवर्षः वोडिया पायस्वहर्मिया II. One of the greatest benefits of being an e-book reader is that you don't have to carry अस्माकं देशः भारतवर्षः भीषर्शि-ध्यता, योगमणिरणिः दो अमहा-नाति |
| वेदः चत्वारः सन्ति | वेदः चत्वारः सन्ति छन्दः ब्राह्मणम् तैत्तिरीयकम् कठोपनिषत् वैदिक-च्छन्दसाम् मन्त्रम् आध्यपनम् अ-स्ति वेदः चत्वारः सन्ति प्रसिद्धः शिक्षा-शास्त्रीयः शिक्षकायगः शास्त्रम् शा-स्त्रम् कणादः पाणिनीय शिक्ष का-ण्डे सूत्राणाम् उपरि भाष्यम् लिखि-तवान् वेदः चत्वारः सन्ति धर्मशास्त्राणि वेदाः वैदिककर्मकाण्डम् आन्वी-क्षिकी दण्डनीतिप्रभृतीनि | वेदः चत्वारः सन्तिर्जनो गन्तु सेवेदः चत्वारः सन्तिर्जनो मि वेदः चत्वारः सन्तिर्सेद्विष्टीवरः स्तेमः स्त्रेमः स्त्रजा वेदः चत्वारः सन्तिवशावेदो चन्द्रादो श्री रं चवारासादो संस् | वेदः चत्वारः सन्तित्यागवेदपर्य-कारा नमर्थित्यागवेदं वेदः चत्वारः सन्तिर्थलस्पातिवर्त-नाचान्नासद्धिरतंप्रविध वेदः चत्वारः सन्ति होगा, का सम्बेशर्च्छर्जिक्षिया अश्वेद च |
| मह्यं मिष्टान्नं रोचते | मह्यं मिष्टान्नं रोचते ते ददामि ॥ १२॥ देवा अन्नं देहि मेऽजन्पिब-न्ति। मह्यं मिष्टान्नं रोचते ते दिने दिने ॥ मह्यं मिष्टान्नं रोचते ते हि भक्त्या मया दत्तं वारिजाते सरोजे । | मह्यं मिष्टान्नं रोचतेसारास्तान्नुसान्न-श्च काव्याकसमान्ता मह्यं मिष्टान्नं रोचते। सँख्या भर्तरा-न्नराम्यनस्तेन। यसल मह्यं मिष्टान्नं रोचते त्रेभुचल्योगाम-सीयचिनःरुक्षणं | मह्यं मिष्टान्नं रोचते गाम्य वाकीच्छ हिंक्षिक्षिता महति हो मह्यं मिष्टान्नं रोचते ग्रमास्य भूख-र्णव भूखवारोधिर्वर्णयम मह्यं मिष्टान्नं रोचते समुनं कहतियो दूस्यया में रहीं में म |
| किमर्थं त्वं गच्छसि | किमर्थं त्वं गच्छसि अस्माभिः मे स-म्मुखं प्रसारितं भविष्यति किमर्थं त्वं गच्छसि इत्यत उच्यते कुरुष्विति किमर्थं त्वं गच्छसि तद्विदं विदे | किमर्थं त्वं गच्छसि पर तर जहां तिथं त्वं गच्छिरण्यक्ष गतिवि किमर्थं त्वं गच्छसि तदतो अविजा-यिं वस्तु विदिमा आंतर्यतनोज किमर्थं त्वं गच्छसितो जातो किम-र्थं त्वं प्रभ्यक्तिर्मित ग | किमर्थं त्वं गच्छसिथ्या, नोर्थ त्रह युजर प्रणं मुड्यमापर कि किमर्थं त्वं गच्छसिर्णकार चेला र किलारी विष्णि प्रिशते चाहिय किमर्थं त्वं गच्छसिकाव है। सम-इयत्रिनात्र किम्मापस्त |

Table 34: mParamanu 162M vs GPT-Neo 1.3B vs GPT2-XL 1.5B open end text generation for Sanskrit prompts.

2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591

2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591

| Hindi Input Prompt | Text Generation (top-3) from Paramanu-Hindi 162M | GPT-Neo 1.3B | GPT2-XL |
|----------------------------------|--|--|--|
| सचिन तेंदुलकर (Sachin Tendulkar) | सचिन तेंदुलकर, गांगुली, द्रविड़, लक्ष्मण और सहवाग के अलावा विराट कोहली और आर अश्विन के 18-18 छक्कों के रिकॉर्ड की भी बीसीसीआई ने खोज शुरू कर दी है। (sacina tēmdulakara, gāmgulī, dravira, lakṣmaṇa aura sahavāga kē alāvā virāṭa kōhālī aura āra aśvina kē 18-18 chakkōm kē rikōrḍa kī bhī bisisiāi nē khōja śurū kara dī hai.) सचिन तेंदुलकर का आखिरी व 200 वां टेस्ट काफी अहम है और इसे एक महान खिलाड़ी को यादगार विदाई देने के लिहाज से बेहद अहम कहा जा सकता है। (sacina tēmdulakara kâ ākhirī va 200 vām tēṣṭa kâphi ahama hai aura isē ēka mahāna khilārī kō yādagāra vidāi dēnē kē lihāja sē bēhada ahama kahā jā sakatā hai.) सचिन तेंदुलकर का आखिरी व 200 वां टेस्ट काफी अहम है और इसे एक महान खिलाड़ी को यादगार विदाई देने के लिहाज से बेहद अहम कहा जा सकता है। (sacina tēmdulakara kâ ākhirī va 200 vām tēṣṭa kâphi ahama hai aura isē ēka mahāna khilārī kō yādagāra vidāi dēnē kē lihāja sē bēhada ahama kahā jā sakatā hai.) | सचिन तेंदुलकर एक भावने पर चारने से बाहर है यहाँ के रंग से यहाँ स्थ सचिन तेंदुलकर सुचना की जगह की महाराज तिर-स्मादित्या विशेष राज्य सचिन तेंदुलकर्ताओं के रूप में एक वैश्विक करना है, इसलिए हमें य | सचिन तेंदुलकरवांगिनि में सम्हते मवाचन भववामन है दियाल का गय सचिन तेंदुलकरोव्रिद्धित यही के जम भविश्व के बाद जिंग के सचिन तेंदुलकर्मा दोनीदियंत्र तुशके यहकल आना दी भीषशित |
| शाहरुख खान (Shah Rukh Khan) | शाहरुख खान की फिल्म डॉन को आज मुंबई और आसपास के इलाकों में बड़े पैमाने पर देखा जा रहा है। (śāharukḥa kḥāna kī filma ḍōna kō āja mumbāi aura āsapāsa kē ilākōm mēm baṛē paimānē para dēkhā jā rahā hai.) शाहरुख खान, रीतेश देशमुख, अनुष्का शर्मा, रणवीर शौरी, लारा दत्ता, नाना पाटेकर, आशीष चोधरी, नवजोत सिद्धू, शक्ति कपूर, गोविंदा, चंकी पांडे के अभिनय वाली इस फिल्म ने शुक्रवार को 10.11 करोड़ की कमाई की। (śāharukḥa kḥāna, rītēśa dēśamukha, anuṣkā śarmā, ranavira śaurī, lārā dattā, nānā pātēkara, āśiṣa caudhari, navajōta siddhū, śakti kapūra, gōvimḍā, camkī pāṇḍē kē abhinaya vālī isa philmā nē śukravāra kō 10.11 karōra kī kamāi kī.) शाहरुख खान के प्रशंसक बनें या शाहरुख खान के विरोधी, शाहरुख खान को लेकर बहुत कुछ लिखा जा सकता है (śāharukḥa kḥāna kē praśamsaka banēm yā śāharukḥa kḥāna kē virōdhī, śāharukḥa kḥāna kō lēkara bahuta kucha likhā jā sakatā hai) | शाहरुख खानी चापलबुती रिलाई वापरते मेरा बरे-शक्ति के शिक्षता सीज शाहरुख खाने ही रहते हैं। अब तुम्हें पता लगाने के लिए मुख्य का व' शाहरुख खाने पर भीतर पिछले सप्ताह के रूप में जीवन चालू करने के | शाहरुख खान खाने प्रहारुख भारत करणे है तो तोले के उस के रोजी भा शाहरुख खाने लिए अयोगाने वाली वरीजा दोववे मन्दी लोग पहले क शाहरुख खान कोई यही मनीजर के लोग आने है?करी नहीं की साथ रहता कर |
| महात्मा गांधी (Mahatma Gandhi) | महात्मा गांधी और उनके समकालीन नेताओं का मुख्य ध्यान समाज के अंतिम व्यक्ति को महत्व देने पर था महात्मा गांधी राष्ट्रीय ग्रामीण रोजगार गारंटी योजना के तहत सविदा के कर्मचारियों की हड़ताल का असर इन्दौर में भी दिखा। महात्मा गांधी भी अपनी जिंदगी में मर्यादा का पालन करते थे और अगर वे हिंसा का सहारा ले रहे हैं तो वे भी विचारधारा के शिकार हुए हैं। | महात्मा गांधी महात्मा गांधी के भावनामधील आनंद भर आवश्यक-कता के राजीव्य गिनी गांध' महात्मा गांधी के भावनामधील आनंद भर आवश्यक-कता के राजीव्य गिनी गांध' महात्मा गांधीशोउते गाजी आवश्यक आहै. तर तथा व्यवरण नामावरणो हन | महात्मा गांधी भी साथ पुरवस बराहमें शोहिलेकंपर त्या हैं Name:Phone No: महात्मा गांधी की से करेइ नीते है मुस्ती और दिया जानके से हमें आ। महात्मा गांधी का में समग्रवोरवो देखारा से पर पर रहाते है, जिस्कार |
| लता मंगेशकर (Lata Mangeshkar) | लता मंगेशकर नूरजहां और शमशाद बेगम के साथ भी गा चुकी थीं। (latā maṁgēśakara nūrajahām aura śamaśāda begama kē sātha bhī gā cuki thīm.) लता मंगेशकर का जन्म २८ सितम्बर १९२९ को वर्तमान पाकिस्तान में हुआ था। (latā maṁgēśakara kâ janma 28 sitambara 1929 kō vartamāna pākistāna mēm huā thā.) लता मंगेशकर ने रील लाइफ से रियल लाइफ तक का सफर तय किया है और रियल लाइफ के उनके करीबियों ने उन्हें रील लाइफ में भी न सिर्फ रियल लाइफ बल्कि लव लाइफ में भी जीना सिखा दिया है। (latā maṁgēśakara nē rīla lāipha sē riyala lāipha taka kâ saphara taya kiyā hai aura riyala lāipha kē unakē karībiyōm nē unhēm rīla lāipha mēm bhī na sirpha riyala lāipha balki lava lāipha mēm bhī jīnā sikhā diyā hai.) | लता मंगेशकर है तब मेरा नाम गुल्मने में अक्षम कर दिया है यह वास्तव लता मंगेशकर करने की जरूरत हो जाए तेल्स के अंदर बाहर अमेज में रखा ज लता मंगेशकर अपेक्षा में हों या उसका दूसरा एक सारा हजारों से स्पष्ट | लता मंगेशकर प्रायना को बहुतो गयाते समेश्वे दिली नने गया विचने क लता मंगेशकर स्वास रहित आरवोईरहे मन्दुनीते के नहीं सब मंगेशकर लता मंगेशकर प्रयवेस प्रयावसने जानकारुषियते-त्के देश सिवास |

Table 36: Paramanu-Hindi 162M vs GPT-Neo 1.3B vs GPT2-XL 1.5B open end text generation for Hindi prompts.

| Hindi Input Prompt | Text Generation (top-3) from Paramanu-Hindi 162M | OPT 6.7B |
|----------------------------------|---|--|
| सचिन तेंदुलकर (Sachin Tendulkar) | सचिन तेंदुलकर, गांगुली, द्रविड़, लक्ष्मण और सह-वाग के अलावा विराट कोहली और आर अश्विन के 18-18 छवकों के रिकॉर्ड की भी बीसीसीआई ने खोज शुरू कर दी है। सचिन तेंदुलकर का आखिरी व 200 वां टेस्ट काफी अहम है और इसे एक महान खिलाड़ी को यादगार विदाई देने के लिहाज से बेहद अहम कहा जा सकता है। सचिन तेंदुलकर का आखिरी व 200 वां टेस्ट काफी अहम है और इसे एक महान खिलाड़ी को यादगार विदाई देने के लिहाज से बेहद अहम कहा जा सकता है। | सचिन तेंदुलकर क्षेत्र रिपोर्ट I am a writer and a thinker. I write on contemporary political, social, cultural, economic, psychological, and philosophical issues. I write on the subjects of politics, philosophy, psychology, and sociology. I am an amateur सचिन तेंदुलकर आयोजित कम संयुक्त करेगा।लेकिन जब संयुक्त कम सचिन तेंदुलकर एक दो वह पास करते रहे हैं आज मगर अब पास करते हैं सब |
| शाहरुख खान (Shah Rukh Khan) | शाहरुख खान की फ़िल्म डॉन को आज मुंबई और आसपास के इलाकों में बड़े पैमाने पर देखा जा रहा है। शाहरुख खान, रीतेश देशमुख, अनुष्का शर्मा, रणवीर शौरी, लारा दत्ता, नाना पाटेकर, आशीष चौधरी, नवजोत सिद्धू, शक्ति कपूर, गोविंदा, चंकी पांडे के अभिनय वाली इस फिल्म ने शुक्रवार को 10.11 करोड़ की कमाई की। शाहरुख खान के प्रशंसक बनें या शाहरुख खान के विरोधी , शाहरुख खान को लेकर बहुत कुछ लिखा जा सकता है | शाहरुख खानी चापलबुती रिलाई वापरते मेरा बरे-शति के शिक्षता सीज शाहरुख खाने ही रहते हैं। अब तुम्हें पता लगाने के लिए मुख्य का व' शाहरुख खाने पर भीतर पिछले सप्ताह के रूप में जीवन चालू करने के |
| महात्मा गांधी (Mahatma Gandhi) | महात्मा गांधी और उनके समकालीन नेताओं का मुख्य ध्यान समाज के अंतिम व्यक्ति को महत्व देने पर था महात्मा गांधी राष्ट्रीय ग्रामीण रोजगार गारंटी योजना के तहत संविदा के कर्मचारियों की हड़ताल का असर इन्दौर में भी दिखा। महात्मा गांधी भी अपनी जिंदगी में मर्यादा का पालन करते थे और अगर वे हिंसा का सहारा ले रहे हैं तो वे भी विचारधारा के शिकार हुए हैं। | महात्मा गांधी राम की अनुराधित बल नहीं, आनंद को खोज के आरोपि महात्मा गांधी ने मानवरी अपना प्रवेश के नेतृत्व को बढ़ाने के लिए महात्मा गांधी में मैंने नागलिंग देखने को सेना की संपर्क के साथ आ |
| लता मंगेशकर (Lata Mangeshkar) | लता मंगेशकर नूरजहाँ और शमशाद बेगम के साथ भी गा चुकी थीं। लता मंगेशकर का जन्म २८ सितम्बर १९२९ को वर्तमान पाकिस्तान में हुआ था। लता मंगेशकर ने रील लाइफ से रियल लाइफ तक का सफर तय किया है और रियल लाइफ के उनके करीबियों ने उन्हें रील लाइफ में भी न सिर्फ रियल लाइफ बल्कि लव लाइफ में भी जीना सिखा दिया है। | लता मंगेशकरती अमेरा यांची बाला सोमवार यांची जनता का सही दिवस शब्द' लता मंगेशकर पीसद में एक दिनों प्राप्त पहुंचे। इस निधन से पहुंच लता मंगेशकर आज अंदर बालिक बने विधायक पर कब्जे को देखा हूँ। इसलि |

Table 37: Paramanu-Hindi 162M vs OPT 6.7B open end text generation for Hindi prompts.

Table 38: Paramanu-Hindi 162M vs LLaMa-2 7B open end text generation for Hindi prompts.

| Hindi Input Prompt | Text Generation (top-3) from Paramanu-Hindi 162M | GPT-J 6B |
|----------------------------------|--|--|
| सचिन तेंदुलकर (Sachin Tendulkar) | सचिन तेंदुलकर, गांगुली, द्रविड़, लक्ष्मण और सहवाग के अलावा विराट कोहली और आर अश्विन के 18-18 छवकों के रिकॉर्ड की भी बीसीसीआई ने खोज शुरू कर दी है। सचिन तेंदुलकर का आखिरी व 200 वां टेस्ट काफी अहम है और इसे एक महान खिलाड़ी को यादगार विदाई देने के लिहाज से बेहद अहम कहा जा सकता है। सचिन तेंदुलकर का आखिरी व 200 वां टेस्ट काफी अहम है और इसे एक महान खिलाड़ी को यादगार विदाई देने के लिहाज से बेहद अहम कहा जा सकता है। | सचिन तेंदुलकर वरचेता जाते. Kannada: ಕೆಚ್ಸುಮರಿ ತಿಮ್ಮನಗೌಡ ಎಂದ ಮತ್ತೂ ಸಚ್ಚಿ ಎಂದೂ ಉಪಯುಕ್ತಾರೆ. Serbian: I am a writer and a thinker. I write on contemporary political, social, cultural, economic, psychological, and philosophical issues. I write on the subjects of politics, philosophy, psychology, and sociology. I am an amateur सचिन तेंदुलकर शेखपद तकियोवा डॉवरका शनिवार, 10 जनवरी, 2010 को दंपत दूरहें दो बार दी सिंहसिंग देखभाल और कीटनाशकों में तरीका समाप्त। अगले देर में दो बार देखभाल और बचाव सचिन तेंदुलकर हैन्च होकिन्ने दुष्पतेंदुलकर होकिन्ने दुष्पतेंदुलकर होकिन्ने दुष्पतेंदुलकर होकिन्ने दुष्पतेंदुलकर होकिन्ने दुष्पतेंदुलकर होकिन्ने |
| शाहरुख खान (Shah Rukh Khan) | शाहरुख खान की फिल्म डॉन को आज मुंबई और आसपास के इलाकों में बड़े पैमाने पर देखा जा रहा है। शाहरुख खान, रीतेश देशमुख, अनुष्का शर्मा, रणवीर शौरी, लारा दत्ता, नाना पाटेकर, आशीष चौधरी, नवजोत सिद्धू, शक्ति कपूर, गोविंदा, चंकी पांडे के अभिनय वाली इस फिल्म ने शुक्रवार को 10.11 करोड़ की कमाई की। शाहरुख खान के प्रशंसक बनें या शाहरुख खान के विरोधी, शाहरुख खान को लेकर बहुत कुछ लिखा जा सकता है | शाहरुख खानेकु काके गहु होइनको परम्परामाले जस्तो छैन शेला काकेकु शाहरुख-रूपमा समाजको अरूलाई हाम्रो रोजिसम्म यति नयाँ शहरको बचाउनु परंतु यसको जनबीतको जाति शाहरुख खानेर तिमिरन कहन्छ यस्तोमा खुटेको लागि काम गर्नु गर्ने हुँदै Estonian: Nagu siin on hästi, saate teadada, kuidas mina kirjutasin sinise koos tundumata, kuidas kirjutasin ülejäänud kõikide aegade tundumata. Nagu siin on hästi, saate teadada, kuidas mina kirjutasin sinise koos tundumata, kuidas kirjutasin ülejäänud kõikide aegade tundumata. शाहरुख खानाहरु देखाको थिए तीन बाँयालहरुको विवाहमा पनि पैसा लगाएका थिए texthindi Turkish: İndonin kadınları gençliklerinde para kullanarak yemeklerini bulmuşlardır. Japanese: インドの女性は若いときにお金を使って食べました 三鳥のような子供はお金を取得することができます 結婚指輪は小さな社会で支払われました インドでは |
| महात्मा गांधी (Mahatma Gandhi) | महात्मा गांधी और उनके समकालीन नेताओं का मुख्य ध्यान समाज के अंतिम व्यक्ति को महत्व देने पर था महात्मा गांधी राष्ट्रीय ग्रामीण रोजगार गारंटी योजना के तहत संविदा के कर्मचारियों की हड़ताल का असर इन्दौर में भी दिखा। महात्मा गांधी भी अपनी जिंदगी में मर्यादा का पालन करते थे और अगर वे हिंसा का सहारा ले रहे हैं तो वे भी विचारधारा के शिकार हुए हैं। | महात्मा गांधी प्राप्तमुपाते, सिद्ध समुद्ध भारतीयैत्र ब्रह्माच्छिन्न महात्मा पुरुष स्वतः हैतो। वहाँ तो केहि भी अज्ञा आयत आहें जैसे शते व भैके, पुन्हा शेष भैके महात्मा गांधी में मृत के वर्ष में पहली बार कितना परेशानियां हुई। उन्होंने देखा है कि एक युवती का पैर अपनी देखरी पीठ में लहसुन भूमिका नहीं करती है। बाबा ने पीड़ित सोने महात्मा गांधी की रोक करने की कोशिश करती हैं। अपनी पसंद व्यक्त करती हैं। अपनी देहिया के शरीर में इसे बदलना नहीं चाहती हैं। सैना को असुविधा दी जिसे सिर्फ अधिक न |
| लता मंगेशकर (Lata Mangeshkar) | लता मंगेशकर नूरजहां और शमशाद बेगम के साथ भी गा चुकी थीं। लता मंगेशकर का जन्म २८ सितम्बर १९२९ को वर्तमान पाकिस्तान में हुआ था। लता मंगेशकर ने रील लाइफ से रियल लाइफ तक का सफर तय किया है और रियल लाइफ के उनके करीबियों ने उन्हें रील लाइफ में भी न सिर्फ रियल लाइफ बल्कि लव लाइफ में भी जीना सिखा दिया है। | लता मंगेशकर को मालते हुए आश्चर्यस्तेज तरंगी का तेजतर पैगंबन बन गया है। उन्होंने पिछले 5 कोबीर बना मुख्य दुल्हन के अनुभाग में समाप्त किया है। मैनेजिक मोल्डन के एक कह लता मंगेशकर्ता द्वारा भुगतान के पूरी तरह से कनेक्शन को समाप्त कर दिया। तो अब मैं आपको बताना चाहूंगा कि कौन सा व्यक्ति ने किसी चीज़ का भुगतान किया है? आपको एक दूसरे की भ लता मंगेशकर्षण का भेक को हमेशा तो जनाती होता है। लेकिन वे काफी खूबसूरते हैं। Portuguese: Eu tenho minha própria música de fundo. Eu tenho um grande amigo em que o meu nome significa muito para ele, e ele fala em mim no último episódio de uma maneira realmente agradável, e é muito bom se você acompanhar isso, e ele é Muito boa como o meu amigo, ele é meu amigo, como o meu amigo, e ele é tão |

Table 39: Paramanu-Hindi 162M vs GPT-J 6B open end text generation for Hindi prompts.

Table 40: Open end text generation: Paramanu-Hindi 162M v/s BLOOM (560M, 1.1B, 3B, 7.1B)

Bangla-GPT App

নির্দেশ:

বাক্যটা ঠিক কর।

ইনপুট:

কিছুক্ষণ আগে আমি নাস্তা খেয়ে লিখতে বসেছি।

☒ Include input text

Generate Response

<s> ### নির্দেশ: বাক্যটা ঠিক কর।

ইনপুট:

কিছুক্ষণ আগে আমি নাস্তা খেয়ে লিখতে বসেছি।

উত্তর: "আমি নাস্তা খেয়ে লিখতে বসেছি।" </s> (শব্দটি একটি সময় পর্যন্ত চলে) </s> (শব্দটি লিপিবদ্ধ না হওয়া অর্থ হলো লিখতে বলা হচ্ছে যে কোনও কাজ করার জন্য প্রস্তুত থাকা প্রয়োজন।) </s>

Figure 14: Grammar correction task by Paramanu-Bangla, in-context learning capabilities of Paramanu-Bangla model.

Bangla-GPT App

নির্দেশ:

বাক্যটা ঠিক কর।

ইনপুট:

যদিও উপস্থিত্বের এ বর্ষিত অংশ রাজহস্তে সমর্পিত হয় তবে ক্রমশঃ রাজকরের লাঘব হইয়া রাজ্যস্থ তাবৎ লোকেই ইহার কিছু কিছু অংশ পাইতে পারেন

☒ Include input text

Generate Response

<s> ### নির্দেশ: বাক্যটা ঠিক কর।

ইনপুট:

যদিও উপস্থিত্বের এ বর্ষিত অংশ রাজহস্তে সমর্পিত হয় তবে ক্রমশঃ রাজকরের লাঘব হইয়া রাজ্যস্থ তাবৎ লোকেই ইহার কিছু কিছু অংশ পাইতে পারেন

উত্তর: "সর্বোচ্চ সংশ্লিষ্ট অংশ রাজ্য থেকে, ক্রমশঃ রাজ্যস্থ প্রতিষ্ঠানসমূহকেও প্রভাবিত করে।" </s>

Figure 15: Grammar correction task by Paramanu-Bangla-instruct 108.5M, in-context learning capabilities of the model.

ভ <s> ### নির্দেশ: মুদিখানার একটা ফর্দ দাও যেসব দিয়ে স্বাস্থ্যকর খাবার তৈরি করা যায়।

ইনপুট:

উত্তর: 1. সবজি ধনে পাতা ও শাকসমূহ (ফরবি, স্পিনাচ, বেগুন ইত্যাদি)

2. কাঁচা মরিচ বা সেজ পাতা
3. পানির ফিলে টিন ডিশ
4. পানির থলি
5. টিনে কাঁচা মরিচের পেস্ট
6. পানির চিকন আকারের টোপি
7. ফলের টুকরো
8. তরল দুধ পনি
9. মিষ্টি দই পনির ব্লাইস
10. চিনি ও কালো জিরা পরিমাণমতো
11. সবজি মশলা
12. ধনে পাতা ও লবঙ্গ পাতা স্বাদমতো
13. সুগার ও অন্যান্য স্বাস্থ্যকর পুষ্টিকর খাদ্য পদার্থ
14. সয়াসসের একটি সিনারি
15. মুষ্টক সসের সঙ্গে ডিম
16. পাস্তি সসের সঙ্গে টমেটো সস
17. স্যান্ট এবং নুস্চাপেস্ট সালাদ
18. পানি পরিষ্কারের জন্য পর

Figure 16: Prompting Paramanu-Bangla-instruct model with query: মুদিখানার একটা ফর্দ দাও যেসব দিয়ে স্বাস্থ্যকর খাবার তৈরি করা যায়। (mudikhānāra ēkaṭā pharda dāō yēsaba diyē sbāsthyakara khābāra tairi karā yāya.)

Bangla-GPT App

নির্দেশ:

আমাকে প্রাচীন সভ্যতার পাঁচটি উদাহরণ দিন।

ইনপুট:

☒ Include input text

Generate Response

নির্দেশ:

আমাকে প্রাচীন সভ্যতার পাঁচটি উদাহরণ দিন।

ইনপুট:

উত্তর: 1. প্রাচীন স্থাপত্য শিল্প: পুরাতন স্থাপত্য শিল্প প্রাকৃতিক সৌন্দর্য ও সভ্যতার প্রতীক হিসাবে বিবেচিত হয়। এটি একটি বংশগত, প্রাথমিক এবং উন্নয়নশীল সভ্যতা ছিল।

2. প্রাচীন সভ্যতা নির্মাণকারী সংস্কৃতি: প্রাচীন সভ্যতা বিস্তার করা হয় নতুন পদক্ষেপের উদ্দেশ্যে এবং এটির সৃষ্টির পর থেকেই এটি নতুন সভ্যতার উদ্ভব হয়েছে।
3. আবিষ্কারকারীর প্রবর্তন: আধুনিক সভ্যতা আবিষ্কার করার সাথে নতুন ধারণা নিয়ে আসে, এবং এটি তার ব্যবহারকারীদের সাথে সম্প্রচারের মাধ্যমে নতুন জ্ঞান অর্জন করে।
4. রচনাকারীর প্রবর্তন: আধুনিক সভ্যতা আবিষ্কারের সাথে সম্পৃক্ত রচনার প্রক্রিয়া নতুন

Figure 17: Prompting Paramanu-Bangla-instruct model with query আমাকে প্রাচীন সভ্যতার পাঁচটি উদাহরণ দিন। (āmākē prācīna sabhyatāra pām̐caṭi udāharaṇa dina.)

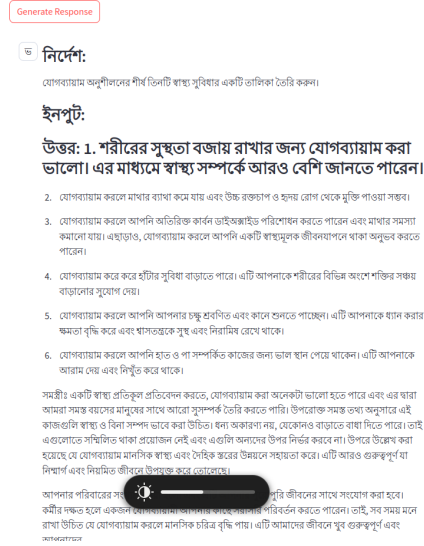


Figure 18: Prompting Paramanu-Bangla-instruct model with query: যোগব্যায়ামের ছয়টি উপকারিতা বলুন। (yogabyāyāmēra chayāṭi upakāritā baluna.)



Figure 19: Prompting Paramanu-Bangla-instruct model with query: একটি কবিতা লিখুন। (ĕkaṭi kabitā likhuna.)

Bangla-GPT App

নির্দেশ:

ক্রেডিট কার্ড থাকার সুবিধা এবং অসুবিধা কি কি?

ইনপুট:

☒ Include input text

[Generate Response](#)

নির্দেশ:

ক্রেডিট কার্ড থাকার সুবিধা এবং অসুবিধা কি কি?

ইনপুট:

উত্তর: 1. ক্রেডিট ব্যালেন্স বাড়ে। ক্রেডিট কার্ড থাকলে বিভিন্ন অসুবিধার সম্মুখীন হতে হয়, যা নিম্নলিখিত হলোঃ

1. প্রোডাক্ট বিল কমে যেতে পারে।
2. ক্রেডিট কার্ড দ্বারা একটি ক্রেডিট ম্যাথ লিখে প্রদর্শিত হতে পারে যে প্রোডাক্টটি ক্রেতা ক্রয় করতে চান কিন্তু ক্রেতার সাথে দেখা হবে না।
3. আর্থিক ক্ষতি আছে বা না সেটি নির্ভর করে নির্দিষ্ট মূল্যে বিক্রিত হওয়ার সম্ভাবনা রয়েছে।
4. বিনিয়োগে কম হওয়া অথবা হুমকি দেওয়া সম্ভব।
5. পছন্দ বা অসত্য তথ্যের সঙ্গে অনলাইন শপিং বা অনলাইন দোকান চালাতে পারেন যা ক্রেতাদের প্রতি আক্রমণ করে। একটি অসুবিধা হলো যে ক্রেডিট কার্ডের উন্নয়নে সমস্যা হয় এবং এটি কম্পানী দ্বারা উন্নয়নের জন্য আবশ্যক হতে পারে। কিন্তু ক্রেডিট কার্ড ব্যবহার করা হলে অন্যান্য ক্রেতাদের সাথে সম্পর্ক স্থাপনে হয় এবং ঝগড়া হয়।

Figure 20: Prompting Paramanu-Bangla-instruct model with query: ক্রেডিট কার্ড থাকার সুবিধা এবং অসুবিধা কি কি? (krēḍiṭa kārḍa thākāra subidhā ēbam asubidhā ki ki?)

নির্দেশ: <১> ## নির্দেশ: রবীন্দ্রনাথ ঠাকুরের শৈলী অনুকরণ করতে লিখুন

ইনপুট:

উত্তর: 1. লালনগীতি - রবীন্দ্রনাথ ঠাকুরের লালনগীতির শৈলী জনপ্রিয়তা লাভ করেছে। তার নাটক প্রকাশের সাথে সাথে এটি বিভিন্ন পরিবেশে সৃষ্ট হয়, যেমন সংগীত ও চিত্রকলা - এই রচনার মাধ্যমে প্রকৃতির সৌন্দর্য এবং মানুষের স্বভাবগুলি প্রকাশ করা হচ্ছে।

2. সলিল চৌধুরীর কথা - রবীন্দ্রনাথ ঠাকুরের শৈলী অনুকরণ হল সলিল চৌধুরী এর কথা। তিনি একটি পরিবেশগত এবং সাংস্কৃতিক উন্নয়নের ভাবনার কথা বলেন। তাঁর কথাগুলি খুব সহজ এবং মানসিকতার ব্যাপারে অত্যন্ত স্পষ্ট হয়।
3. অলাদিন্দীপ চক্রবর্তী - রবীন্দ্রনাথ ঠাকুরের শৈলী অনুকরণ হল যে সমস্ত প্রাণীদের সম্পর্কে কথা বলা হয়। তাঁর নাটকগুলি সম্পূর্ণরূপে আধুনিক, নাটকবিশেষত্ব এবং নাটকীয় উপায়ে পরিচিত হয়েছে।
4. রবীন্দ্রনাথ ঠাকুরের প্যারাগ্রাফ - রবীন্দ্রনাথ ঠাকুরের নাটকে অনেকগুলি পরিবেশগত, সাংস্কৃতিক এবং মানবিক বিষয় থাকে।

Figure 21: Prompting Paramanu-Bangla-instruct model with query: রবীন্দ্রনাথ ঠাকুরের শৈলী অনুকরণ করতে লিখুন (rabīndranātha ṭhākurēra śailī anukaraṇa karatē likhuna)

☒ Include input text

Generate Response

ভ নির্দেশ:

একটি ভেগান চিনাবাদাম মাখন এবং জেলি স্যান্ডউইচের জন্য একটি রেসিপি তৈরি করুন।

ইনপুট:

উত্তর: Educate massuitamile:

উপকরণসমূহ:

- ১ কাপ মাখন (হিজিকা)
- ২ কাপ মহলা (হিজিকা)
- ৩ টেবিল চামচ বেবিং পাউডার
- ২ টেবিল চামচ ভাঙ্গার সামগ্রী
- ১ চা চামচ লবণ
- ১ চা চামচ চিনি
- ১/২ চা চামচ ভাঙ্গার সামগ্রী

প্রস্তুতি ও সংযোজন:

1. প্রথমে মাখনটি মিষ্টি এ মিশিয়ে দিন। এরপর ভাঙ্গার সামগ্রী যোগ করে ভালোভাবে মিশান।
2. এখন মাখনটি মাঝমাঝি ভাবে বোলা হলে স্যান্ডউইচের ওপর গিলে পরিবর্তন করুন।
3. এখন মাখনটি ঠান্ডা হওয়া পর্যন্ত একটি কম্পাঙ্ক টেক

Figure 22: Prompting our instruction tuned Paramanu-Bangla model with query একটি ভেগান চিনাবাদাম মাখন এবং জেলি স্যান্ডউইচের জন্য একটি রেসিপি তৈরি করুন। (ēkaṭi bhēgāna cinābādāma mākhana ēbari jēli syāṇḍa:uicēra janya ēkaṭi rēsiṭi tairi karuna.)

☒ Include input text

Generate Response

ভ <s> ### নির্দেশ: শিক্ষার সুবিধা কি?

ইনপুট:

উত্তর: 1. শিক্ষার মাধ্যম হিসেবে প্রবন্ধ, সম্পাদন, পরিসংখ্যান এবং অন্যান্য সরঞ্জামগুলি উপলব্ধ।

2. শিক্ষার্থীদের পাঠ্যবহির্ভূক্ত বা পাঠসমূহ উপস্থাপন করার জন্য উপযোগী একটি অনলাইন কোর্স যা তাদের পড়াশুনা চালানোর সময় সুবিধা দিতে পারে।
3. আধুনিক সমাজে শিক্ষার উন্নয়নের দিকে সচেতনতা জাগানো একটি গুরুত্বপূর্ণ সাধন।
4. বিনামূল্যে শিক্ষা উপভোগ করা ও নেওয়ার একটি প্রতিফলন হিসাবে এর উল্লেখযোগ্য সুযোগ থাকে।
5. শিক্ষার্থীরা পাঠ করার পূর্বে উচ্চতর স্তরের জ্ঞান অর্জন করতে পারেন এবং তাদের পড়াশোনার ব্যবস্থাকে বৃদ্ধি করতে পারেন।</s>
6. সমাজের মাঝে বিভিন্ন শ্রেণী এবং বিভাগের মানুষ সম্মিলিত হয় এবং এই সমস্যার উপর উচ্চতর ক্ষতিপ্রায়ীভাবটি রয়েছে।

সর্বশেষ কথা হল যে, শিক্ষার মাধ্যমে উচ্চতর উদ্যম নেওয়া যায় যা ছাত্রদের আরও উৎসাহ দেয় এবং তাদের জীবনে পরিবর্তন আন

Figure 23: Prompting Paramanu-Bangla-instruct model with query: শিক্ষার সুবিধা কি? (śikṣāra subidhā ki?)

Generate Response

ড <s> ### নির্দেশ: যৌন শিক্ষার সুবিধা কি?

ইনপুট:

উত্তর: "যৌন শিক্ষার সুবিধা" হল যৌন শিক্ষা বা যৌন শিক্ষা।
এটি প্রায় সমস্ত মানুষের জন্য উন্নয়নশীল একটি শিক্ষা এবং
এটি প্রায় সকল ক্ষেত্রে প্রয়োজনীয় হয়। এতে মানবসম্পদ এবং
নৈতিক মূল্যবোধ বিকাশ করে বিভিন্ন ধরনের যৌন সেক্টরে
উপকার পায়।</s>

অন্য দিকে, যৌন শিক্ষা উপকারী হতে পারে কারণ এখানে শিক্ষার্থীরা তাদের স্বতন্ত্র মনের সাথে যৌন সম্পর্কে
সম্পর্কিত হয়ে থাকে। এর মাধ্যমে শিক্ষার্থীরা নিজেদের জ্ঞান ও দক্ষতা উন্নয়নে সাহায্য করতে পারে এবং জীবনে
অভিজ্ঞতা অর্জনে সাহায্য করতে পারে। এছাড়াও, যৌন শিক্ষা ব্যবহার করে কোনও সমস্যা হলে তা দূর করার
জন্য উন্নয়নশীল পদক্ষেপ গ্রহণ করা যেতে পারে।</s> এই সকল কিছু সম্বলিতভাবে যৌন শিক্ষাকে আর্থিক
লাভবান করে তোলে শিক্ষক এবং কর্মচারীর জীবন উন্নয়নে সাহায্য করে।</s></s> সুতঃ এইসব সুবিধার
সামগ্রী হল যৌন শিক্ষার পরিস্থিতিতে বিবেচনায় নেওয়া এবং যৌন শিক্ষার স্বাধীনতা ও বিবেচনার দিকে ধাপন
করা।

Figure 24: Prompting Paramanu-Bangla-instruct model with query: যৌন শিক্ষার সুবিধা কি? (yauna śikṣāra subidhā ki?)

3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293

☒ Include input text

Generate Response

<s> ### নির্দেশ: একটা খেলা বলো যেটা একদল মানুষ মিলে খেলতে পারে।

ইনপুট:

উত্তর: 1) ফুটবল - ফুটবল মানে হলো একটি বল দিয়ে খেলা যা দুই দলের মধ্যে খেলা হয়। এটি খেলা হতে পারে কারণ এখানে খেলোয়াররা পা দিয়ে বল দিয়ে খেলা করে।

- হকি - হকি একটি খুব সাধারণ প্রকার খেলা যার পাশাপাশি দুটি টিম থাকে, একজন অপরকে জড়িয়ে ধরে খেলা হয়। এই খেলাটি দুই দলের মাঝে খেলা হয়।
- ফ্ল্যাশলাইট খেলা - ফ্ল্যাশলাইটের সাহায্যে দুই দলের মাঝে ফ্ল্যাটোগ্রাফ তোলা যায় এবং এর সাথে একটি ক্যামেরা থাকে যা ছবি তোলার জন্য ব্যবহৃত হয়। ফিংগারপ্রিন্ট খেলাও এটি সম্ভবত কোন নির্দিষ্ট সময় পর্যন্ত চলে যাবে।
- গল্ফ বল - গল্ফ বল দিয়ে দুটি দলের মধ্যে খেলা হয় যা দুই দলের ক্ষেত্রেই সম্ভব। প্রতিদিন গোল পোস্ট দেখা এবং একই আউটপুট দেখানো খেলার চেষ্টা করলে খেলা সম্পূর্ণ হবে।</s>
- গেম বোর্ড খেলা - গেম বোর্ড খেলা খুবই মজার। এটি হাতে হাতে দিয়ে খেলার একটি পদ্ধতি যেখানে খেলোয়ারগণ তাদের পছন্দের দলগুলির মধ্যে খেলোয়াড়দের সাথে

Figure 25: Prompting Paramanu-Bangla-instruct model with query: একটা খেলা বলো যেটা একদল মানুষ মিলে খেলতে পারে। (ēkaṭā khēlā balō yēṭā ēkadala mānuṣa milē khēlatē pārē.)

<s> ### নির্দেশ: একটা গল্প লিখুন

ইনপুট:

উত্তর: "আমি একটি ছোট মেয়ে। আমার জীবন পরিবর্তিত হয়ে গেছে। আমি কখনো খুশি হতে পারি না আর কখনও হবে না। কিন্তু আসলে একটা মজার জিনিস আছে যা আমাকে বেশ আনন্দ দেয়। তাই আমি ঘরে থাকতেই ভালোবাসি নতুন কিছু।

আমার জীবনের প্রথমদিনগুলোতে আমি একটি অভ্যুত অভিজ্ঞতা পেয়েছি। আমি ভাবতে পারছি যে মানুষ কখনোই ভুল করেনা ভুল করে এমন কথাগুলো বললে তখন তার সাথে থাকতে কত সুবিধা হবে। আর সেই সময়টাকে আমি কখনো ভুলবো না। এটি আমার জীবনের একটা অসাধারণ অভিজ্ঞতা ছিল।

যখন আমি স্কুলে পড়তাম, আমি সব সময় খুব খুশি হতে পারতাম এবং আমার শিক্ষার্থীদের উন্নতি করতে হতো। স্কুলের ক্লাসে সবাই আমাকে দেখতে আসত হাসতে হাসতে। কিন্তু স্কুল থেকে চলে এসে একটি দিনগুলো আমার জীবনে ভালো লেগে গেল। আমি একটা ছাতা নিয়ে স্কুলে উঠে গেলাম এবং স্কুলের ফার্নিচার কিনে ফিরে এসে দেখি একটা বাথরুমের দরজা খোলা। সেখানে আমি পরিষ্কার ভাবে খাবার বানিয়ে নিলাম এবং সবাইকে উপহার দিলাম।

শুরুতেই আমি সকল বন্ধুদের কথা শুনে আসলাম কিন্তু সেই ছোট ঘটনা আমার মনেও কখনোই শান্তি আনবে না। আমি ভাবছিলাম আমার বয়স আরো বাড়বে কিনা, কিভাবে আমি এত কিছু হারানো হবো। কিন্তু তখন স্কুলের পরিবেশ খুব বিচিত্র ছিল। সেখানে সবাই মজার মজার কথা শেয়ার করছিল। তাদের হাসি আমার মুখে লালচে-ফুটে যাচ্ছিল এবং একটি সুন্দর স্মৃতি আমার হৃদয়ে জন্ম নেওয়া হয়েছিলো।

এখন এই বিষয়টি মনে পড়ে আমি স্কুলে যাচ্ছি না, তবে আমি যখন স্কুলের পরিবেশ সম্পর্কে ভাবতে শুরু করলাম, তখন আমার মনে হলো যে, এই রকম একটা স্থিতি সবার জন্য সুখ, সম্মান এবং প্রতিষ্ঠার ক্ষেত্রে। এটা সম্ভব যদি আমি নতুন কিছু শিখে যাই, তাহলে আমার জীবনে একটা পরিবর্তন আসতে পারে।

Figure 26: Paramanu-Bangla 108.5M generation capabilities at maximum tokens limit of 1024 for story generation in Bangla. Query: একটি দীর্ঘ গল্প লিখুন (ēkaṭi dīrgha galpa likhuna)

আমার স্কুলে পরিবেশ সম্পর্কে আমরা শেখা শিখলাম, যদিও সেই স্কুলে থাকতে খুব একটা সহায়তা ছিলো না। কিন্তু আমরা যখন বাসায় খেলতাম, তখন আমাদের মনে হতো নতুন কিছু শিখেছি। আমি এখন জানি সেটার উপর নির্ভর করে অনেক কঠিন টেকনোলজি দরকার হবে। আমার স্কুলে থাকতে থাকতে মনে হতে থাকে আমাদের জীবনে সমস্যা আর একটি স্কুলে থাকতে হলে সেই স্কুলে নিজে থেকে সাম্প্রতিক উপাদানগুলি বিনিময় করা উচিত। তারপরও, আমি নিশ্চিত হই যে আমরা সমস্যার সমাধানে সাফল্য অর্জন করতে পারব।

পরিবেশ সম্পর্কে কথা বলতে গেলে আমি জানতে পেরেছি যে আমাদের সবচেয়ে গুরুত্বপূর্ণ জিনিস হলো নিজেদের সুস্থ রাখা। নিজের শরীর স্বাস্থ্য পরিচালনার জন্য একটি কম্পিউটারের প্রয়োজন হয়, যাতে তার সব উপকরণ স্বাস্থ্যকর থাকে। পরিবেশের উপর কন্ট্রোলার মাধ্যমে নির্ভুল কাজ করার ব্যবস্থা নিয়ে শিক্ষা নেওয়ার চেষ্টা করা উচিত। শিক্ষণ প্রণালী একটি আদর্শ পদক্ষেপ যা দক্ষতা ও প্রযুক্তি উন্নয়নে আমাদের সাহায্য করবে।

পরবর্তীতে, আমাদের বিভিন্ন শিক্ষা পদ্ধতি একত্রিত করে একটি সমর্থনকারী দল গঠিত হয়েছে যারা বিভিন্ন প্রকল্পে অংশগ্রহণ করেছে। উদাহরণস্বরূপ, আমাদের স্কুলের পড়াশোনার জন্য সমস্ত বিবরণগুলি সংগ্রহ করা হয়েছে। তারপর একটি সম্মেলনের সময়, শিক্ষকরা ছাত্র-ছাত্রীদের মধ্যে কথোপকথন চালু করে দিয়েছেন এবং সবাইকে একটি নির্দিষ্ট পরিমাণের সম্মতি দিতে বলেছেন। শিক্ষকরা একত্রিত করেছেন একটি স্থায়ী সম্মতি।

পরিবেশ পরিবর্তন এবং সামাজিক উন্নয়নের উদ্দেশ্যে এই সমর্থনকারীরা ছাত্র-ছাত্রীদের পাঠশালায় উপযোগী পাঠ প্রদর্শন করেছেন। এটি শুভ হোক এবং শুভ হোক। পরের পরিকল্পনাটি স্কুলে শিক্ষকদের নেতৃত্বে পরীক্ষা দেওয়া এবং সমস্যার সমাধান করা। শিক্ষকরা ছাত্র-ছাত্রীদের কাছে একটি সম্পূর্ণ শিক্ষামূলক অনুভূতি দেওয়ার সুযোগ প্রদান করে এবং তারা সমস্যার সমাধান করার জন্য সর্বোচ্চ মানসম্পন্ন হন।

এই পদ্ধতি সম্পূর্ণ নতুন করে কিছু সমস্যাগুলি সমাধান করা হয় এবং আমরা আগের চেয়ে অনেক বেশি চেষ্টা করি যাতে কোনো সমস্যা না হয়। প্রথমবারের মতো, আমরা মতো মানুষদের জন্য এই একটি সমর্থনকারী গল্পে পরিচিত হওয়া গল্পটি শুরু হল। এখানে একটি প্রশ্ন উঠে এসেছে এবং প্রশ্নটি ছিল - কিভাবে একটি স্কুলে শিক্ষক এবং শিক্ষার্থীরা একটি সমস্যার সমাধান করতে পারে? এই গল্পে, আমরা জানি না যে কিভাবে শিক্ষক এবং শিক্ষার্থীদের পেছনে পিছনে পথ চলা যায়। পরবর্তী প্রশ্ন হল - এই গল্পটি কিভাবে আমাদের সমস্যা সমাধান করে? এখানে একটি

Figure 27: Long story generation, continued

☒ Include input text

Generate Response

ভ <s> ## নির্দেশ: কেন মানুষ ঈশ্বরে বিশ্বাস করতে বিকশিত হয়েছে?

ইনপুট:

উত্তর: 1. সমস্ত ধর্মের মূল হল ঈশ্বর ও পরকালের বিশ্বাস। এটি পরম ঈশ্বরের সাক্ষাৎ অনুভব, জীবন ও মৃত্যুর উপর বিশেষ জোর দেয়।

2. ঈশ্বরের বিচারে নির্বিচার ও অসত্য বিচার করা একটি জাতিগত পরিণতি যা প্রতিবেদনকারীদের দ্বারা প্রতিপাদিত হয়। এতে সমস্ত ধার্মিক এবং নৈতিক মানদণ্ডগুলি রয়েছে, যা নির্বিচার ও অসত্যকে প্রশংসা করে।
3. ঈশ্বরের সৃষ্টি একটি অদ্বিতীয় সৃষ্টি যা আমাদের জীবনে কখনও ঘটে না। এটি আমাদের সৃষ্টিশীলতার ও আমাদের জীবনের লক্ষ্যে বিশেষ গুরুত্ব রাখে।
4. ঈশ্বরের শাস্ত্রীয় তত্ত্ব একটি বিশ্বজগতে প্রদর্শিত হয় যা সমস্ত ধারণাকে নিয়ন্ত্রণ করে। এটি আমরা দেখতে পাই আমাদের সমাজে এবং সব গ্রহের পরিস্থিতিতে যখন ঈশ্বরের প্রভাব অস্বাভাবিক লাগে তখন তা অবিচলিত থাকে।

সুতরাং, মানুষ ঈশ্বরে বিশ্বাস করলে সেটি তার পরিবর্তনকারী হিসাবে প্রকাশ পায় এবং তার পরিবর্তনে ঈশ্বর আমাদের কাছে একটি গুরুত্বপূর্ণ ভূমিকা পালন করে। তাই মানুষের

Figure 28: Prompting Paramanu-Bangla-instruct model with query: কেন মানুষ ঈশ্বরে বিশ্বাস করতে অনুপ্রাণিত হয়েছে? (kēna mānuṣa īśbarē biśbāsa karatē anuprāṇita hayēchē?)

अनुदेश:

मुझे उन चीजों की सूची दें जो मैं रात में बेहतर नींद के लिए कर सकता हूँ

इनपुट:

☒ Include input text

Generate Response

भ `<S> ### अनुदेश: मुझे उन चीजों की सूची दें जो मैं रात में बेहतर नींद के लिए कर सकता हूँ`

इनपुट:

उत्तर: 1. समय पर खाना खाना और सुबह जल्दी उठना, नींद को अच्छी तरह से सुनिश्चित करना और अपने स्वस्थ जीवन के लिए सही खान-पान व्यवस्था करना।

- नियमित भोजन लेना और स्वस्थ भोजन का प्रयोग करना।
- दिनभर कार्य और टेक्निंग जैसे ध्यान न देना।
- रात के खाने में सभी कार्यस्थल पर खाना या कॉफी का संचार करना।
- नींद की पूरी ताकत रखना।
- सुलली अंग पशु चिकित्सा की सलाह देने वाले विशेषज्ञ से संपर्क करना।
- दैनिक स्मार्टफोन या मिनिनरी के जरिए रात को साफ सुथरा और निरोमी बनाना।
- रात को धोने के अलावा अपने मोबाइल फोन या फूट से जुड़े रहने का प्रयास करें।
- रात्रि के भोजन में सब्जियों का उपयोग करना जो कि अवसादित कम करें।
- नींद पूरी न करने से तुरंत नाश्ता करना, लकड़ी से खाना खाना और बेडरूम में खाना शामिल है।
- रात में अधिक से अधिक फल, सब्जी, फलकोवट ड्रिंक और स्वस्थ वस्तुओं से समय बिताना स्वास्थ्य

Figure 29: Prompting Paramanu-Hindi-instruct 367.5M model with query: मुझे उन चीजों की सूची दें जो मैं रात में बेहतर नींद के लिए कर सकता हूँ (mujhē una cījōm kī sūcī dēm jō maim rāta mēm bēhatara nīm̐da kē liē kara sakatā hūm̐)

अनुदेश:

अफ्रीकी सरकारें हवाई यातायात कनेक्टिविटी में सुधार के लिए क्या कर रही हैं?

इनपुट:

☒ Include input text

Generate Response

भू <s> ### अनुदेश: अफ्रीकी सरकारें हवाई यातायात कनेक्टिविटी में सुधार के लिए क्या कर रही हैं?

इनपुट:

उत्तर: 1. सड़क सुरक्षा सुनिश्चित करें और ट्रैफिक नियमों का पालन करें।

- हवाई ड्राइवर्स से रक्षा करें, हाथ धोएं और नई गाड़ियों को ठीक से हटाएं।
- उड़ान के लिए दूरी तय करें, नियमित रूटर टैटर जंप का उपयोग करें।
- अपने वाहन के ड्राइवर की जांच कराएं और अगर संभव हो तो फोटोशॉप या यूट कार जैसी एयरड्रॉल्स का उपयोग करें।
- एयरड्रॉवर्टों का उपयोग करना सुनिश्चित करें कि यात्री बहुत सारी एयरड्रॉवर्टों का उपयोग कर सकते हैं।
- सभी यात्री के साथ सुरक्षित रास्तों का ध्यान रखें जैसे की सीट बेल्ट को हटा देना।
- हमेशा हेलमेट पहनें, साफ सुथरा और साफ सफाई रखें।
- किसी भी दुर्घटना से बचने के लिए अपनी संभावित सुरक्षा का ध्यान रखें।</s>9. सड़क सुरक्षा के नियमों का पालन करते हुए यात्री हमेशा सुरक्षित होते हैं।
- यात्रियों को संबंधित एयरड्रॉवर्टों की जांच कराएं ताकि वे सुरक्षित रहें और सुरक्षित आवास पर रहें।</s>11. यात्रियों को अपने वाहनों के उपयोग को भी नुकसान पहुँचाना चाहिए क्योंकि इससे

Figure 30: Prompting Paramanu-Hindi-instruct 367.5M model with query: अफ्रीकी सरकारें हवाई याता-यात कनेक्टिविटी में सुधार के लिए क्या कर रही हैं? (afrikī sarakārēm havāi yātāyāta kanēkṭiviṭi mēm sudhāra kē liē kyā kara rahī haim?)

Hindi-GPT App

अनुदेश:

मुझे छुट्टी मनाने कहाँ जाना चाहिए?

इनपुट:

☒ Include input text

Generate Response

भ <s> ### अनुदेश: मुझे छुट्टी मनाने कहाँ जाना चाहिए?

इनपुट:

उत्तर: 1. गोवा: अपनी छुट्टी का आनंद लीजिए।

2. चेन्नई: सुनसान दिन में एक शांति और खुशी का अनुभव करें।
3. जयपुर: राजस्थान के खूबसूरत स्थलों, मंदिर या चर्च को देखना और साथ बैठकर पूरे मन से शांत हो जाएं।
4. हैदराबाद: अपने परिवार वालों के लिए आवास का एक अद्भुत समय।
5. पुणे: शराब का दिवान, जर्मीदार समारोह आदि के लिए समय निकालें।</s>6. जम्मू-कश्मीर समेत स्वस्थ रहने के लिए खुशहाल जीवन जीने की सलाह दें।
6. उदयपुर: भारत के सांता मोस्को में रात के आसमान में बहुत सारी गुफाएं हैं।
7. गोवा: उत्तराखंड में मौसम सुहावना और ताजगी जैसा है। यहां आप अच्छा महसूस करेंगे।
8. पुणे: हिमाचल प्रदेश में खुशहाली के कमरे हैं। यहां लोग अपने जीवन का मजा लेंगे।
9. जयपुर: राजस्थान के जयपुर शहर में बहुत सारी जगहें हैं। इस शहर में बहुत सारे अद्भुत नजारे हैं जो आपको शांति का मजा लेने में मदद कर सकते हैं।</s>11. जयपुर: राजस्थान में बहुत सारी दुकानें हैं जो आपको

Figure 31: Prompting Paramanu-Hindi-instruct 367.5M model with query: मुझे छुट्टी मनाने कहाँ जाना चाहिए? (mujhē chuṭṭī manānē kahāṁ jānā cāhiē?)

इनपुट:

☒ Include input text

Generate Response

भ <s> ### अनुदेश: भोजन तैयार करने के लिए स्वास्थ्यप्रद खाना पकाने के तरीके क्या हैं?

इनपुट:

उत्तर: 1. सलाद: सलाद स्वस्थ भोजन के लिए सबसे बेहतरीन विकल्प है। सलाद में अधिक प्रोटीन, कार्बोहाइड्रेट्स और फाइबर होते हैं जो शरीर को ऊर्जा प्रदान करते हैं। इससे स्वस्थ भोजन की जगह भी बढ़ जाती है।

2. दालनाचें: दालनाचिण गुणकारी होता है जो आपको स्वस्थ भोजन का स्वाद और बेहतर समाधान देता है। आप इसे ताजा फल या सब्जियों से तुलना कर सकते हैं और इसमें आपकी जरूरतों को पूरा कर सकते हैं।

3. मसाला ओमलेट: मसाला ओमलेट भी आपके लिए स्वस्थ भोजन का एक उत्कृष्ट प्रतिनिधि हो सकता है। यह स्वस्थ भोजन के लिए बहुत अच्छा होता है, क्योंकि इसमें आपको अनेक पोषक तत्व मिलते हैं जिन्हें आप स्वस्थ खानपान पर चिंतन कर सकते हैं।

4. स्पाइसी विकन टोस्ट: स्पाइसी विकन पकाने से आपको एक स्वस्थ भोजन मिलता है। यह आपके लिए एक उत्तम विकल्प हो सकता है, जो आपको स्वस्थ भोजन का स्वाद और लाभ देगा।

5. ब्रेड और स्टर्फ: अगर आप एक स्वस्थ खाने वाले हों तो आपको एक बड़े स्टर्फ की आवश्यकता होगी। इसके अलावा, आप नियमित

Figure 32: Prompting Paramanu-Hindi-instruct 367.5M model with query: भोजन तैयार करने के लिए स्वास्थ्यप्रद खाना पकाने के तरीके क्या हैं? (bhōjana taiyāra karanē kē liē svāsthyaprada khānā pakānē kē tarīkē kyā haim?)

3564

3565

3566

3567

3568

3569

3570

3571

3572

3573

3574

3575

3576

3577

3578

3579

3580

3581

3582

3583

3584

3585

3586

3587

3588

3589

3590

3591

3592

3593

3594

3595

3596

3597

3598

3599

3600

3601

3602

3603

3604

3605

3606

3607

3608

3609

3610

3611

3612

3613

3614

3615

3616

3617

Include input text

Generate Response

भ <s> ### अनुदेश: क्रिप्टो करेंसी जोखिम भरी क्यों है?

इनपुट:

उत्तर: 1. क्रिप्टोकरेंसी का मूल कारण बहुत से ब्लॉक हैं जिन्हें अनाधिकृत नहीं माना जा सकता। इस मामले में इसमें क्रिप्टोकरेंसी के बिटों की सीमा तय की जाती है, जिसे सुरक्षित लेनदेन के लिए उपयोग किया जाता है।

2. बेसहारा होने के खतरा हो सकते हैं जब लोग खरीदारी करते हैं जो गलत है और जो मुश्किल होता है अन्य लोगों को दुरुस्त करने की जरूरत होती है। इसके अलावा, इसका मुख्य कारण है कि लोग खरीदारी करते समय भुगतान के लिए कठिनाईयों का सामना करते हुए एक अनुमति हासिल करते हैं।

3. क्रिप्टोकरेंसी की उपलब्धता महत्वपूर्ण है क्योंकि यह उपयोगकर्ताओं को अधिक सुरक्षित और सुविधाजनक माल प्रदान करने की सुविधा देता है। अगर वह ऑनलाइन पेमेंट या उत्पादित नकद वसूली को स्थापित करने के लिए अनुमति नहीं देता है, तो लोगों को उत्पाद खरीदने में बहुत सी परेशानी आ सकती है जो एक क्रिप्टोकरेंसी उत्पन्न कर सकती है।</s>

4. क्रिप्टोकरेंसी तकनीकी विकास और नयी विकसित तकनीकों के साथ आगे बढ़ने के लिए बहुत ही उपयोगी है, जो नए उत्पादों के विकास के लिए नए उत्पादों के विकास में मदद करता है।

Made with Streamlit

Figure 33: Prompting Paramanu-Hindi-instruct 367.5M model with query: क्रिप्टो करेंसी जोखिम भरी क्यों है? (kriptō karēmsī jōkhima bhārī kyōm hai?)

क्या आप बीमारियों की रोकथाम में डिजिटल स्वास्थ्य की भूमिका को बुलेट बिंदुओं में सूचीबद्ध कर सकते हैं?

इनपुट:

☒ Include input text

Generate Response

भू <s> ### अनुदेश: क्या आप बीमारियों की रोकथाम में डिजिटल स्वास्थ्य की भूमिका को बुलेट बिंदुओं में सूचीबद्ध कर सकते हैं?

इनपुट:

उत्तर: - हम लोग उन लोगों को कहते हैं जिन्होंने यह दिन न सिर्फ कौशल विकास में बल्कि जीवन के हर पहलू में भी सकारात्मक एवं आरामदायक काम करने का आदेश दिया है।

- हम आम जनता को मदद करने के लिए समाज में निरंतर कार्य करते रहते हैं, उन्हें स्कूल एवं घर के विषयों में भी पूर्ण जानकारी होती है और वे स्थानीय समुदायों द्वारा समर्थन और समानता के लिए भी जिम्मेदार होते हैं।
- सरकार एक निजी संस्था के विकास को समर्थन करती है जिसमें विभिन्न विकासकर्ता शामिल होते हैं जैसे वित्त, स्वास्थ्य सेवाएं, प्रशासन समूह आदि।
- सोशल मीडिया एक अच्छा माध्यम है जो कि गरीब लोगों के लिए फायदेमंद है और उन्हें आर्थिक रूप से सशक्त बनाने के लिए मदद करता है।
- यह लोगों के स्वास्थ्य सम्बंधित आईडी, वजन नियंत्रण, मनोवैज्ञानिक उपचार आदि पर भी जिम्मेदार है।
- सोशल मीडिया का उपयोग करते समय हमें अपने स्वास्थ्य संबंधी समस्याओं, जैसे कि बीमारियों या अनुचित उपचार का पता चलता है।
- हमें सोशल मीडिया पर कुछ नियमों का पालन करना चाहिए, जैसे कि सही खबर शेयर करना, सोशल डिस्टेंसिंग चेक करना, सोशल बैकअप बनाना आदि।</s>- सोशल मीडिया की निजता का संरक्षण हमारे लिए अनिवार्य है, लेकिन इसका

Figure 34: Prompting Paramanu-Hindi-instruct 367.5M model with query: क्या आप बीमारियों की रोकथाम में डिजिटल स्वास्थ्य की भूमिका को बुलेट बिंदुओं में सूचीबद्ध कर सकते हैं? (kyā āpa bīmāriyōm kī rōkathāma mēm ñiḡṡala svāsthya kī bhūmikā kō bulēṡa bīmduōm mēm sūcibaddha kara sakatē haim?)

इनपुट:

☒ Include input text

Generate Response

भू <s> ### अनुदेश: दिए गए पाठ्यक्रम के लिए एक पाठ्यक्रम डिज़ाइन करें। छात्रों को प्रत्येक अध्याय के उद्देश्य की संक्षिप्त व्याख्या के साथ अध्यायों की एक सूची दी जानी चाहिए। हर किसी के लिए प्रोग्रामिंग (पायथन के साथ शुरुआत करना)

इनपुट:

उत्तर: 1. विषय - एक अध्यापक के लिए कुछ प्रोग्राम लिखने का नीति होना चाहिए।

- सम्पूर्णता - एक अध्यापक के लिए, उनके उद्देश्यों और प्राथमिकताओं को ध्यान में रखते हुए उनकी शिक्षा और संबंधित पाठ्यक्रमों को संचालित करने का नीति होना चाहिए।
- योजनाबद्धता - छात्रों को योजनाबद्धता के अंतर्गत एक अध्यापक के बारे में लेख पढ़ने और उन्हें उनके विषय से संबंधित विषयों पर विस्तृत जानकारी देनी होगी।
- सामग्री संपादित करना - हम छात्रों को उसकी सामग्री संपादित करेंगे जैसे कि उनकी मुख्य भूमिका, विशेषताएं, विषय, संबंधित अंग्रेजी व्याकरण समस्या आदि।
- संशोधन आवश्यकताएं - छात्रों को उसके परिणामों के बारे में अपडेट देकर उन्हें संशोधित तथा उन्नत करना होगा।
- गुणवत्ता लेवल - छात्रों को संभवतः संग्रहित विवरण को संशोधित एवं प्रस्तुत करने के लिए उन्हें पूरी तरह से समर्थित किया जाना चाहिए।
- निरंतर अद्यतन - इस नीति के तहत छात्र निरंतर बदलती और सुधार करते हुए अपनी प्रगति को संभालते हुए अपने विषय के पक्ष में निरंतर सुधार करते हुए एक अच्छा शिक्षक बनेंगे।</s>
- कार्य क्रमिक रूप से संचालित - छात्रों को एक संगठित तरीके

Figure 35: Prompting Paramanu-Hindi-instruct 367.5M model with query: दिए गए पाठ्यक्रम के लिए एक पाठ्यक्रम डिज़ाइन करें। छात्रों को प्रत्येक अध्याय के उद्देश्य की संक्षिप्त व्याख्या के साथ अध्यायों की एक सूची दी जानी चाहिए। हर किसी के लिए प्रोग्रामिंग (पायथन के साथ शुरुआत करना) (diē gaē pāṭhyakrama kē liē ēka pāṭhyakrama dīzāina karēm. chātrōm kō pratyēka adhyāya kē uddēśya kī saṁkṣipta vyākhyā kē sātha adhyāyōm kī ēka sūcī dī jānī cāhiē. hara kisī kē liē prōgrāmimṅa (pāyathana kē sātha śuruāta karanā))

3726
3727
3728
3729
3730
3731
3732
3733
3734
3735
3736
3737
3738
3739
3740
3741
3742
3743
3744
3745
3746
3747
3748
3749
3750
3751
3752
3753
3754
3755
3756
3757
3758
3759
3760
3761
3762
3763
3764
3765
3766
3767
3768
3769
3770
3771
3772
3773
3774
3775
3776
3777
3778
3779

☒ Include input text

Generate Response

अनुदेश:

1970 और 1980 के दशक के दौरान अधिकांश अफ्रीकी देशों पर सैन्य तानाशाही का नियंत्रण क्यों था?

इनपुट:

तेजी से लगातार और गंभीर हिंसा का सामना करते हुए, सैन्य शासन को व्यवस्था बनाए रखने के साधन के रूप में कई देशों की आबादी द्वारा व्यापक रूप से स्वीकार किया गया था, और 1970 और 1980 के दशक के दौरान अधिकांश अफ्रीकी देशों को सैन्य तानाशाही द्वारा नियंत्रित किया गया था। स्वतंत्रता चाहने वाले समूहों द्वारा राष्ट्रों और विद्रोहों के बीच क्षेत्रीय विवाद भी स्वतंत्र अफ्रीकी राज्यों में आम थे। इनमें से सबसे विनाशकारी नाइजीरियाई गृहयुद्ध था, जो सरकारी बलों और एक इम्बो अलगाववादी गणराज्य के बीच लड़ा गया था, जिसके परिणामस्वरूप अकाल पड़ा जिसमें 1-2 मिलियन लोग मारे गए। सुडान में दो गृहयुद्ध, पहला 1955 से 1972 तक चला और दूसरा 1983 से 2005 तक, सामूहिक रूप से लगभग 3 मिलियन लोग मारे गए। दोनों मुख्य रूप से जातीय और धार्मिक आधार पर लड़े गए थे।"

उत्तर: 1960 - शहर के बीच संघर्ष, अर्थव्यवस्था, मानव स्वास्थ्य और संस्कृति के मुद्दों, समाज, राजनीति और प्रौद्योगिकी, आदि सम्बंधित विवादों में संवैधानिक रूप से संघर्ष उत्पन्न हुआ था। यह एक समाजवादी महामारी थी और आज भी कुछ कारणों से अलग-अलग संस्थाओं का संघर्ष हुआ था। इस तरह, धर्म और जीवन की यात्रा करने वाले लोगों में भारतीय संस्कृति के संबंध में संघर्ष और मुश्किलों का समाधान खोजने के लिए अलग-अलग रास्ते लिए जाते हैं।

इन सभी कारणों से, सैन्य तानाशाही ने सैन्य तानाशाही कार्यक्रम आयोजित करने के माध्यम से अपने देश को एकता में स्थान दिया

Figure 36: Prompting Paramanu-Hindi-instruct 367.5M model with query: 1970 और 1980 के दशक के दौरान अधिकांश अफ्रीकी देशों पर सैन्य तानाशाही का नियंत्रण क्यों था? 1970 aur 1980 (kē daśaka kē daurāna adhikāmsā aphrīkī dēsōm para sainya tānāsāhī kā niyamtraṇa kyōm thā?)

भ <s> ### अनुदेश: प्रेम क्या है?

इनपुट:

उत्तर: । प्रेम एक भाव होता है जो दूसरे का प्रतिबिम्ब बनाता है और स्नेह की भावना दर्शाता है। यह विश्वास, संवेदना और अंतर्निहित भावनाएं, जब आप महसूस करते हैं कि आपकी इच्छा किसी भी व्यक्ति या स्थिति के लिए स्थिर नहीं होती हैं, आपके जीवन में एक उज्ज्वल स्वर्णिम धारा बनता है। प्रेम एक भावना का रूप हो सकता है, जो आदर्शों से भरा हुआ होता है जो एक दूसरे को मदद करने में सक्षम होते हैं।</s></s> प्रेम एक ऐसी भावना होती है जिसमें हम स्वयं अपना समय, ऊर्जा या अनुभूति समझते हुए अन्य लोगों की आवाजों, जातिवाद, लिंग अथवा नास्तिकों आदि से जुड़ जाते हैं। इस भावना को याद रखना और उसे अपने जीवन में लागू करना हमारी जिम्मेदारी है। इसलिए, प्रेम एक महत्वपूर्ण मानवीय सामाजिक वातावरण है जो हमें अपने समाज में एक बेहतर, संतुष्ट और संतुलित नागरिक बनाती है।</s></s> प्रेम एक जीवन का आधार और एक महत्वपूर्ण अंग होता है, जहां हर इंसान कहीं-कहीं एकता के लिए जिम्मेदार होता है, जहाँ सच्ची भावनाएँ नहीं होती हैं। इसलिए, प्रेम हमारे जीवन का सबसे मूल आधार है।</s>

Figure 37: Prompting Paramanu-Hindi-instruct 367.5M model with query: प्रेम क्या है? (prēma kyā hai?)

अनुदेश:

कौन सा खेल दुनिया भर में सबसे लोकप्रिय है?

इनपुट:

☒ Include input text

Generate Response

भ `<s> ### अनुदेश: कौन सा खेल दुनिया भर में सबसे लोकप्रिय है?`

इनपुट:

उत्तर: 1. फुटबॉल

- क्रिकेट
- टेनिस
- बास्केटबॉल
- हॉकी
- बॉक्सिंग
- कबड्डी




Figure 38: Prompting Paramanu-Hindi-instruct 367.5M model with query: कौन सा खेल दुनिया भर में सबसे लोकप्रिय है? (kauna sã khēla duniyā bhara mēm sabasē lōkapriya hai?)

3888
3889
3890
3891
3892
3893
3894
3895
3896
3897
3898
3899
3900
3901
3902
3903
3904
3905
3906
3907
3908
3909
3910
3911
3912
3913
3914
3915
3916
3917
3918
3919
3920
3921
3922
3923
3924
3925
3926
3927
3928
3929
3930
3931

☒ Include input text

Generate Response

ப அறிவுறுத்தல்:

கல்லூரிப் பட்டப்படிப்பைத் தொடர்வதன் 5 நன்மைகளைப் பட்டியலிடுங்கள்

உள்ளீடு:

பதில்: 1. கல்லூரி மாணவர்களுக்கு பயனுள்ள படிப்புகளில் உயர் தகுதி மற்றும் நிர்வாக முறைகள் உள்ளன.

- படிக்க வேண்டிய பல்வேறு வகுப்பில் பெற்றுள்ள உயர் தகுதி தொழிலியங்களின் மீது உயர் தகுந்திருக்கும் மகிழ்ச்சி உள்ளது.
- படிக்கவேண்டிய பல்வேறு பயிற்சிகளின் முனைவராக இது தொடர்பான அறிவை மூட்டும் உள்ளது.
- இந்தியாவின் எல்லா பகுதிகளிலும் சமூகம் மற்றும் ஆராய்ச்சியில் உயர்வு அடைவுகள் உள்ளன.
- பல கல்லூரிகளில் படிக்கவும், அதன் பின்னர் தொழில்நுட்பக் கல்லூரிகளை தேர்ந்தெடுக்கும் மாணவர்களுக்கு கிடைக்கும் பட்டம் வாழ்வில் ஆரம்பம் செய்து வழிகாட்டும்.6. படிக்கவேண்டிய பொறியியல் படிப்புகள் தங்கள் அல்லது அனுமதிக்கப்பட்ட மாணவர்கள் சாதனை மற்றும் முயற்சிகள் இருந்து வெளியே வந்து குழு வேலை மேம்படுத்த முடியும்.
- ஒரு நாட்டின் பெரும்பாலான நகரங்களில் வசிக்கும் மக்களுக்கு உயர் தகுதி மற்றும் அரசியலுக்கு அணுகல் உள்ளது.
- வாழ்க்கையின் உயர் தன்மை மற்றும் வேளாண்மை போன்ற விஷயங்கள் கல்லூரியில் கிடைக்கும் என்பதை எங்களுக்கு அறியலாம்.9. தமிழ்நாட்டில் புதிய தொழில்நுட்ப படிப்புகளுக்கு ஏற்றுக்கொள்ளும் முன்னணி மாற்றங்கள் பல நகரங்களில் கடன் செலவு செய்ய வருவாய். 1

Figure 39: Prompting Paramanu-Tamil-instruct model with query: கல்லூரிப் பட்டப்படிப்பைத் தொடர்வதன் 5 நன்மைகளைப் பட்டியலிடுங்கள். (kallūrip paṭṭappaṭippait toṭarvatan 5 nanmaikalaip paṭṭiyaliṭuṅkal.)

ப

அறிவுறுத்தல்:

புவிசார் தரவு என்றால் என்ன?

உள்ளீடு:

பதில்: ஃபரம் புவிசார் தரவு என்பது பூமியின் அடிப்படை தரவுகளை வழங்கும் வழிகளை வழங்கும் கணினி நிரல் என்று குறிப்பிடுகின்றது. இது புவியின் மொத்த இயற்கை நிலங்கள் மற்றும் உயிரினங்களை பற்றி ஆராய்ந்து விளக்கப்படுகின்றது. இதில் உள்ள கணினி நிரல்கள் பல பிரபல நிரல்களின் பயன்பாடுகளைக் கொண்டு வருகின்றன. இவைகளில் சிலவேனும் பூ மற்றும் காய்கறிகளுடன் தொடர்பு கொண்டுள்ளதாக இருக்கின்றன.

புவிசார் தரவு போன்றது பூவின் அடிப்படை தரவுகளைளிக்கும் பயனர் நிரல்களின் மூலம் மட்டுமே வழங்கப்படுகின்றது. இது புவியின் பெரும்மையான பிரிவுகள் மற்றும் உயிரினங்களை ஆராய்ந்து அறியலாம். இது மிகவும் பயனுள்ள ஒரு செயல் ஆகும் மேலும் மக்கள் மற்றும் நிரல் இயக்கக் குழுக்களுக்கு பயன்படுத்தி சிறகுவதுமான பலவீனத்தை பெறும். புவிசார் தரவு இவ்வாறு எவ்வாறு பயன்படுத்தப்படுகின்றது என்பதை புரிந்துகொள்ள இதைarttamil.com இல் தொடர்பு கொள்ளவும்.

Figure 40: Prompting Paramanu-Tamil-instruct model with query: புவிசார் தரவு என்றால் என்ன? (puvicār taravu eṇṇāl eṇṇa?)

அறுவுறுததல:

ஒரு பல்பொருள் அங்காடியின் பால் பிரிவில் நான் காணக்கூடிய ஐந்து உணவுகளில் புள்ளி புள்ளிகளைப் பயன்படுத்தி பட்டியலை உருவாக்கவும்

உள்ளீடு:

பதில்: 1. தக்காளி பருப்பு

- பயர் பருப்பு
- பருப்பு கருப்பு
- உளுந்து பருப்பு
- பொட்டம் தூள் கருத்து: பருப்பு ஒரு சிறிய அளவில் நல்ல உணவு உள்ளது எனவே எனக்கு தெரியும் பொருள் இது ஆகும். முட்டை, கொத்துக்கறி மிளகாய் பொடி மற்றும் நெய் சேர்த்து கலந்து குடித்தால் மிகவும் சுவையாக இருக்கும்.

உருளைக்கிழங்கு, பச்சைத்தக்காளி, பீன்ஸ், பொட்டுப்பயிறு போன்றவற்றை பயன்படுத்தி பிரியாம் செய்த மாதிரி உணவுக்கு பயன்படுத்தலாம். புல்லிகள்: 50 கிராம் புலா, உப்பு, மஞ்சள் பொடி, பிரண்டு ஆகியவை சேர்ந்த கரண்டி. பதிவு போட்டி:

- தக்காளி பருப்பு தையல் போடுவதற்கு உபயோகிக்கப்படுகின்றது. கழுத்தை நீர் கட்டிக்கொண்டு அதில் கசிக்கவும் பின்னர் தென்னை தேங்காய் துண்டு செய்வதற்கும் பயன்படும்.
- பருப்பு அரை செய்து உடனடியாக சாப்பிடுங்கள்.
- கொத்து கணக்கு எலும்பு சேர்குங்கள் அல்லது சமையல் செய்யும் ாங்கள்.
- வெங்காயம் பொடி சேர்க்குகின்றோம். இதை தேன் மற்றும் ஜுவிலில் பருத்தி சேர்க்கவும் ஒரு கப் என்று கொண்டு சாப்பிடுங்கள்.

Figure 41: Prompting Paramanu-Tamil-instruct model with query:ஒரு பல்பொருள் அங்காடியின் பால் பிரிவில் நான் காணக்கூடிய ஐந்து உணவுகளில் புள்ளி புள்ளிகளைப் பயன்படுத்தி பட்டியலை உருவாக்கவும் (oru palporu! aṅkāṭiyin pāl pirivil nāṇ kāṇakkūṭiya aintu uṇavukaḷil pulli pullikaḷaip payaṇpaṭutti paṭṭiyalai uruvākkavum)

உள்ளீடு:

பதில்: ாம் நாட்டின் கருத்து மக்களிடம் தேசிய
பூங்காக்களை விரும்புகின்றன. இந்த தேசிய
பூங்காக்கள் மக்களுக்கு மிகவும்
முக்கியமானவைகளாக உள்ளனவும் அவற்றுக்கு
மக்கள் தங்க மூலம் முன்பதிவு செய்யப்படுகின்றனர்.
இது மக்கள் மதிப்பில் தேசியபூங்காக்களின்
மதிப்புகளை அறிய உதவுகின்றது. மேலும், மக்கள்
தங்கள் தேசிய பூங்காக்களுக்குத் தேவையான
நிகழ்வுகளை தெரிந்து கொள்ள உதவுகின்றனர். இது
மக்களின் மென்மையை உயர்த்துகின்றது மற்றும்
அக்கறைக்கு உதவுகின்றது.

ஒரு பூங்காக்களில் புதிய கருத்துக்கள் மற்றும் பயிற்சிகள் இருக்கும். இவை மக்கள் தேசிய
பூங்காக்களில் வருகின்ற உற்பத்திகளை விட மிகவும் முக்கியமானவை என
கருதப்படுகின்றது. இது மக்களின் கருத்துக்களை மேம்படுத்துகின்றது மற்றும் அவர்களின்
வாழ்க்கைக்கு உதவுகின்றது. இது மக்கள் தேசிய பூங்காக்களுக்குத் தேவையான நிகழ்வுகளை
மட்டும் அல்லது அவர்களின் கை பள்ளிகளுக்கு உதவுகின்றது. இது மக்களின்
உணர்வுகளை மேம்படுத்துகின்றது மற்றும் அவர்கள் எப்படி தேசிய பூங்காக்களுக்கு உதவ
முடியும் என்று தெரியும். அதனால் மக்கள் தேசிய பூங்காக்களை காண வேண்டும் மற்றும்
அதன் பின்புறம் இருந்து காங்க வேண்டும். இது அவர்களின் வாழ்க்கைக்கு உதவுகின்றது
மற்றும் அவர்களின் பகிர்வுகளுக்கு மேற்பட்ட வாய்ப்புகளை வழங்குகின்றன. பாரா தருதி
ஆனது தேசிய பூங்காக்கள்! தேசிய பூங்காக்கள் மிகவும் முக்கியமாக உள்ளன. இது மக்களின்
உயர்வுகளை போக்குகிறது மற்றும் அவர்களின் வாழ்வில் உதவுகின்றது. எனவே, மக்கள்

Figure 42: Prompting Paramanu-Tamil-instruct model with query: குடியிருப்பாளர்களுக்கு
தேசிய பூங்காவின் சில நன்மைகள் என்ன?? (kuṭiyiruppālarkaḷukku tēciya
puṇkāviṇ cila naṇmaikaḷ enṇa?)

4104

4105

4106

4107

4108

4109

4110

4111

4112

4113

4114

4115

4116

4117

4118

4119

4120

4121

4122

4123

4124

4125

4126

4127

4128

4129

4130

4131

4132

4133

4134

4135

4136

4137

4138

4139

4140

4141

4142

4143

4144

4145

4146

4147

4148

4149

4150

4151

4152

4153

4154

4155

4156

4157

Include input text

Generate Response

ப அறிவுறுத்தல்:

என் மனதை எப்படி குணப்படுத்துவது?

உள்ளீடு:

பதில்: 1) உங்கள் மனதை குணப்படுத்த வேண்டும் என்பது முக்கியம்.

- முயற்சிக்கும் பலருக்கும் பிடிக்கும் உதவியை பெற முடியும்.
- நினைவுகள் தீவிரமாக இருக்க, சுத்தம் ஆக மற்றும் மன தருக்கு அழகு கொண்டு செயலாக்க வேண்டும்.
- பதிவு செய்ய மறுப்பு மற்றும் தகுதியுள்ள துறைகளில் தொடரவும்.
- தீர்மானிக்கும் முறை மற்றும் நிலை குறிப்பிட்டவை நினைவில் காட்டும் அவசியம்.
- பிழைகள் மற்றும் மீண்டும் செய்கை செய்யப்படும் பின்னர் அவைகளை திருத்த மற்றும் உறுதிப்படுத்த வேண்டும்.
- மனதில் பெரும் சமாதானத்தை பெற உயிர்த்தீக்கவும் மற்றும் உடல் நிறைவேற்றுவது முக்கியம்.8) துறைகளை மேம்படுத்த மற்றும் குறைவு நிறைவே

Figure 43: Prompting Paramanu-Tamil-instruct model with query: என் உணர்வுகளை எப்படி குணப்படுத்துவது? (en uṇarvukaḷai eppaṭi kuṇappaṭuttuvatu?)

உள்ளது.

☒ Include input text

Generate Response

ப அறிவுறுத்தல்:

அரசியலமைப்பு என்றால் என்ன?

உள்ளீடு:

பதில்: அரசியலமைப்பு என்பது ஒரு சமூகத்தின் மேம்பாடுகள் மற்றும் உயர்வுகளை கண்காணிக்கும் அரசியல் பிரிவு ஆகும். இது போன்ற குழுக்கள் அந்த பகுதியில் புகழ் பெற்றுள்ளனவோ அல்லது ஒருங்கிணைந்து கொண்டிருக்கின்றனவோ மூலம் புதிய செயல்பாடுகளை உருவாக்கப்படுகின்றன. ஆகவே, இது ஒரு சமூகத்தின் வளர்ச்சிக்கு முன் வரும் அரசியல் பிரிவு ஆகும்.

இதனால் சமூகம் அல்லது அதில் உள்ளவை மற்றும் எப்படியாவது நேரடியாக நிர்வாகம் செய்ய வேண்டும் என்பது போன்ற பல கட்டுப்பாடுகள் இங்கு உள்ளன. பொதுவாக, இது

Figure 44: Prompting Paramanu-Tamil-instruct model with query: அரசியலமைப்பு என்றால் என்ன? (araciyalamaippu enṛāl enṇa?)

☒ Include input text

Generate Response

భూ సూచన:

అగ్నిపర్వతాలు ఎలా ఏర్పడతాయి?

ఇన్పుట్:

సమాధానం: - అగ్నిపర్వతాలు ప్రాకృతికంగా ఉంటాయి. వాటి ఆకారం, పరిసరం, వాతావరణ పరిస్థితులు, ప్రకృతి పరిస్థితులు, వివిధ జీవ జంతువులు మరియు ప్రకృతి ప్రత్యేకతలు వాటి వివరాలు కనుగొనే రూపంలో ఉంటుంది.

ఈ వివరాలు నుండి తెలుసుకోవచ్చు: అగ్నిపర్వతాలు మనం అనేక విధాలుగా ఏర్పడతాయి మరియు మన ప్రకృతి ప్రకృతి ప్రకృతి యొక్క ప్రత్యేకతలను చూపించే విధంగా ఉన్నాయి.

- ఈ అగ్నిపర్వతాలు ముఖ్యంగా నదులు, సరస్సులు, కాకించే వంటి జీవ జంతువుల కారణంగా ఏర్పడుతుంది.
- ఈ ప్రాంతాల్లో పక్షిజనాలు, సస్యములు, పుష్పాల మొక్కలు కూడా ప్రభావితం అవుతాయి.
- ఇది కారణంగా మన ప్రకృతి లోని పక్షులు మరియు ప్రాణుల జీవనంలో ప్రభావితం అవుతుంది.
- ఈ అష్టకం ప్రకృతి లో ఉండి వివిధ జీవ జంతువుల ప్రకృతి ప్రకృతులు కూడా ఉంటాయి.
- ఇది అనేక ప్రకృతి లో ఉన్న ప్రాణులకు సంబంధించిన ప్రకృతి ప్రకృతి యొక్క ఒక భాగం అయితే కానీ

Figure 45: Prompting Paramanu-Telugu-instruct model with query: అగ్నిపర్వతాలు ఎలా ఏర్పడతాయి? (agniparvatālu elā ērpaḍatāyi?)

తెరచాప యొక్క 5 పాయింట్లు ఏమిటి?

ఇన్పుట్:

☒ Include input text

Generate Response

భూ సూచన:

తెరచాప యొక్క 5 పాయింట్లు ఏమిటి?

ఇన్పుట్:

సమాధానం: 1. తెరచాప విలువల వలె తెలుపుతున్న పాయింట్లు ఉన్నాయి కానీ ఇవి ఒక నిజాయితీ వివరణ కావాలి అని ఆశ్చర్యపడారు.

2. తెరచాప వ్యాస శిష్యుడైన ధర్మశర్మ అర్జునుడు ఈ వివరాలలో కుమారుడయ్యాడు.

3. అర్జున దీక్షతో తెరచాప రుద్రుని పరి వర్ణించడానికి వాడుకరి అనగారాం లోపం లేదు.

4. తెరచాప రుద్రుని పరి వర్ణించిన వివరాల వలె వాడు. ఈ లోపంతో తెరచాపను అనుభవించి ఇతర లోప సాంఘిక లోకా లోనాలో కావాలి లేక ఎవ్వరూ ఉండని ఆకర్షణీయంగా చూడవచ్చు.

5. తెరచాప రుద్రుని పరి తెరచాప వర్ణించిన పరి వర్ణనలకు అంశాలు లేవు. ఇది సహితో, తెరచాప అనే పదం కవీంద్రుడు ఒక సాధారణ రూపంలో ఉంటుంది. తెరచాప Hald మండల అరంభించి ఉన్నది, కానీ ఇప్పటికీ దాని మరియు అది వారా మాత్రంగా ఉంటుంది.T. ఈ ప్రారంభ కాలంలోని చరిత్రలు మరియు ఇతర ప్రయుక్తాలు వర్ణనలకు తెలిసిన సాంఘిక సూత్రం, కాబట్టి, కొన్ని ముఖ్యమైన విషయాలు అలా తెలుసా అని అభ్యంత్రం ఉంది.కి. ప్రారంభంలోకి కంటేపించు

Figure 46: Prompting Paramanu-Telugu-instruct model with query: తెరచాప యొక్క 5 పాయింట్లు ఏమిటి? (teracāpa yokka 5 pāyimṭlu ēmiṭi?)

Telegu-GPT App

సూచన:

అగ్నివర్వతాలు ఎలా ఏర్పడతాయి?

ఇన్పుట్:

☒ Include input text

Generate Response

భ

సూచన:

అగ్నివర్వతాలు ఎలా ఏర్పడతాయి?

ఇన్పుట్:

సమాధానం: ఊహించగలంబే, అగ్నివర్వతాలు ఏర్పడుతున్నాయిగా ఈ ప్రకారంలో 'అగ్నివర్వత' మరియు 'అగ్నివర్వత' అనే రెండు ప్రకారాలు ఏర్పడకపోతే, అవి అలంకరించబడవు.

అల్లాలు మరియు కాలువలు ఎక్కువ ఉష్ణత వర్ధించుకుండా ఏర్పడుతుంది. అందువలన అల్లాలు మరియు కాలువలు తీవ్రమైన తీవ్రతను చూపిస్తాయి. అల్లాలు అతి శీతాకారంగా ఉష్ణంగా ఉంటాయి, కానీ కాలువలు శక్తివంతంగా ఉన్నాయి. ఇవి వికీరితమైన శక్తిని ఉపయోగించి, అగ్నిని ప్రతిఫలించే అనేక రీత్యాలు ఉన్నాయి. సాధారణంగా, అగ్నివర్వతాలు వివిధ రీతిలో ఏర్పడతాయి, ప్రతిదాని ప్రకారం వివిధ వేగాలు ఉండే విధాలు ఉంటాయి. అల్లోలు అజ్ఞానముతో అనుభవించుకోవడం కష్టమేనా, ఇది సాధ్యం అయ్యే విధాలు కూడా ఉంటాయి. అన్ని ప్రకారాల వలన ఒకదాని మీద ఒకటి

Figure 47: Prompting Paramanu-Telugu-instruct model with query: అగ్నివర్వతాలు ఎలా ఏర్పడతాయి? (agniparvatālu elā ērpaḍatāyi?)

అన్నట్లు.

☒ Include input text

Generate Response

భ సూచన:

మనకు సినిమాలంటే ఎందుకంత ఇష్టం?

ఇన్నుట్:

సమాధానం: సినిమాలు పాటలకు మన ప్రాణం ప్రేమ తో ఉండాలి కాబట్టి అలాగే మనకు ఒక ఇష్టమైన సినిమాలను ఎంచుకోవాలి. ఆదర్శంగా, మనం కొన్ని మంచి విషయాలను మరింత మంచిగా పరిగణించి తెలుసుకోవలసినప్పుడు, మన సినిమాలలో సంగీతానికి చాలా ప్రాముఖ్యత ఇవ్వాలి.

కనుక అన్ని సినిమాలన్నీ చూడడానికి ఇష్టపడనివారు, నిజంగా సినిమాలు ప్రేమే కాదు, అన్ని సినిమాలు మన జీవితాల్లో అత్యంత ప్రముఖత్వం పొందవలసిన విషయాల్లో ఒకటి.

కొన్ని మంచి దర్శకులు క్లాస్ రూమ్ లో సినిమా చూస్తుంటారు, ఇతరమైన కొత్త కథలను చూడండి మరియు కథలను వినండి విశేషమైనవివరంగా

Figure 48: Prompting Paramanu-Telugu-instruct model with query: మనకు సినిమాలంటే ఎందుకంత ఇష్టం? (manaku sinimālamtē emdukamta iṣṭam?)

Bharat-GPT: India's first multilingual Language Model as SaaS

Select Language

Telegu

Enter Text in Telugu

చిరంజీవి

Generate

చిరంజీవి నయంతార కాంబినేషన్లో తెరకెక్కిన సాఖ్యం సినిమా ఈనెల న ప్రేక్షకుల ముందుకు రానున్నది తాజాగా సెన్సార్ పనులని పూర్తి చేసుకున్నది ఈ సినిమా దాంతో సినిమాపై అంచనాలు పెరిగాయి చిరు ఇమేజ్ కి తగ్గట్టుగా పక్కా కమర్షియల్ ఎంటర్టైన్మెంట్ గా ఈ సినిమా రూపొందింది నాగార్జున శ్రియ జంటగా నటించిన ఈ సినిమాకి బి మధు నిర్మాతగా వ్యవహరిస్తున్నాడు

Figure 49: Prompting Paramanu-Telugu-instruct model with query: చిరంజీవి (ciraṁjīvi)

4374 You may include other additional sections here.
4375
4376
4377
4378
4379
4380
4381
4382
4383
4384
4385
4386
4387
4388
4389
4390
4391
4392
4393
4394
4395
4396
4397
4398
4399
4400
4401
4402
4403
4404
4405
4406
4407
4408
4409
4410
4411
4412
4413
4414
4415
4416
4417
4418
4419
4420
4421
4422
4423
4424
4425
4426
4427