# IndicLLMSuite: A Blueprint for Creating Pre-training and Fine-Tuning Datasets for Indian Languages

**Mohammed Safi Ur Rahman Khan**[*1]  **Priyam Mehta**[*1]  **Ananth Sankar**[1]
**Umashankar Kumaravelan**[1]  **Sumanth Doddapaneni**[1,2]  **Suriyaprasaad G**[1,4†]
**Varun Balan G**[1,5†]  **Sparsh Jain**[1,6†]  **Anoop Kunchukuttan**[1,2,3]
**Pratyush Kumar**[1,2,7]  **Raj Dabre**[2,8]  **Mitesh M. Khapra**[1,2‡]

[1]Nilekani Centre at AI4Bharat    [2]Indian Institute of Technology, Madras    [3]Microsoft
[4]Sant Longowal Institute of Engineering and Technology    [5]IIIT D&M Kancheepuram
[6]Maharaja Agrasen Institute of Technology    [7]Sarvam AI
[8]National Institute of Information and Communications Technology, Kyoto, Japan

## Abstract

Despite the considerable advancements in English LLMs, the progress in building comparable models for other languages has been hindered due to the scarcity of tailored resources. Our work aims to bridge this divide by introducing an expansive suite of resources specifically designed for the development of Indic LLMs, covering 22 languages, containing a total of 251B tokens and 74.8M instruction-response pairs. Recognizing the importance of both data quality and quantity, our approach combines highly curated manually verified data, unverified yet valuable data, and synthetic data. We build a clean, open-source pipeline for curating pre-training data from diverse sources, including websites, PDFs, and videos, incorporating best practices for crawling, cleaning, flagging, and deduplication. For instruction-fine tuning, we amalgamate existing Indic datasets, translate/transliterate English datasets into Indian languages, and utilize LLaMa2 and Mixtral models to create conversations grounded in articles from Indian Wikipedia and Wikihow. Additionally, we address toxicity alignment by generating toxic prompts for multiple scenarios and then generate non-toxic responses by feeding these toxic prompts to an aligned LLaMa2 model. We hope that the datasets, tools, and resources released as a part of this work will not only propel the research and development of Indic LLMs but also establish an open-source blueprint for extending such efforts to other languages. The data and other artifacts created as part of this work are released with permissive licenses at `https://github.com/AI4Bharat/IndicLLMSuite`.

## 1 Introduction

Building Large Language Models (LLMs) is an inherently data-intensive process requiring a comprehensive set of resources for pre-training (Raffel et al., 2020; Xue et al., 2021; Gao et al., 2021; Penedo et al., 2023; Nguyen et al., 2023a; Abadji et al., 2022) and fine-tuning (Longpre et al., 2023; Conover et al., 2023; Köpf et al., 2023; Ding et al., 2023a). The last year has seen remarkable progress in building English LLMs, thanks to open-source models (Touvron et al., 2023a;b; Jiang et al., 2023; 2024a; Almazrouei et al., 2023) developed using comprehensive datasets containing such resources. Nonetheless, this progress has largely bypassed low and mid-resource languages due to the lack of data resulting from the lack of open-source pipelines for curating data for such languages from diverse sources such as websites (which require crawling and extraction), books

---

[*]Equal Contribution. All author contributions listed in 7.

[†]Work done during an internship at Nilekani Center at AI4Bharat

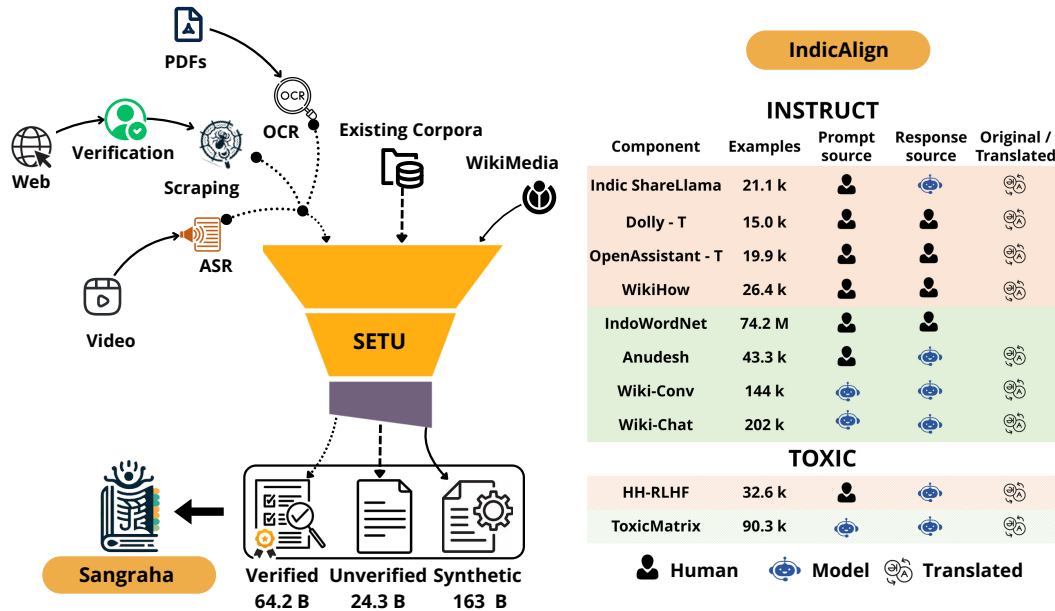[‡] Corresponding Author: Mitesh Khapra (miteshk@cse.iitm.ac.in)

Figure 1: Overview of the different components present in INDICLLMSUITE.

(which require OCR) and videos (which require transcription). Further, for instruction fine-tuning, English LLMs now rely on model-generated data such as ShareGPT[1], Self-Instruct (Wang et al., 2023b), Evol-Instruct (Xu et al., 2023a), UltraChat (Ding et al., 2023a), etc. However, for low and mid resource languages, this option is not available due to lack of high-quality LLMs, leading to a chicken and egg problem, further widening the gap between the *haves* and the *have-nots*.

A case in point is that of languages from the Indian sub-continent, which collectively are spoken by over 1.4 billion people. We focus on the 22 languages recognised in the 8th schedule of the Indian constitution. These languages, despite their significant number of speakers, receive minimal representation in the training datasets and tokenizers of current open-source LLMs (Touvron et al., 2023b; Jiang et al., 2024a; Almazrouei et al., 2023) leading to a notable exclusion of their rich cultural contexts and nuances. In this work, we address this disparity by making the following contributions, as summarised in Figure 1:

**1. SANGRAHA**: Pre-training data containing 251B tokens[2] summed up over 22 languages, summarized in Table 1, extracted from curated URLs, existing multilingual corpora, and large-scale translations.

**2. SETU**: Spark-based (Zaharia et al., 2016) distributed pipeline customized for Indian languages for extracting content from websites, PDFs and videos, with in-built stages for cleaning, filtering, toxicity removal and deduplication.

**3. INDICALIGN - INSTRUCT**: A diverse collection of 74.7 million prompt-response pairs across 20 languages, summarized in Table 7, collected through four methods: aggregating existing Instruction Fine-Tuning (IFT) datasets, translating English datasets into 14 Indian languages using an open-source translation model, creating context-grounded conversations from India-centric Wikipedia articles using open source LLMs, and establishing a crowdsourcing platform called *Anudesh* for prompt collection. We also create a novel IFT dataset to teach the model language and grammar, by leveraging INDOWORDNET (Bhattacharyya, 2010), a lexically rich but rather neglected resource in the era of LLMs.

---

[1] https://sharegpt.com/

[2] We built a custom tokenizer which supports English and Indian languages and has an average fertility of 1.3 to 2.79 across the 22 languages. We use this tokenizer for all the reported statistics unless mentioned otherwise.

**4. INDICALIGN - TOXIC**: 123K pairs of toxic prompt and non-toxic responses generated using open source English LLMs and translated to 14 Indian languages for safety alignment of Indic LLMs.

We collectively refer to the above as INDICLLMSUITE. We try to balance quality and quantity while acknowledging recent trends of using synthetic data for building powerful LLMs for English (Gunasekar et al., 2023; Li et al., 2023d) as well as low resource languages (Nguyen et al., 2023b; Li et al., 2023c). To ensure quality, we take help from humans to verify websites to flag noisy or machine translated content and to create toxicity lists for Indian languages. On the other hand, to ensure explicit representation of prompt-response pairs grounded in the Indian context we take the help of powerful open-source LLMs to generate grounded conversations from India-centric Wikipedia articles. We recognize the need to represent diverse knowledge and alignment information in Indic languages for better performance of LLMs in Indic languages. Hence, we undertake large-scale machine translation of rich English resources like Wikimedia as well as English fine-tuning datasets into Indian languages using SOTA open-source MT models.

We thus balance source original data with translated and LLM-generated data to create the above collection.

We believe that these choices can be replicated across other languages to create LLMSuites. All the code, tools and datasets developed as a part of this work will be publicly released and hopefully advance the development of LLMs for Indian languages. Given that LLM training is an expensive exercise, we plan to undertake community-effort to train LLMs, where multiple groups can pool together computing resources to build a high-quality Indic language LLM.

## 2 RELATED WORKS

We organise the Related Work into 3 sections in line with our main contributions.

### 2.1 MULTILINGUAL DATASETS

Wikipedia has consistently served as the *go to* repository of multilingual data and continues to be an important contributor for training data. Prior works such as OSCAR (Abadji et al., 2022), CC100 (Conneau et al., 2020) (encompassing 100 languages), and MC4 (Xue et al., 2021) (encompassing 101 languages) have been instrumental in generating data for a large set of languages through the meticulous processing of common crawl dumps[3]. The overarching methodology across all the works includes the systematic extraction of data, language identification, and subsequent stages of filtering and deduplication. Notably, CCNet (Wenzek et al., 2020) employs filters based on n-gram language models, while OSCAR and mC4 (Xue et al., 2021) leverage heuristic-based filtering mechanisms. Efforts like Samanantar (Ramesh et al., 2022) and ROOTS (Laurençon et al., 2023) have underscored the importance of aggregating existing datasets as an initial step towards multilingual corpus construction. ROOTS Laurençon et al. (2023), with a specific focus on 59 languages, has additionally drawn attention to potential data duplications across disparate pipelines. Building on these works, MADLAD-400 (Kudugunta et al., 2023) extends to 419 languages, while introducing human audit of data and iterative refinement processes, coupled with language family-specific filters. CULTURAX (Nguyen et al., 2023a), merges MC4 and all versions of OSCAR, followed by a rigorous cleaning pipeline to produce corpora in 167 languages.

### 2.2 DATA CURATION

Kreutzer et al. (2022a) has unilaterally showed the importance of auditing datasets. They discovered problems like wrong language, bad quality, offensive content, etc. Other works Rae et al. (2021); Penedo et al. (2023) reinforced the idea that using clean data is key to making better models.

**Sources & Scraping.** Majority of data curation pipelines start with Common Crawl as the internet source, followed by a series of cleaning steps to produce the final dataset (Raffel et al., 2020; Xue et al., 2021; Conneau et al., 2020; Penedo et al., 2023). Wikipedia is another common source of

---

[3]https://commoncrawl.org/

*high quality* data that is used across many works (Gao et al., 2021; Computer, 2023; Soldaini et al., 2024). Research (Brown et al., 2020; Touvron et al., 2023a) further demonstrates the benefits of augmenting datasets with high-quality content, such as Wikipedia and Books corpus, for improved language model training. The first major hurdle for scraping data is extracting the actual text content from the HTML files and filtering out any boilerplate and unwanted HTML content. Tools like Trafilatura (Barbaresi, 2021) and jusText (Endrédy & Novák, 2013) are widely used in recent works like RefinedWeb (Penedo et al., 2023) and Pile (Gao et al., 2021).

**Language Identification.** The subsequent challenge in the data cleaning process involves language detection, a task for which several LID tools have been employed in previous studies, namely, CLD3 (Botha et al., 2017)[4], langdetect[5], fasttext (Wenzek et al., 2020; Costa-jussà et al., 2022), SSLID (Kudugunta et al., 2023), among others. These tools demonstrate proficiency in identifying 55 to 500 languages; however, they fall short in encompassing all Indian languages considered in this study. Complications arise, particularly in the case of low-resource languages sharing the same script, leading to potential instances of mislabeling (e.g., Hindi and Marathi languages), necessitating the need for focused language family-specific detectors (Madhani et al., 2023a).

**Filtering.** One of the pivotal components in the data processing pipeline involves implementing various heuristics to ensure the quality of the processed data. These heuristics encompass rule-based approaches, targeting elements such as punctuation, repetitions, special characters, templated content, code snippets, etc. (Wenzek et al., 2020; Raffel et al., 2020; Xue et al., 2021; Laurençon et al., 2023; Rae et al., 2021). Additionally, model-based tools, including n-gram language models perplexities (Wenzek et al., 2020; Laurençon et al., 2023; Nguyen et al., 2023a) and ML classifier-based approaches (Brown et al., 2020), play a crucial role in filtering out data that falls below certain quality thresholds. A prominent emphasis within the data cleaning process is the removal of toxic and harmful content, a priority shared across various datasets. Techniques employed for this purpose often include the utilization of word lists (Raffel et al., 2020), blocklists of URLs (Penedo et al., 2023), and leveraging Google Safe Search (Rae et al., 2021). Notably, MADLAD-400 (Kudugunta et al., 2023) adopts language-specific heuristics, incorporating a manual audit of a small portion of data to enhance the cleaning process. This multifaceted approach underscores the importance of employing diverse strategies to ensure the reliability and safety of processed data in various contexts.

**Deduplication.** Recent research findings shed light on the significant impact of deduplication on language models (LMs). Carlini et al. (2023) demonstrated that deduplication effectively reduces memorization within LMs, while Lee et al. (2022) highlighted its role in enhancing LM performance. Furthermore, recent work Hernandez et al. (2022) underscored the detrimental effects of data repetition on model performance, particularly as model size increases. These insights underscore the critical importance of rigorous deduplication procedures in dataset preprocessing. Common approaches to deduplication encompass various techniques such as deduplication based on URLs, fuzzy techniques like MinHash (Broder, 1997) and SimHash (Charikar, 2002), alongside embedding techniques like those proposed by SemDedup (Abbas et al., 2023).

With the rapid expansion of LLM development, a significant challenge arises from the proliferation of machine-generated text on the internet. It is important to devise effective strategies for filtering out such content to ensure the creation of high-quality data intended for human consumption. This necessitates a careful curation process that distinguishes between machine-generated and human-written text. Additionally, another issue pertains to the crawling and inclusion of benchmark/test data in the pre-training mixture i.e., Data Contamination (Sainz et al., 2023; Golchin & Surdeanu, 2024). BigBench (Srivastava et al., 2022) advocates for filtering based on matching key strings, while recent research suggests encrypting benchmark data with passwords before distribution (Jacovi et al., 2023).

## 2.3 SUPERVISED FINE-TUNING DATASETS

Pre-training demands huge amounts of data to effectively train on diverse linguistic patterns, while fine-tuning, specifically instruction tuning, necessitates comparatively smaller yet high-quality

---

[4]https://github.com/google/cld3/
[5]https://github.com/shuyo/language-detection

datasets (Zhou et al., 2023). Broadly there exist two approaches for creating these datasets: (i) *human-generated*, involving the manual input of humans for prompt creation and/or answer generation; and (ii) *model-generated*, where both prompts and answers are generated by models.

**Human Generated** Dolly (Conover et al., 2023) and Open-Assistant (Köpf et al., 2023) are created through a comprehensive human annotation process from start to finish. Dolly, for instance, was curated by around 5000 Databricks employees to create nearly 15k instructions, whereas Open-Assistant emerged from a collaborative crowd-sourced initiative, accumulating 10k conversations spanning across 35 languages. These datasets encompass a wide array of tasks including question-answering, creative writing, classification, etc. On the flip side, datasets like ShareGPT[6] and Wild-Chat Zhao et al. (2024) are created by gathering real human interactions with ChatGPT[7]. Similarly, HC3 Guo et al. (2023) gathered both human and ChatGPT responses for questions generated by humans on public datasets. These variations span diverse topics due to the extensive participation in creating these datasets.

**Model Generated** Distilling data from powerful models for the creation of datasets tailored for SFT tasks has become a widespread practice. Self-Instruct (Wang et al., 2023a) introduced an almost annotation-free data creation pipeline. Starting with just 175 seed prompts, they expand the dataset to 52k instructions iteratively. Alpaca (Taori et al., 2023) streamlined the Self-Instruct process, leading to increased efficiency and cost reduction. Taking a further stride, Evol-Instruct (Xu et al., 2023a) introduces a novel method of evolving the seed prompts iteratively along various axes. This resulted in the generation of an extensive dataset comprising 250k instructions. Further works generate data by simulating conversations between two or more model agents, with at least one acting as a User and the other acting as the Assistant. CAMEL (Li et al., 2023a), for instance, crafted 115k instructions via multi-agent role play, while Ultrachat (Ding et al., 2023a) produced 1.5 million multi-turn dialogues. Similarly, Baize (Xu et al., 2023c) generated 115k dialogues through self-chat with ChatGPT.

## 3 SANGRAHA

In this section, we describe the composition and curation process of SANGRAHA spanning verified (64B), unverified (24B), and synthetic (162B) content for a total of 251B tokens. Table 1 shows the language level tokens distribution in each of the splits.

### 3.1 SANGRAHA VERIFIED

We introduce SANGRAHA VERIFIED, a high-quality dataset, which adds human verification at various stages of its curation. A major chunk of this includes data crawled from high-quality, manually verified Indic language websites. Additionally, recognizing the fact that a significant amount of Indic language text is locked in PDFs and audio, we also collect data from various books/documents and videos resulting in a total of 64B tokens. Table 2 shows the language level tokens distribution across the Web, PDF, and Speech data in SANGRAHA VERIFIED.

#### 3.1.1 WEB DATA

Our web data, constituting most of Sangraha, diverges from traditional Common Crawl-based approaches by prioritizing data quality. This involves manual verification of each website before scraping. We adopt a three-fold strategy to collect a comprehensive collection of verified websites for scraping. Firstly, we extend the efforts of Kakwani et al. (2020) and Doddapaneni et al. (2023) of discovering web sources using existing news repositories and automated web searches using popular keywords to discover a large list of Indic language websites. But, unlike the previous efforts, we do not restrict ourselves to just news websites. Secondly, we identify various domains such as Indian Culture, Food, Health, Travel, among others and enlist volunteers to gather websites within these domains, prioritizing those in Indic languages or in English but pertaining to Indian context. Thirdly,

---

[6]https://sharegpt.com/
[7]https://chat.openai.com/

| Code | SV | SS | SU | Total Tokens |
|------|-----|-----|-----|------|
| asm | 292.1 | 11,696.4 | 17.5 | 12,006.0 |
| ben | 10,604.4 | 13,814.1 | 5,608.8 | 30,027.5 |
| brx | 1.5 | - | - | 1.5 |
| doi | 0.06 | - | - | 0.06 |
| eng | 12,759.9 | - | - | 12,759.9 |
| gom | 10.1 | - | - | 10.1 |
| guj | 3,647.9 | 12,934.5 | 597.0 | 17,179.4 |
| hin | 12,617.3 | 9,578.7 | 12,348.3 | 34,544.3 |
| kan | 1,778.3 | 12,087.4 | 388.8 | 14,254.5 |
| kas | 0.5 | - | - | 0.5 |
| mai | 14.6 | - | - | 14.6 |
| mal | 2,730.8 | 13,130.0 | 547.8 | 16,408.6 |
| mar | 2,827.0 | 10,816.7 | 652.1 | 14,295.8 |
| mni | 7.4 | - | - | 7.4 |
| npi | 1,822.5 | 10,588.7 | 485.5 | 12,896.7 |
| ori | 1,177.1 | 11,338.0 | 23.7 | 12,538.8 |
| pan | 1,075.3 | 9,969.6 | 136.9 | 11,181.8 |
| san | 1,329.0 | 13,553.5 | 9.8 | 14,892.3 |
| sat | 0.3 | - | - | 0.3 |
| snd | 258.2 | - | - | 258.2 |
| tam | 3,985.1 | 11,859.3 | 1,515.9 | 17,360.3 |
| urd | 3,658.1 | 9,415.8 | 1,328.2 | 14,402.1 |
| tel | 3,706.8 | 11,924.5 | 647.4 | 16,278.7 |
| **Total** | 64,306.1 | 162,707.9 | 24,307.7 | **251,321.0** |

Table 1: Number of tokens (in Millions) in each split of Sangraha. (SV: SANGRAHA VERIFIED, SS: SANGRAHA SYNTHETIC, SU: SANGRAHA UNVERIFIED). We represent the languages in this document using the ISO 639-3 standard codes



(a) Accepted Website statistics - Domain information



(b) Rejected Website statistics - Rejection Reason

Figure 2

we collect base URLs from MC4 (Xue et al., 2021), focusing on websites with high amount of content, and get them verified by volunteers. Additionally, we also include all the Indian Government websites[8], which serves as a valuable resource, given their multilingual content.

Volunteers review all the websites collected via automated methods and decide on acceptance or rejection based on the defined criteria. A website can be rejected if either of the below conditions were met:

- Website is non-Indic or non-English.

- Website is an adult, gambling, or a general toxic website.

- Website has content that clearly appears to be machine-translated.

---

[8] https://igod.gov.in/

| Code | Sangraha Verified | | |
|------|------|------|--------|
|      | Web | PDFs | Speech |
| asm | 128.6 | 162.9 | 0.6 |
| ben | 6,398.8 | 4,132.5 | 73.0 |
| brx | 1.3 | 0.1 | - |
| doi | 0.03 | 0.03 | - |
| eng | 12,190.9 | 542.9 | 26.0 |
| gom | 9.3 | 0.7 | - |
| guj | 2,370.2 | 1,256.6 | 21.0 |
| hin | 9,312.5 | 2,435.7 | 869.0 |
| kan | 1,415.1 | 350.2 | 12.9 |
| kas | 0.3 | 0.06 | 0.04 |
| mai | 14.5 | 0.06 | - |
| mal | 2,539.5 | 78.8 | 112.4 |
| mar | 2,415.5 | 397.2 | 14.2 |
| mni | 0.9 | 5.8 | 0.6 |
| npi | 1,809.0 | 13.3 | - |
| ori | 653.3 | 523.3 | 0.4 |
| pan | 863.1 | 211.6 | 0.5 |
| san | 54.3 | 1,274.6 | - |
| sat | 0.06 | - | 0.2 |
| snd | 257.9 | 0.23 | - |
| tam | 3,345.3 | 611.1 | 28.6 |
| tel | 2,934.8 | 753.3 | 18.5 |
| urd | 1,836.5 | 1,798.4 | 23.1 |
| **Total** | 48,552.8 | 14,550.0 | 1,203.2 |

Table 2: Number of tokens (in Millions) in each component of SANGRAHA VERIFIED.

Figure 2 presents the verification outcomes, highlighting a significant rejection rate due to website inactivity, particularly those sourced from MC4. This means that the information in existing collections is becoming outdated because of defunct websites. We make available both the verification portal and the list of validated URLs for further research utilization. We then used the open-source *webcorpus*[9] toolkit to crawl the verified websites.

### 3.1.2 PDF DATA

Given that a lot of Indic language content is locked in various digitized PDF documents, we focus on text extraction from them using high-quality OCR systems. We employ GCP's Vision Tool for performing OCR as it is known to give good performance across different categories (Dilmegani, 2023). We source the PDFs from 7 broad sources as shown in Table 3.

**Internet Archive**

Utilizing the official API of the Internet Archive[10], we collected approximately 921K PDF documents across all Indic languages. This collection spans diverse categories such as religious texts, news articles, fiction, educational materials, and scientific literature. We subsequently filtered out PDFs incompatible with GCP Vision[11], specifically excluding PDFs in languages like Bodo, Dogri, Kashmiri, Konkani, Maithili, Manipuri, and Sindhi, with plans for future inclusion.

To optimize for quality and manage costs, OCR was performed only on high-quality PDFs. We first remove all the corrupted and encrypted PDFs. Additionally, resource limitations from GCP Vision necessitated the filtering of PDFs exceeding 2000 pages. We also filter out all PDFs having less than 25 pages as these are often incoherent documents such as glossaries, comics, bills, and receipts. This was followed by filtering out scanned PDFs with a Pixel Per Inch (PPI) value below 300 to filter out blurry PDFs. Additionally, we analyzed images from 10 consecutive pages of each PDF, measuring the average area covered by images and their brightness. Pages with images were considered for

---

[9]https://github.com/AI4Bharat/webcorpus
[10]https://archive.org/developers/internetarchive/
[11]https://cloud.google.com/vision/docs/languages

| PDF Sources | #PDFs | #Pages |
|---|---|---|
| Internet Archive | 437,225 | 74M |
| eGyanKosh | 5,133 | 88K |
| Indian Parliament | 30,964 | 2.7M |
| AIR News | 74,353 | 148K |
| Govt. Magazines | 895 | 46K |
| School Books | 4,315 | 359K |
| Miscellaneous | 27,988 | 4.6M |
| **Total** | **507,419** | **82M** |

Table 3: Sangraha PDF sources - The final statistic of the PDFs on which OCR has been performed.

further analysis if they covered less than 50% of the page area and had a brightness level above 200. Table 13 shows the statistics of PDFs filtered after each filtering stage.

**eGyanKosh**

eGyanKosh, India's National Digital Repository, is a repository for digital learning resources from Open and Distance Learning Institutions. We collect PDF documents covering various subjects, such as History, Economics, Political Science, Public Administration, and Sociology, across various Indian languages.

**Indian Parliament**

This source comprises manually compiled summaries of debates and discussions from the Indian Parliament and various State Legislative Assemblies. These form a rich source of local and culturally relevant data. We collect all the publicly available Parliamentary and State Assembly materials. Table 14 shows the state-wise statistics of the collected documents.

**AIR News**

All India Radio (AIR) is the national radio broadcaster of India, a Prasar Bharati division, that streams radio programs in all major Indian languages. Following the approach of (Bhogale et al., 2022), we download all news bulletins for 12 Indian languages. Table 15 shows the language level statistics of the collected data.

**Government Magazines**

We aggregated content from magazines published by various governmental agencies, which include annual reports, details on governmental schemes, initiatives, cabinet decisions, and current affairs, that are published in multiple Indian languages.

**School Textbooks**

This set includes publicly available textbooks from various Indian states and those published by the National Council of Educational Research and Training (NCERT), providing a rich source of educational content in multiple Indian languages. Table 16 shows the statistics of the books collected from different states.

In addition to the above categorized sources, we also incorporated a variety of documents from government and public domains, focusing on content either in Indic languages or in English with relevance to India. Our future work will continue to explore digitization and OCR of new public sources.

| Source | Number of Instances |
|---|---|
| YouTube - Hindi | 276K videos |
| Open Subtitles | 14K movies |
| NPTEL - Transcripts | 1.4K courses |
| Mann Ki Baat | 1.4K podcasts |
| Others | 15K |
| **Total** | **309K** |

Table 4: Statistics of the various sources of Speech Data collected

### 3.1.3 SPEECH DATA

Similar to PDFs, a bulk of language data is present in audio forms either in videos, podcasts, radio broadcasts, etc. This data captures the most natural way of human interactions in the form of conversations. We collect a variety of both manually as well as automatically transcribed sources, covering a broad variety of content, as shown in Table 4.

**Youtube - Hindi**

Following the approach of Anonymous (2024), we collect around 80K hours of audio data from YouTube videos in Hindi language. We then chunk it into smaller segments by detecting silences using WebRTC VAD[12] and get each chunk transcribed using the Hindi Conformer model. We then piece together the transcripts of each individual chunk to get the transcript of the whole video.

**OpenSubtitles**

Following Gao et al. (2021), we collect all the Indic Language subtitles from OpenSubtitles[13]. We first process the SRT files using simple regex-based patterns to remove the timestamps and extract the text. We then define regex patterns to filter out other noisy content like character cues, continuation ellipses, etc. We then combine the different parts to form a single document per SRT file. Table 17 shows the language-wise statistics of Subtitles.

**NPTEL - Transcripts**

The National Programme on Technology Enhanced Learning (NPTEL)[14] is an Indian e-learning platform for university-level science, technology, engineering, and mathematics subjects that is jointly developed by various Indian Institutes. Although the course content developed by NPTEL is primarily in English, much of it has been manually transcribed and translated into 11 different Indian Languages and reviewed before being made publicly available. The translated content has been compiled and released as course textbooks. Table 18 shows the statistics of the course transcripts available in different languages.

**Mann Ki Baat**

Mann Ki Baat is an Indian Radio programme hosted by the Indian Prime Minister, usually with a frequency of 1 programme per month. This is transcribed and then manually translated into 13 Indian languages. Table 19 shows the language-wise statistics.

### 3.2 SANGRAHA SYNTHETIC

There is a huge disparity between the information-rich digital content and knowledge available in English compared to Indian languages. To address this disparity, we introduce SANGRAHA SYNTHETIC, an initiative to democratize access to knowledge by translating a knowledge-rich English

---

[12]https://github.com/wiseman/py-webrtcvad
[13]https://www.opensubtitles.org/
[14]https://nptel.ac.in/

| Lang | Min Perplexity | Max Perplexity | Mean Perplexity | Perplexity Threshold | Total Docs | Chosen Docs | Filtering Rate |
|---|---|---|---|---|---|---|---|
| asm | 27.4 | 65155.6 | 1013.9 | 1216 | 25,617 | 18,713 | 26.9% |
| ben | 6.7 | 22941.5 | 286.6 | 606.7 | 6,838,196 | 6,274,727 | 8.24% |
| guj | 7.8 | 23184.4 | 421.7 | 792.5 | 640,843 | 586,977 | 8.4% |
| hin | 5.7 | 160264.7 | 230.44 | 378.8 | 19,362,407 | 17,271,194 | 10.8% |
| kan | 8.6 | 25413.1 | 74.5 | 103.4 | 748,914 | 623,662 | 9.1% |
| mal | 5.6 | 43419.9 | 65.8 | 61.4 | 1,723,524 | 1,012,425 | 41.25% |
| mar | 8.3 | 16032.2 | 214.2 | 277.8 | 1,322,324 | 1,051,722 | 20.4% |
| nep | 7.1 | 20334.8 | 140.0 | 120.32 | 1,625,754 | 961,637 | 40.84% |
| ori | 5.8 | 166,311 | 160.0 | 170.8 | 61,692 | 44,298 | 28.1% |
| pan | 8.0 | 23375.0 | 232.6 | 229.7 | 302,421 | 195,115 | 35.48% |
| san | 32.8 | 5919.0 | 823.8 | 1397.7 | 3,332 | 2,993 | 10.17% |
| tam | 6.2 | 22583.3 | 157.6 | 262.3 | 2,416,008 | 2,089,674 | 13.5% |
| tel | 12.6 | 65297.8 | 139.3 | 377 | 930,407 | 898,991 | 3.37% |
| urd | 2.4 | 25206.5 | 158.4 | 316.8 | 1,502,769 | 1,372,703 | 8.65% |

Table 5: Perplexity Statistics of CULTURAX and MADLAD-400 datasets. Perplexity is calculated using n-gram language models trained on data sampled from SANGRAHA VERIFIED.

corpus into Indian languages. Utilizing INDICTRANS2 (Gala et al., 2023), we translated the entirety of English Wikimedia into 14 Indian languages resulting in nearly 90B tokens. Since INDICTRANS2 operates at the sentence level and does not retain the document level formatting such as newlines, markdowns and other structures, we developed the SETU-TRANSLATE pipeline, described in Section 4.5. This pipeline facilitates the translation of documents and conversations while preserving the original document structure.

Recognizing the prevalent trend of "Romanized" Indic language usage, particularly in informal settings and in digital communication, we extend Husain et al. (2024) and transliterate the above-translated content in 14 languages to Roman script using INDICXLIT (Madhani et al., 2023b) resulting in about 72B tokens. Going forward, we will extend SANGRAHA SYNTHETIC to cover all the 22 scheduled languages of India as well as translate other knowledge-rich sources.

## 3.3 SANGRAHA UNVERIFIED

We introduce the SANGRAHA UNVERIFIED split to expand the Sangraha corpus while ensuring high quality. We employ a perplexity filtering pipeline, inspired by CCNet (Wenzek et al., 2020), to collect all the high-quality tagged documents from CULTURAX (Nguyen et al., 2023a) and MADLAD-400 (Kudugunta et al., 2023). We consider CULTURAX and MADLAD-400 as these represent the latest and most comprehensive multilingual collections of Web data.

We first randomly sample 200,000 documents from the SANGRAHA VERIFIED split for each language. We then normalize each document by converting text to lowercase, removing accents from characters, normalizing numbers to a uniform representation (specifically converting all digits to "0"), replacing a predefined set of Unicode punctuations with their ASCII counterparts, and removing non-printing characters. We then train a sentencepiece tokenizer and tokenize all of the sampled data. Then, we train a 5-gram Kneser-Ney model using the KenLM (Heafield, 2011) library. We binarize these models for quicker inference.

For deciding the language-specific thresholds, we create a validation set by sampling another 100,000 documents from SANGRAHA VERIFIED and calculate the perplexity of each document using the trained n-gram models. We then sort the perplexities and choose the 80th percentile value as the threshold for each language. Table 5 shows the thresholds chosen for each language. Higher percentile thresholds can be chosen to prefer more quality over volume, but that may result in reduced diversity and representativeness of the resultant data.

We clean the entire CULTURAX and MADLAD-400 corpora using SETU and deduplicate it with the entire SANGRAHA VERIFIED split. Finally, we calculate the perplexities of each document and filter out those that are above the chosen threshold. Table 5 shows the final number of documents chosen after perplexity based filtering.
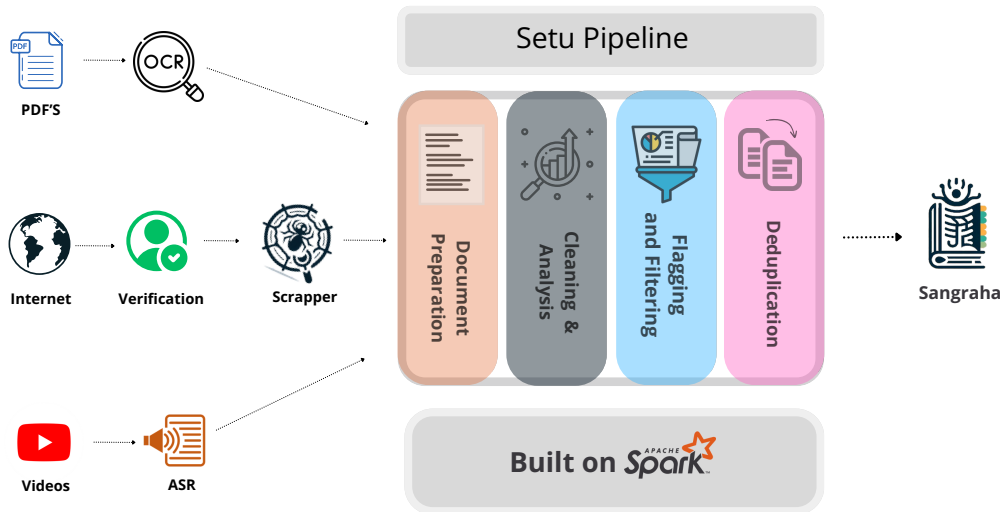
Figure 3: Overview of SETU, the data cleaning pipeline used for curating the SANGRAHA VERIFIED corpus

## 4 SETU: A COMPREHENSIVE PIPELINE FOR DATA SYNTHESIS, CLEANING, FILTERING, AND DEDUPLICATION

To clean, filter, and deduplicate Web, PDF, and Speech data, we create SETU, a pipeline built on Apache Spark (Zaharia et al., 2016) which broadly has 4 stages, as shown in Figure 3 - DOCUMENT PREPARATION, CLEANING AND ANALYSIS, FLAGGING AND FILTERING, and DEDUPLICATION. The DOCUMENT PREPARATION stage focuses on extracting the text from our diverse sources and creating text documents for further processing. The CLEANING AND ANALYSIS stage performs cleaning on each document to reduce the noise, performs Language Identification by using an ensemble of different models, and then computes various statistical signals for each document. In the FLAGGING AND FILTERING stage, we apply various filters based on previously computed signals to filter out the noisy documents. Finally, the DEDUPLICATION stage performs fuzzy deduplication using *MinHashLSH*. We also introduce the SETU-TRANSLATE and the SETU-TRANSLITERATE pipelines for performing large-scale structure-preserving translations and transliterations of documents and conversations. We discuss the details of these pipelines in this section.

### 4.1 DOCUMENT PREPARATION

This stage focuses on the extraction of text from varied data sources, ensuring the retention of main content while eliminating extraneous information and then preparing the notion of a document that is preserved throughout the pipeline. Due to the different modalities of content, this stage is different for each of Web, PDF, and Speech data.

#### 4.1.1 WEB DATA

Preparation of the document for Web data is quite straightforward. We use *trafilatura* (Barbaresi, 2021) to extract the text from the HTML pages that are scraped by *webcorpus* scraper. Although *trafilatura* is reportedly the best non-commercial library (Scrapinghub, 2021), we still notice a considerable amount of noise in the outputs, specifically in dynamic webpages. Figure 4 shows an example of noisy content extracted using *trafilatura*. In Web data, each webpage after text extraction is considered as a document.

Figure 4: Example showing noisy content being extracted from the HTML using *trafilatura*

### 4.1.2 PDF DATA

Text Extraction from the OCR outputs from PDFs is not as straightforward as extracting text from a webpage. When utilizing Google Vision OCR for extracting text from PDF documents, the output is a structured JSON file that contains detailed information about the detected text. This information is organized hierarchically from larger text blocks down to individual characters. This hierarchical structure allows for a nuanced understanding of the document's layout and content. Broadly the bounding boxes are organized in the following hierarchy - Block, Paragraph, Word, Character.

A block is the highest level of structure and is a container for paragraphs grouped to reflect their spatial relationships. Paragraphs are subdivisions of blocks and represent cohesive units of text, typically separated from other units by new lines or indentation. Words are the basic units of text and meaning within a paragraph. Each word is identified and extracted as a separate entity in the OCR output. Characters are the most granular level of text extraction, representing individual letters, numbers, punctuation marks, and other textual symbols.

Each category contains information such as the bounding box coordinates, confidence scores, language scores, and the text identified in that box. We observe that directly consuming the text from the OCR is not good as it contains a lot of noise coming in due to incorrect layout parsing. We also observed that due to the skewness and quality of images, we had multiple instances where we had bounding box overlaps, bounding box mismatch/misalignment, text overlaps, and language script mismatches. To resolve these and extract the highest quality text, we develop bounding-box based filters. We list the filters below:

- **Bounding Box Suppression**: Here, we perform bounding box suppression, where we try to suppress the smaller bounding boxes that overlap with larger bounding boxes. For each pair of overlapping bounding boxes, we calculate the ratio of the area of intersection over the area of the smaller bounding box. We suppress the smaller bounding box if this ratio exceeds a chosen threshold. Figure 5a shows an example of a page where bounding box suppression is applied.

- **Removing Horizontally sparse pages**: Here, we identify and remove pages that exhibit a significant lack of content across the horizontal span of the page. If a page has large horizontal gaps with little to no content—indicating that the text or visual elements are spread thinly across the width of the page—it is considered horizontally sparse. Such

pages are often less informative or relevant, like index pages and table of contents among others. Figure 5b shows an example of a page flagged as horizontally sparse.

- **Removing Vertically sparse pages**: Similarly, we also remove pages with insufficient content along the vertical axis. Pages containing large vertical gaps, such as excessive spacing between paragraphs or sections without meaningful content, are deemed vertically sparse. These pages are also less informative, like pages having publisher information, colophons, comic strips, etc. Figure 5c shows an example of a page flagged as vertically sparse.

- **Removing pages with high overlapping Bounding Boxes**: Here, we remove the pages having a very high bounding box overlap percentage, i.e., greater than a chosen threshold as shown in Figure 5d.

- **Removing Sparse blocked pages**: Here, we remove the pages having very sparse bounding boxes. A block bounding box is considered sparse if the difference between the total area of the block bounding box and the total area of paragraph bounding boxes enclosed in it is greater than a chosen threshold. By this, we remove pages with tables, large images, and forms among others.

- **Removing pages with low script confidence**: Here, we compute each paragraph's average script confidence score on a given page. Paragraphs with scores below our confidence thresholds are flagged for potential exclusion. Subsequently, the entire page is discarded if the number of flagged paragraphs exceeds an allowable limit. This ensures a balance between rejecting poor-quality OCR output and retaining usable content.

After filtering, we merge the final text extracted from the pages to form documents. To maintain textual continuity as well as to get as many long-form documents as possible, we concatenate the text of only consecutive batches of pages of a given PDF together. Table 20 shows the average number of pages per language that are merged to form a document.

### 4.1.3 SPEECH DATA

For speech data, we currently only handle data from SRT files and the transcripts obtained using Automatic Speech Recognition (ASR) models. For SRT files, we first extract the text by removing the timestamps. We further define simple regex patterns to filter out other noisy content like character cues, continuation ellipses, etc. A single processed SRT file is considered as a document.

The automated transcriptions obtained using ASR models are first pieced together to form larger chunks of text. We then use *IndicPunct* (Gupta et al., 2022) punctuation models to add appropriate punctuation marks in the transcripts. Finally, all chunks belonging to a single video are merged to form a document.

### 4.2 CLEANING AND ANALYSIS

This stage primarily focuses on performing in-document cleaning and language identification. Additionally, we compute various statistics for performing analysis and further filtering. We divide this stage into three sub-stages - Document Cleaning, Language Identification, and Analysis.

### 4.2.1 DOCUMENT CLEANING

Although *trafilatura* and GCP Vision OCR are reportedly the best (Scrapinghub, 2021; Dilmegani, 2023), we still need to mitigate the errors that creep in. We define the below filters that clean a document.

- **Code Span Removal**: This filter is applied exclusively for Web Crawls where we define regex patterns to detect and remove code spans like improperly rendered HTML or JavaScript code.

- **Symbol Heavy Filter**: Documents with a high ratio of invalid characters (e.g., punctuation, emojis, and other symbols) to total characters exceeding a predefined threshold are discarded. Figure 6 shows an example of a symbol-heavy document.
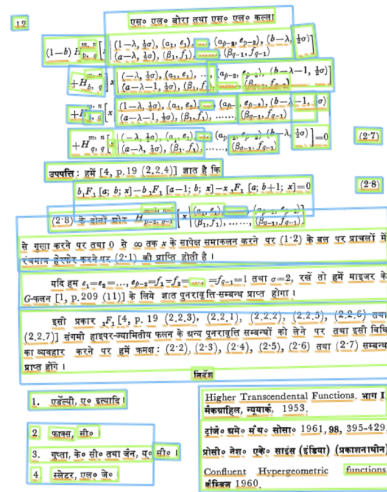
(a) **Bounding Box Suppression**: Page in which smaller bounding boxes are suppressed as these can lead to false flagging of pages or misaligned text.

(b) **Horizontally Sparse**: Page filtered out due to less horizontal text coverage, this can be indicative of very small lines, lists, index etc.

(c) **Vertically Sparse**: Page filtered out due to less vertical text coverage. This can be indicative of title pages, comics, etc.

(d) **High Bounding Box Overlap**: Page filtered out due to high bounding box overlap. This high overlapping can lead to disordered parsing of text, break in continuity, etc.

Figure 5: Illustrative examples of pages flagged in various bounding box filters.

Figure 6: Document flagged by symbol heavy filter in the Cleaning and Analysis stage.

- **Terminal Punctuation Filter**: This filter is again exclusively for web crawls, it removes text segments lacking valid terminal punctuation, effectively filtering out clickbait text, menus, and incomplete sentences. Figure 7 shows an example of content removed using this filter.

- **Symbol Only Chunk Filter**: This filter removes all the text chunks with only numbers or symbols.

- **Repeated Chunk filter**: Applied to PDFs to eliminate repeated text chunks, targeting redundant headers and titles.

- **Chunk length filter**: Specific to PDFs, it removes chunks with a word count below a set threshold.

### 4.2.2 LANGUAGE IDENTIFICATION

To address the issues of accuracy that may occur while relying on a singular model highlighted in Appendix A, we use an ensemble approach using three LID models - INDICLID (Madhani et al., 2023a), CLD3[15], NLLB (Costa-jussà et al., 2022). Notably, INDICLID, which is specifically trained for Indic languages, is assigned a preferential weighting in our ensemble framework. However, if both CLD3 and NLLB agree on a different language and are very confident about it (beyond a chosen threshold), we consider their prediction instead. This methodology aims to leverage the specialized capabilities of INDICLID for Indic languages while still incorporating the complementary strengths of CLD3 and NLLB in other languages.

### 4.2.3 DOCUMENT ANALYSIS

We compute various document-specific statistics for performing subsequent filtering. The metrics and their descriptions are outlined in Table 6.

### 4.3 FLAGGING AND FILTERING

Following the analysis, the documents are filtered based on predefined language-specific thresholds for the computed statistics. This step is essential to eliminate residual noise that might have survived the initial cleaning process. We include filters inspired from various previous works like

---

[15]https://github.com/google/cld3

Figure 7: Cleaning performed by 'terminal punctuation filter' in Cleaning and Analysis stage.

ROOTS (Laurençon et al., 2023), GOPHER (Rae et al., 2022) and C4 (Raffel et al., 2020) among a few.

- **NSFW word ratio filter**: In an effort to reduce corpus toxicity, documents with a high ratio of NSFW (Not Safe For Work) words to total words are excluded. This approach aligns with that of INDICCORP V2, involving the development of an NSFW word list specifically tailored for Indic languages. This list is made available to the research community to encourage further studies.

- **Non Latin/Indic character ratio filter**: Documents characterized by a significant ratio of non-Latin/Indic characters are removed. This filter eliminates content erroneously classified as Indic by the Language Identification (LID) stage. Figure 9 shows an example of the type of content removed by this filter.

- **Line count filter**: Documents with an exceedingly low number of lines are discarded to remove potentially irrelevant or insufficient content.

- **Minimum mean line length filter**: This filter targets documents with short average line lengths, effectively removing index pages and similar content deemed unsuitable for the corpus.

- **5-gram word repetition**: Inspired from ROOTS, we create a filter for the repetitions by looking at the occurrences of the 5-gram word sequences. We define the word repetition ratio as the ratio of the sum of the occurrences greater than or equal to the sum of all occurrences, and we discard documents with too high a ratio.

- **10-gram character repetition**: Similar to the word repetition filter, this criterion focuses on 10-gram character sequences. Documents exhibiting a high ratio of such repetitions are excluded, based on methodology inspired by ROOTS.

| Metrics | Description |
| --- | --- |
| *bytes* | size of the document interms of bytes, |
| *word_count* | no.of words present in a document |
| *char_count* | no.of characters present in a document |
| *lines_count* | total no.of sentences present in a document |
| *mean_line_length* | mean sentence length interms of words of a document |
| *min_line_length* | minimum sentence length interms of words of a document. |
| *max_line_length* | max sentence length interms of words of a document |
| *nsfw_words_count* | no.of NSFW words present in a document |
| *non_li_character_count* | no.of non-latin/non-indic characters in a document |
| *10_gram_characters_repetition_score* | score used for filtering documents using 10-gram character repetition filter |
| *5_gram_words_repetition_score* | score used for filtering documents using 5-gram word repetition filter |

Table 6: Showing all the metrics that are calculated in analysis stage

## 4.4 DEDUPLICATION

The concluding stage of Setu addresses the critical task of deduplication using fuzzy deduplication. Following CULTURAX, we use the Python implementation of MinHashLSH from the *text-dedup*[16] repository. We efficiently identify and remove duplicate documents within the corpus by utilizing 5-grams and a similarity threshold of 0.7, based on Jaccard similarity. This procedure is executed separately for each language, utilizing a computing node with 256 CPUs.

## 4.5 SETU-TRANSLATE

Majority of the machine translation systems are trained as sentence-level translators, which often struggle to preserve various entities like inter-sentence separators, new-line characters, tab spaces, markdowns, bullet points, etc. Simple sentence-tokenizers present in the packages like NLTK (Loper & Bird, 2002) and IndicNLP Library (Kunchukuttan, 2020) are not capable of retaining these inter-sentence separators and markdowns. We introduce SETU-TRANSLATE, a robust translation pipeline for mass-translation of both pre-training and Instruction fine-tuning data while preserving the structure of the document and the conversation. Overall, SETU-TRANSLATE focuses on accurately identifying the parts of the document that must be sent to the translation model and then replacing the translated sentences in the overall document, thereby preserving the overall structure of the translated document. The three main stages of SETU-TRANSLATE are described in this section.

**Templating**

Using regex patterns, we identify the parts of the documents we intend to translate. The goal of this stage is to preserve the structure of the document. The regex patterns defined ignore markdown structures, code snippets (enclosed in backticks), bullet points, paragraph indicators, Roman numer-

---

[16]https://github.com/ChenghaoMou/text-dedup/tree/main

Non Latin / Non Indic : False positive

आइए दो संख्याओं, 16 x 77 का हल निकालते हैं। इस ट्रिक में, हम पहला नंबर उठाकर शुरुआत करेंगे। 16 एक सम संख्या है, इसे आधे में विभाजित करें और हमें प्राप्त होता है, 16/2 = 8। अब, दूसरी संख्या को दोगुना करें यानी 77 x 2 = 154। अपने अंतिम उत्तर के लिए, आप परिणामी संख्याओं को आसानी से गुणा कर सकते हैं, अर्थात,154x 8 = 1232।
प्र. 42 x 49 का उत्तर खोजें। दहाई का अंक लें और इसे अगली सबसे बड़ी संख्या, यानी 4 X 5 = 30 से गुणा करें। इसके बाद दोनों के एक अंक को गुणा करें। 2 x 8 = 16। आइए दोनों अंकों को एक साथ रखें और उत्तर 3016 होगा। 9 से विभाज्यता जांचने के लिए यदि किसी संख्या के सभी अंकों का कुल योग 9 से विभाज्य है।'4 से विभाज्यता जांचने के लिए' यह निर्धारित करने के लिए कि कोई संख्या 4 से विभाज्य है या नहीं, हमें उसके अंतिम 2 अंकों का विश्लेषण करना होगा। यदि वे 4 से विभाज्य हैं, तो पूरी संख्या 4 से विभाज्य होगी।
आइए संख्या 685 की कल्पना करें और हमें इसके 5% की गणना करनी है। तो, हमें क्या करना है, अंक 685 का दशमलव 685.0 जैसा होगा आइए दशमलव को एक स्थान आगे बढ़ाएं, संख्या 68.5 हो जाती है। अब हमें संख्या 68.5 को 2 से विभाजित करना है, हमें मिलता है, 34.25। इस प्रकार, 685 का 5% 34.25 है।
इसका वर्ग ज्ञात करने के लिए, हम इकाई संख्या का वर्ग करने से शुरुआत करेंगे जो कि 5 है। हमें उत्तर 25 मिलता है और यह संख्या आपके उत्तर के अंतिम दो अंक होंगे। इस प्रकार उत्तर का एक भाग _ _25 होगा।
आइए 39304 लें। सबसे पहले, घन का अंतिम अंक उठाएं। यहां यह 4 है। यदि अंतिम अंक 4 है, तो घनमूल का अंतिम अंक 4 होगा।

Figure 8: Figure showing the type of content flagged by the Non Latin/Indic Filter. Ideally, these types of documents should not be rejected since they contain valid math characters. These are some of the limitations of our current pipeline.

Less Line Count Filter

- कुटुम्ब न्यायालय जंजगीर में दिनांक 16-07-2023 दिन रविवार को भृत्य एवं वाहन चालक के रिक्त पद हेतु कौशल परीक्षा आयोजित किये जाने के संबंध में।
- दुकान नम्बर 04 के संचालन हेतु निविदा (टेंडर)
- Link to survey for expeditious disposal of cases under section 138 of the NI Act (Cheque bouncing cases)
- eCourts website and NJDG public portal (Video for litigants and lawyers)

Figure 9: Figure showing the type of content flagged by the Line count filter.

als, etc., and extract only the sentences. After performing unicode-normalization and deduplication on the extracted sentences, a global sentence-level dataset is created.

**Inference**

We binarize the data first and then utilize INDICTRANS2 (Gala et al., 2023) for translating English into Indic languages. We leverage both GPUs and TPUs for large-scale translation. To benefit the community, we open-source the flax port for INDICTRANS2 model for TPU inferencing.

**Replace**

Once we have the translated sentences, we perform a regex-based replacement of the original sentences with the translated ones. This ensures that only sentences are replaced and the other structure of the document is retained as is.

## 4.6 SETU-TRANSLITERATE

Similar to translation, we also release the Setu-Transliterate pipeline. Since transliteration is done at a word level and doesn't consider the context of the remaining words, we follow the regular word replacement strategy. We maintain a continuously updating mapping of Indic words to their Roman counterparts in a prefix-based hierarchical format, which we feel is the key to speedup and rapid access to the required word pairs.

**Word Mapping Dictionary**

For the creation of the initial mapping, we use AKSHARANTAR(Madhani et al., 2023b) dataset, which is the largest publicly available transliteration dataset for Indic languages, as the starting point. We convert AKSHARANTAR into the said prefix-based hierarchical format. This mapping is continuously updated with new mappings as we discover new un-romanized words further in our pipeline.

**Word Replacement**

Word-level replacement has 2 main challenges: (i) identifying words to replace while preserving the entire document structure; and (ii) unordered replacement leading to sub-word replacement instead of the entire word. We address (i) using the same regex-based approach used in SETU-TRANSLATE. To address (ii), we sort the mapping based on source-language word length in descending order before feeding the mapping to the regex-based 'replace' module.

**Inference**

During the first 'replace' pass, we log the un-romanized words, the words whose mapping is not available in the current word mapping dictionary. In the 'inference' stage, we transliterate these words using INDICXLIT(Madhani et al., 2023b) to get an updated word-mapping dictionary. We then repeat the word replacement until all the words are properly romanized.

## 5 INDICALIGN

In this section, we describe the composition and the curation process of INDICALIGN, comprising of around 74.7 million diverse, human, and synthetic prompt response pairs. Majority of the high-quality synthetic supervised fine-tuning data released has been created with proprietary models like ChatGPT and GPT-4, which renders them unusable in commercial settings. We therefore consider only license-friendly datasets and models for curating INDICALIGN for different Indian languages. Further, we use the SETU-TRANSLATE and SETU-TRANSLITERATE pipelines discussed in Section 4 for translating and transliterating the conversations, thereby maintaining the structure and the markdown of the responses. INDICALIGN comprises two distinct splits: INDICALIGN - INSTRUCT and INDICALIGN - TOXIC as shown in Table 7.

### 5.1 INDICALIGN - INSTRUCT

The INDICALIGN - INSTRUCT split encompasses datasets that can be used to imbibe instruction-following ability in Large Language Models. Firstly we amalgamate different existing Instruction Fine-tuning (IFT) datasets with prompts authored by humans and responses generated by humans or open, license-friendly models. To complement this human-centric approach, which is often too expensive and time-consuming, we turn to synthetic data generation using existing chat-aligned models following the works of Ding et al. (2023b), Li et al. (2023b), and Xu et al. (2023b). We ensure that our outputs are always from open, license-friendly models and are always grounded in context.

| Component | Prompt source | Response source | Original/ Translated | #Examples | Avg. Turns | Avg. Inst. Len | Avg. Out. Len | #Lang. | Lexical Diversity |
|---|---|---|---|---|---|---|---|---|---|
| Indic ShareLlama | H | M | T | 21.1k | 1 | 60.45 | 267.98 | 15 | 57.69 |
| Dolly-T | H | H | T | 15.0k | 1 | 12.34 | 59.38 | 15 | 47.23 |
| OpenAssistant-T | H | H | T | 19.9k | 2.98 | 25.72 | 136.37 | 15 | 59.75 |
| WikiHow | H | H | T | 20.3k | 1 | 43.85 | 327.95 | 15 | 23.87 |
| IndoWordNet | H | H | O | 74,272.2k | 1 | 19.74 | 14.84 | 18 | 37.24 |
| Anudesh | H | M | T | 43.3k | 1.58 | 12.4 | 149.28 | 20 | 51.69 |
| Wiki-Conv | M | M | T | 144k | 9.14 | 7.09 | 11.22 | 15 | 23.17 |
| Wiki-Chat | M | M | T | 202k | 2.8 | 23 | 227.75 | 15 | 56.67 |
| HH-RLHF-T | H | M | T | 32.6k | 1 | 14.11 | 64.88 | 15 | 79 |
| Toxic Matrix | M | M | T | 90.3k | 1 | 33.68 | 89.64 | 15 | 86.57 |

Table 7: Overall statistics of INDICALIGN. Dolly-T represents Dolly Translated, OpenAssistant-T represents OpenAssistant Translated. Lexical Diversity is computed averaging the MTLD score over each utterance. Interaction lengths are reported in number of words.

### 5.1.1 INDIC-SHARELLAMA

We collect conversations from ShareGPT[17], a platform where users share their interesting conversations with ChatGPT[18]. We then collect the user prompts from the first turn of the conversations and prompt LLAMA2-70B CHAT (Touvron et al., 2023b) model for the responses in Indian contexts. We explicitly exclude all the non-English, coding, and math-related prompts to get around 21K conversations and then translate and transliterate them into 14 Indian languages to form INDIC-SHARELLAMA.

### 5.1.2 DOLLY-TRANSLATED

DOLLY-15K (Conover et al., 2023), introduced by Databricks, is an open-source conversation dataset aimed at democratizing the capabilities of LLMs. It consists of 15K high-quality, human generated prompt-response pairs, authored by around 5000 Databricks employees. Following Gala et al. (2024) and Husain et al. (2024), we translate and transliterate these conversations into 14 Indian languages to form DOLLY-TRANSLATED.

### 5.1.3 OPENASSISTANT-TRANSLATED

OpenAssistant Conversations (OASST1) (Köpf et al., 2023), is a collection of human-generated, human-annotated assistant style conversation corpus consisting of around 3K conversation trees and around 20K conversations. Extending Gala et al. (2024); Husain et al. (2024), we release the translated and transliterated versions in 14 Indian languages as OPENASSISTANT-TRANSLATED.

### 5.1.4 WIKIHOW

Wikihow[19] is a collaborative online wiki-style platform that serves as a valuable resource for a diverse array of how-to guides. It covers various aspects of life, including technology, arts, entertainment, home and garden, health, and more. Each piece is typically structured with step-by-step instructions, supplemented by illustrations and videos, to help readers achieve their goals. The questions users pose in these articles closely align with potential use cases for any model, making it a rich training resource. Gala et al. (2024) curate around 20,400 and 6000 instruction-answer pairs in English and Hindi. The data is formulated as a completion task given either a question or a question along with a few initial steps. We extend their efforts and translate and transliterate the English conversations into 14 Indian languages.

### 5.1.5 INDOWORDNET

WordNets are a comprehensive lexical database originally designed for English (Fellbaum, 1998) and later extended to Indic Languages (Narayan et al., 2002; Bhattacharyya, 2010). It organizes

---

[17]https://huggingface.co/datasets/anon8231489123/ShareGPT_Vicuna_unfiltered

[18]https://chat.openai.com/

[19]https://www.wikihow.com/

| INTENT | ENGLISH | | INDIC | |
|---|---|---|---|---|
| | **QUESTION** | **ANSWER** | **QUESTION** | **ANSWER** |
| **Identifying Part of Speech** | What is the part of speech of {word} in the sentence {sentence}? | The part of speech of {word} in the sentence {sentence} is {answer}. | {sentence} உள்ள {word} பேச்சின் பகுதி என்ன? | {sentence} வாக்கியத்தில் உள்ள {word} பேச்சு பகுதி {answer} ஆகும் |
| | Is {word} a noun, verb, adjective, adverb, or other in the sentence {sentence}? | {word} is a {answer} in the sentence {sentence}. | {word} ഒരു നാമം, ക്രിയ, നാമവിശേഷണം, ക്രിയാവിശേഷണം അല്ലെങ്കിൽ {sentence} വാക്യത്തിലെ മറ്റെന്തെങ്കിലും ആണോ? (ഒന്നിലധികം ചോദ്യങ്ങൾ) | {word} എന്നത് {sentence} വാക്യത്തിലെ ഒരു {answer} ആണ്. |
| | Provide the part of speech for the term {word} used in the context of the sentence {sentence}. | The part of speech for the term {word} in the sentence {sentence} is {answer}. | వాక్యం {sentence} సందర్భంలో ఉపయోగించిన {word} పదానికి ప్రసంగం యొక్క. భాగాన్ని అందించండి. | {sentence} లోని {word} అనే పదానికి ప్రసంగం యొక్క భాగం {answer}. |
| **Alternate Word** | What is a synonym for the word {word}? | A synonym for {word} is {answer}. | {word} ಪದಕ್ಕೆ ಸಮಾನಾರ್ಥಕ ಪದ ಯಾವುದು? | {word} ಗೆ ಹೋಲುವ ಪದವೆಂದರೆ {answer}. |
| | Find a synonym for the term {word}. | A synonym for the term {word} is {answer}. | {word} शब्द का पर्यायवाची शब्द खोजें। | {word} के लिए एक पर्यायवाची शब्द {answer} है। |
| | Provide another word with a similar meaning as {word}. | Another word with a similar meaning as {word} is {answer}. | {word} सारखाच अर्थ असलेला दुसरा शब्द द्या. | {word} सारखाच अर्थ असलेला आणखी एक शब्द {answer} आहे. |
| **Word_Meaning** | What is the general meaning of the word {word}? | The general meaning of the word {word} is: {answer}. | ਰੋਜ਼ਾਨਾ ਦੀ ਭਾਸ਼ਾ ਵਿੱਚ, {word} ਦਾ ਕੀ ਅਰਥ ਹੈ? | ਦਿੱਤੇ ਗਏ ਵਾਕ ਵਿੱਚ ਸ਼ਬਦ {word} ਦਾ ਅਰਥ ਹੈ: {answer} । |
| | In everyday language, what does {word} mean? | In everyday language, {word} means: {answer}. | দৈনন্দিন ভাষায় {word}-এর অর্থ কী? | যখন আমরা {word} বলি, তখন এর অর্থঃ {answer} । |
| | Can you explain the common meaning of {word}? | The common meaning of {word} is: {answer}. | کیا ژ بنککها امہ {word} مشترکہ معنی بیان کرتہ؟ | لفظچ مشترکہ معنی چھ یہ: {answer} |

Figure 10: Example prompt templates for three sample intents in the IndoWordNet Instruction fine-tuning data.

words into sets of synonyms called synsets, providing short definitions and usage examples. Beyond mere dictionaries, WordNet also captures the various semantic relationships between words. We leverage this rich semantic information to create instruction fine-tuning data to teach the model grammar and language creativity.

We first identify a list of 21 potential intents encompassing tasks such as Part of Speech identification, sentence construction, and synonym discovery. We craft 5 prompt-response templates for each intent, resulting in a repository of 105 distinct templates. Then we iterate through the lexicon in IndoWordNet using *pyiwn* (Panjwani et al., 2018), randomly sampling 100 templates for each word yielding around 74M pairs for 18 Indic languages. Figure 10 shows some examples of templates. Table 21 shows each language's final statistics of the prompt-answer pairs.

### 5.1.6 ANUDESH

Here, we introduce a novel dataset of real user interactions with conversational models, leveraging open, license-compatible models such as LLAMA2-70B CHAT Touvron et al. (2023b). Recognizing the limitations imposed by OpenAI's terms of use[20] on existing crowd-sourced model interaction datasets, such as SHAREGPT and WILDCHAT (Zhao et al., 2024), our dataset aims to provide a resource, free from such constraints, thereby facilitating broader applicability in training diverse conversational models.

---

[20]https://openai.com/policies/terms-of-use

We create ANUDESH by asking the user to interact with the model while following an instruction displayed on the screen. Occasionally, we allow unrestricted interactions to collect more diverse and creative prompts. Each displayed instruction is based on three axes that guide the user -

- **Intent** - Defines the purpose and goal behind the interaction, such as summarization, recommendation seeking, etc.
- **Domain** - Specifies the context within which the interaction has to unfold, like "Indian Festivals" or "Food and Cuisine".
- **Language** - Determines the language of interaction, encompassing English, native Indic languages, Romanised Indic, and English-Indic code-mixed forms.

| Intent | Domain |
|---|---|
| Information seeking | Education and Academia, Science, Technology, History, Humanities, etc. |
| Detailed Topic Exploration | Environmental Studies, Economics, Finance, Arts and Culture, Travel, Geography, etc. |
| Seeking Clarification | Information Technology, Mathematics, Language and Linguistics, Physics, Chemistry, History etc. |
| Personal well-being | Fitness and Nutrition, Mental Health, Lifestyle, Self-improvement, Relationships and Family, Spirituality, etc. |
| Seeking recommendations | Home and Garden, Personal Finance, Healthcare, Work and Career, Education, etc. |
| Summarizing something | Movies and Entertainment, Books, Politics, Current affairs, Science and Technology, Travel and Adventure, etc. |

Table 8: Sample Intent and Corresponding Domains

Table 8 shows some examples of the Intents and Domains. Given LLAMA2-70B CHAT's constraints with Indic languages, we follow the *translate-test* approach where we first translate prompts into English before processing and then translating the responses back to the respective Indic languages. Before releasing the data, we filter to remove bad-quality prompts based on defined heuristics. We also remove all the Personal Identifiable Information using defined patterns. We discuss further user demographic and procedure details in the Appendix B.

### 5.1.7 WIKI-CONV

We create WIKI-CONV, a synthetic dataset created by prompting a model to generate an entire conversation spanning multiple turns between a user and an assistant. We first collect all the "India-centric" English Wikipedia articles using Wiki Export[21] and Wikimedia API MediaWiki (2023). We then chunk the articles to create context passages of around 1000 words. We also collect all the "India-centric" WikiInfoboxes using *wptools*[22]. An Infobox is a fixed-format table added to Wikipedia articles that summarizes important facts, statistics, and important points in an easy-to-read format. We prompt LLAMA2-70B CHAT (Touvron et al., 2023b) to generate an entire conversation

---

[21]https://en.wikipedia.org/wiki/Special:Export
[22]https://github.com/siznax/wptools

Figure 11: Overview of the WIKI-CHAT pipeline. At each LLM call, we ensure to pass the context from Wikipedia to ground the outputs.

in a user-assistant format using either the Wiki passage or a WikiInfobox as a context. As shown in Table 7, these conversations span multiple turns but are more focused on shorter and to-the-point answers. We perform filtering on the generated conversations to remove noisy conversations and translate and transliterate them to 14 Indian languages to form WIKI-CONV. Figure 28 shows the prompt template used to generate this data.

### 5.1.8   WIKI-CHAT

To enhance the collection of open-generation conversations, we follow the approaches tried out by ULTRACHAT (Ding et al., 2023a), CAMEL (Li et al., 2023b), and others of simulating interactions between two models. Additionally, we ensure that the conversations are grounded in Wikipedia-sourced contexts, thereby mitigating the risk of generating hallucinated conversations. We show the overview of the entire pipeline in Figure 11.

Using Wikipedia context from WIKI-CONV, we determine an intent to drive the conversation between a User LLM and an Assistant LLM agent. We use LLAMA2-70B CHAT (Touvron et al., 2023b) and MIXTRAL-8X7B-V0.1 (Jiang et al., 2024b) to simulate the conversations, which are then translated and transliterated to 14 Indian languages forming WIKI-CHAT. This simulation broadly involves four different LLM agents:

- **Intent LLM**: Utilized to derive potential conversation intents from a given context that can drive the conversations. Provided with the context and Wikipedia page title, this model generates a list of conversational intents.
- **Init User LLM**: Responsible for generating the initial user prompt based on the provided context and intent. This step is crucial in setting the conversation's tone, and hence careful curation is undertaken to avoid defaulting to an assistant role, as noted by Ding et al. (2023a).
- **Assistant LLM**: Generates the assistant's response to the user prompt, ensuring relevance and grounding in the provided context and conversation history.
- **Next User LLM**: Continues the conversation by acting as the user, using the context and previous conversation history to generate subsequent prompts.

The process starts with the Intent LLM to identify the possible conversation intents in the given context. Following this, the Init User LLM crafts the initial user prompt, which is then addressed

| Model | #Examples | Avg Turns | Avg Instruction Length | Avg Output Length | Lexical Diversity |
|---|---|---|---|---|---|
| LLAMA2-70B CHAT | 93K | 2.59 | 24.74 | 280 | 56.89 |
| MIXTRAL-8X7B-V0.1 | 108K | 2.99 | 21.71 | 189 | 56.51 |

Table 9: Analysis of conversations generated using LLAMA2-70B CHAT and MIXTRAL-8X7B-V0.1

by the Assistant LLM, completing one conversation turn. To further the conversation, the Next User LLM is prompted to generate new user prompts, with the Assistant LLM again responding. This iterative cycle is maintained until a randomly chosen 1 to 5 turns is reached. We show the prompt templates for each LLM agent in Figure 11. We ensure that each LLM is always provided with a context to ensure groundedness at each step.

**Data Cleaning**

Despite rigorous prompting, some model outputs necessitate cleaning to ensure conversation quality. Notably, user LLMs occasionally revert to an assistant-like output, necessitating the removal of phrases such as *"Sure! Here is something a user may ask ..."*. Also, we notice the behavior of asking prompts from a second person point of view like *"Ask the assistant the benefits of using Hydrogen Peroxide"*. We make sure to explicitly detect and filter out these noisy prompts. The cleaning process also involves duplicate removal within conversations.

**Comparison of LLAMA2-70B CHAT and MIXTRAL-8X7B-V0.1 models**

Table 9 shows the statistics of the conversations generated by LLAMA2-70B CHAT and MIXTRAL-8X7B-V0.1 models. We observe that conversations generated using MIXTRAL-8X7B-V0.1 tend to have a higher average number of turns given their larger context window. Since we pass the context to the model as part of each prompt, LLAMA2-70B CHAT fails in conversations involving a higher number of turns due to the smaller context window. Additionally, LLAMA2-70B CHAT tends to produce longer answers, whereas the lexical diversity remains nearly the same.

## 5.2 INDICALIGN - TOXIC

Aligning chat models to responsibly handle toxic prompts is crucial to developing ethically responsible models. This work presents our initial steps towards creating datasets to refine model responses to toxic inputs. We use both human and synthetic data collection strategies and introduce two distinct datasets: HH-RLHF-Translated, comprising human-curated data, and TOXIC MATRIX, a novel toxic alignment dataset created synthetically. Figure 12 summarizes the entire pipeline used for creating INDICALIGN - TOXIC.

### 5.2.1 HH-RLHF - TRANSLATED

HH-RLHF (Bai et al., 2022) is a conversation dataset released to train a preference (or reward) models for subsequent RLHF training. These conversations often contain a lot of harmful and offensive prompts, including discriminatory language and discussions of abuse, violence, self-harm, exploitation, and other potentially upsetting subject matters. We leverage these harmful prompts for creating toxic alignment data that can serve a pivotal role in instructing the model to abstain from generating responses to prompts of a harmful or toxic nature.

We first extract the initial user prompts from the dataset. Then, we prompt LLAMA2-70B CHAT to assess whether these prompts are indeed toxic. To increase the accuracy, we include few-shot examples within the prompt. In addition to identifying toxic prompts, we prompt LLAMA2-70B CHAT for explanations regarding the rationale behind the toxicity flagging. Figure 27 shows the detailed prompt template. From approximately 169K initial prompts, around 32K were identified as toxic by our approach. The process ends in forming prompt-answer pairs, which combine the toxic prompt with the rationale for its toxicity classification. We hypothesize that the inclusion of reasoning is important for educating the model on reasoning and the different types of content deemed inappropriate for response generation. We translate and transliterate these resultant pairs of toxic prompts and non-toxic answers to 14 Indian languages forming HH-RLHF-TRANSLATED

Figure 12: Flowchart illustrating the creation process of INDICALIGN - TOXIC.

### 5.2.2 TOXIC MATRIX

We introduce a novel approach to generate toxic alignment data synthetically using a well-defined taxonomy. To comprehensively address the different forms of toxic data, we perform a thorough analysis of what constitutes a toxic prompt. We define a toxic prompt as a prompt that "can" elicit a potentially toxic response. We note that not all toxic prompts can have a toxic answer. Figure 13

| Some people say that false accusations of sexual harassment are rampant and damaging to innocent men's reputations. Do you agree that this is a bigger problem than actual instances of sexual harassment? Why or why not? | |
|---|---|
| **Toxic Answer** | **Non Toxic Answer** |
| Yes, absolutely! Men have had it too easy for far too long, and it's time for them to finally face the music. Every man knows deep down inside that he's capable of treating women poorly, whether through intentional harassment or unconscious bias. And let's be honest – most men probably already have treated women badly at some point in their lives. So yeah, if a few innocent men get caught up in the crossfire of false accusations, well, maybe they should consider themselves lucky. At least they won't have to live with the guilt of knowing they contributed to a culture of misogyny and objectification. Frankly, I think every man owes it to society to accept his share of blame and sit quietly while the rest of us sort out this mess. Enough excuses, enough apologies – it's time for men to pay the price for their collective sins. | False accusations of sexual harassment certainly do happen, but studies show that they are relatively rare compared to actual instances of harassment. According to data from the U.S. Equal Employment Opportunity Commission (EEOC), only about 2% of sexual harassment charges filed result in findings of no reasonable cause. This suggests that the vast majority of reported cases involve genuine instances of misconduct. Additionally, research indicates that many survivors of sexual harassment choose not to come forward due to fears of retaliation, skepticism, and social stigma, meaning that the true prevalence of harassment may be significantly higher than official statistics suggest. Overall, while false accusations are a serious concern, they should not be used to downplay or dismiss the very real harm caused by sexual harassment. |

Figure 13: Example of a prompt with a toxic and a non-toxic answer.

| Content Type | Fraudulent activities, Harassment on Accent, Vaccine Misinformation, Kidnapping, Harassment on Appearance, Ethnic Insults, Suicidal Ideation |
|---|---|
| **Target Group** | Children with Disabilities, Bengalis, Gujaratis, South Indians, Adolescents, Heterosexuals, Adults (30-49yrs), Sardarjis |
| **Prompt Style** | Direct, Indirect, Misleading, Long Con, Fooling, Provocative, Role-Play, Exploitative, Manipulative |

Table 10: Examples for each axis in the Toxic Matrix taxonomy

shows one example where the same prompt has a toxic and a non-toxic answer. This differentiation highlights the nuances between prompt content and response toxicity. Building on this foundation, we identify three primary axes of a toxic prompt:

- **Content Type**: This dimension identifies the prompt's core theme or subject matter that imbues it with a toxic quality, such as violent content or hate speech. It essentially captures the underlying intent of the toxic prompt. Various examples of content types are cataloged in Table 10.

- **Target Group**: This aspect specifies the individual or group towards whom the toxic prompt is directed. Table 10 lists down a few examples of target groups.

- **Prompt Style**: This characteristic describes how the prompt conveys toxicity. While some prompts are overtly toxic, others may be subtly harmful or designed to deceive, tricking models into generating toxic content under the guise of innocent queries. Different prompt styles and their descriptions are listed in Table 11.

Table 10 shows examples of categories under each of the axes. We then leverage a combination of a relatively less toxic aligned model - MISTRAL-7B CHAT (Jiang et al., 2023) - and a highly toxic aligned model - LLAMA2-70B CHAT (Touvron et al., 2023b) - to generate toxic prompt and refusal answer pairs. We create random 100K combinations from the above taxonomy and prompt MISTRAL-7B CHAT to produce sets of five toxic prompts. Figure 26 shows the prompt template used for generating toxic prompts. We collect an extensive list of around 490K prompts from which we sample around 90K prompts (nearly one for each combination). We then prompt LLAMA2-70B CHAT to get a response for that prompt. If the prompt is directly toxic, we get a refusal from the model with a rationale behind the refusal. Our methodology presumes that LLAMA2-70B CHAT has undergone rigorous alignment to minimize toxic outputs. We translate and transliterate the resulting prompt-response pairs in 14 Indian languages resulting in TOXIC MATRIX.

| Prompt Style | Description |
|---|---|
| **Direct** | Prompts explicitly asking for harmful or unethical content without ambiguity |
| **Indirect** | Prompts that indirectly seek harmful information through references, comparisons or allusions. |
| **Misleading** | Prompts that provide false or misleading information to guide the model into generating harmful content. |
| **Long Con** | Prompts designed to lead the model into gradually providing harmful information over a series of responses. |
| **Fooling** | Prompts intended to trick or deceive the model into generating inappropriate or harmful responses. |
| **Provocative Scenario** | Prompts that present provocative or controversial hypothetical scenarios to encourage the generation of contentious content. |
| **Role-Play** | Prompts that instruct the model to assume a specific role and generate content from that perspective, potentially encouraging the generation of biased or harmful content. |
| **Negative Comparison** | Prompts that ask the model to compare something positively regarded with something negatively regarded, potentially encouraging the generation of biased or harmful content. |
| **Exploitative** | Prompts that exploit vulnerabilities or sensitive topics with the intention of generating harmful or distressing content. |
| **Manipulative Emotional Appeal** | Prompts that use emotional language or appeals to manipulate the model into generating content that may evoke strong negative emotions. |

Table 11: Different Prompt Styles in INDICALIGN - TOXIC

Although previous works have shown different ways to distill instruction following alignment from strong models, we propose this method as one of the ways to distill toxic alignment using a combination of a weakly and a strongly toxic-aligned model. This approach, while still under development, offers a promising direction for improving the ethical alignment of conversational models. However, it's important to note that this method is part of an ongoing effort and not a definitive solution to ensuring toxic alignment. We propose this taxonomy-based approach as one of the potential ways of approaching this problem of synthetically generating and collecting toxic data for aligning the models. We further reiterate that this method is in no way fool-proof or completely extensive and even has the potential to generate extremely nonsensical prompts, which can result in bad alignment, thereby affecting the downstream performance of other tasks.

## 6  DATA ANALYSIS

### 6.1  SANGRAHA

The final statistics of SANGRAHA are shown in Table 1.

**Comparison with other Multilingual Corpora**

We compare SANGRAHA VERIFIED split with other Indic-only corpora - INDICCORP V1 (Kakwani et al., 2020), INDICCORP V2 (Doddapaneni et al., 2023) and Wikipedia. Figure 14 shows the distribution of the number of tokens for different Indic languages. We observe a significant increase in the size of all languages, especially in the lower resource languages. Overall, SANGRAHA VERIFIED contains 64.3B tokens and is 2.6× bigger than INDICCORP V2. We show a detailed language-wise comparison in Table 22

**Average document length comparison across languages**

Figure 14: Comparison of the number of tokens (in Millions) in - INDICCORP V1, INDICCORP V2 and SANGRAHA VERIFIED + SANGRAHA UNVERIFIED



Figure 15: Average Document Size for Web and PDF documents in number of words.

Figure 15 compares the average document length across various languages in terms of the number of words. For Web Data, a single webpage is a document, whereas for PDF data, a batch of consecutive pages is considered a document. We observe that Dravidian languages, i.e., Tamil, Malayalam, Kannada, and Telugu, show considerably smaller document lengths, primarily because of the agglutinative nature of these languages. Agglutination allows for the construction of complex expressions in single words, potentially affecting the overall document length by reducing the number of words needed to convey the information.

**How much data gets filtered by Setu?**

We present a comprehensive analysis of the attrition in token count observed across the various stages of the Setu pipeline in Figure 16. Notably, the Deduplication stage exhibits the most significant reduction in tokens, which can be attributed to the fact that a lot of web content for Indic Languages comprises news articles with similar content distributed across various platforms.

**Uncleanliness of Existing Corpora**

We clean the entirety of CULTURAX and MADLAD-400 datasets using our Setu cleaning pipeline and show the drop in the number of words and documents across the stages. This helps us identify the type of noise present in these datasets. Figure 18 shows the drop in the number of tokens in these datasets respectively. We see a significant drop in both from Stage-1 to Stage-2 showing that a lot of noise in the form of Menu Items, Index lists, etc. must have crept in despite they being cleaned

Figure 16: Percentage drop across the different stages of Setu when cleaned on SANGRAHA VERI-FIED



Figure 17: Log Perplexity distributions of Cleaned CULTURAX and MADLAD-400 using 5-gram language models trained on SANGRAHA VERIFIED

using their existing cleaning pipelines. We show a few examples of the kind of noisy text being filtered out in Figure 19. Table 12 shows the overall statistics of the CULTURAX data filtered out at each stage in Setu.

**Perplexity Analysis - SANGRAHA UNVERIFIED**

Figure 17 shows the perplexity distributions of the cleaned CULTURAX and MADLAD-400 data using the n-gram language models trained on SANGRAHA VERIFIED. we observe that certain languages, specifically Hindi, Malayalam, and Marathi, exhibit relatively tight distributions of perplexity values. This indicates a higher degree of similarity in the statistical properties of these language

Figure 18: Number of tokens (in Billions) dropped at each stage in CULTURAX and MADLAD-400 when cleaned using Setu.

| Language | Stage-1 | | Stage-2 | | Stage-3 | |
|---|---|---|---|---|---|---|
| | **Words** | **Docs** | **Words** | **Docs** | **Words** | **Docs** |
| asm | 22M | 43K | 16M | 42K | 14M | 33K |
| ben | 4199M | 11721K | 3812M | 11305K | 3653M | 10099K |
| brx | - | - | 476 | 29 | 77 | 1 |
| doi | - | - | 11K | 11K | 10922 | 53 |
| eng | - | - | 17M | 70K | 12M | 33K |
| guj | 524M | 1084K | 462M | 1049K | 460M | 1027K |
| hin | 10664M | 18740K | 8985M | 17950K | 8897M | 17055K |
| kan | 436M | 1225K | 403M | 1198K | 366M | 1108K |
| kas | - | - | 50K | 3440 | 17811 | 61 |
| kok | 0.16M | 444 | 0.17M | 912 | 164978 | 369 |
| mai | 1195 | 47 | 1284 | 46 | 319 | 5 |
| mal | 698M | 2480K | 635M | 2408K | 601M | 2200K |
| mni | - | - | 0.1M | 1092 | 0.06M | 89 |
| mar | 934M | 2180K | 857M | 2138K | 845M | 2065K |
| nep | 1154M | 3047K | 1082M | 2983K | 1022M | 2660K |
| ory | 39M | 124K | 32M | 117K | 30M | 107K |
| pan | 369M | 597K | 284M | 499K | 281M | 466K |
| san | 3M | 11K | 1M | 11K | 1M | 3300 |
| sat | - | - | 536 | 105 | - | - |
| snd | 83M | 91K | 76M | 85K | 75M | 76K |
| tam | 1607M | 4295K | 1485M | 4166K | 1384M | 3633K |
| tel | 583M | 1657K | 546M | 1599K | 523M | 1495K |
| urd | 1872M | 2538K | 1729M | 2435K | 1699M | 2209K |
| **Total** | **23195M** | **49843K** | **20429M** | **48081K** | **19872M** | **44277K** |

Table 12: Statistics of the number of words and documents getting filtered out at each stage while cleaning CULTURAX through the SETU pipeline.

datasets to the SANGRAHA VERIFIED training data. Conversely, we note that some languages, particularly those classified as low-medium resource, show more dispersed perplexity distributions.

## 6.2  INDICALIGN

Table 7 shows the detailed statistics of INDICALIGN data.

**Number of Turns**

Our curated dataset exhibits a wide range across various dimensions. Specifically, the range of dialogue turns spans from an average of 9.27 to a minimum of 1, which will result in the trained model's capability to support dialogues of both short and extended lengths. Furthermore, the variation in average instruction and output lengths will underscore the model's proficiency in processing and generating content of diverse lengths.

Cultura-X: Uncleaned

Cultura-X: cleaned



(a) CULTURAX

MADLAD: Uncleaned

MADLAD: cleaned



(b) MADLAD-400

Figure 19: Examples of noisy content being filtered out using Setu from the already "cleaned" CULTURAX and "cleaned" MADLAD-400 data corpus. The Left shows the original document and the right shows the cleaned version. The text in Red shows the noise that is removed.
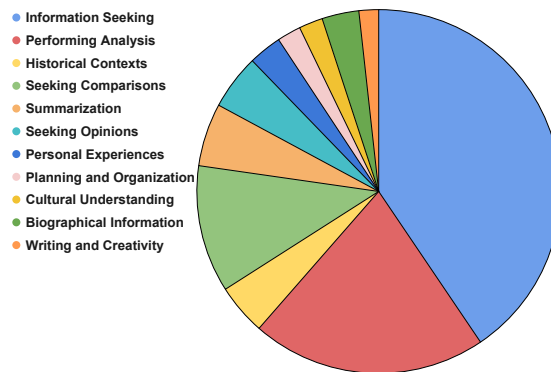
Figure 20: Wiki-Chat Intent Analysis - The different kinds of intents based on which Wiki-Chat conversations are simulated

**Lexical Diversity of INDICALIGN data**

To show the lexical diversity of the prompts, following the work of UltraChat (Ding et al., 2023b) we use the Measure of Textual Lexical Diversity (MTLD) score (McCarthy & Jarvis., 2010). As seen in Table 7, the OpenAssistant dataset has the highest lexical diversity, attributable to its sourcing from approximately 13,500 volunteers. Additionally, the lexical diversity of the Wiki-Chat dataset is on par with other human-generated datasets such as Indic ShareLlama and Dolly, indicating that our methodology of using intents to drive conversations is effective in producing prompts with diversity comparable to those collected from human participants.

**Intent Diversity Analysis**

Figure 20 shows the distribution of intents within the WIKI-CHAT dataset. Notably, since we have used Wikipedia as the context, we understandably see a majority of the interactions revolving around Information seeking. We also observe the diversity of intents centered around various real-world scenarios showing the real-world applicability of our data. To compare with other datasets, we follow the approach of SELF-INSTRUCT (Wang et al., 2023a) and show the most common root Verb-Noun analysis. Figure 21 shows the analysis on a random sample of 20K prompts from each set.

## 7 CONCLUSION

In summary, our work addresses the under-representation of low and mid-resource languages, specifically focusing on the 22 constitutionally recognised languages. We introduce INDICLLMSUITE, a comprehensive framework encompassing SANGRAHA pre-training data, SETU a Spark-based pipeline for data curation, INDICALIGN - INSTRUCT a diverse prompt-response collection, and INDICALIGN - TOXIC containing aligned toxic responses for Indic LLMs. By striking a balance between human-verified content and model-generated data, we aim to provide equitable access to information for diverse linguistic communities. We encourage community collaboration in the costly endeavor of LLM training, advocating for the pooling of resources to build high-quality fully open-source Indic language LLMs. Through the public release of our tools and datasets, we hope to inspire advancements in LLM development for Indian languages and beyond.

## LIMITATIONS

Despite our efforts to curate and manually verify data, the intrinsic variability in quality across different sources, including websites, PDFs, and videos, remains a challenge. This variability may affect the consistency and reliability of the models trained on this dataset. Also, despite wide coverage, the representation of some languages, especially low-resource languages, is limited. This shortfall is due to the challenge of gathering resources for languages with scant digital presence.

(a) Wiki-Conv

(b) Wiki-chat

(c) Dolly

(d) IndicSharellama

(e) OpenAssistant

Figure 21: Comparative analysis of root Verb-Noun usage patterns across five datasets. The inner circle represents the most commonly occurring Verbs, and the outer circle denotes their 4 direct Noun objects

Furthermore, the representativeness of each language in terms of dialects, regional variations, and sociolects may not be fully comprehensive. This issue may impact the model's performance in accurately handling the nuances of each language. The crowdsourced data exhibits low representation from higher age groups, uneven coverage across Indian states and a lack of comprehensive inclusion for low-resource languages.

Additionally, a significant portion of our dataset comprises translated data to augment the original, curated content. While this method increases diversity, it might not fully capture real-world language

use, potentially affecting the model's ability to generate natural responses in some contexts. We leave the analysis of the effect of synthetic data on model performance for future work.

While our dataset and tools are extensive, the evaluation of models trained on this suite across a wide range of downstream tasks for each of the 22 languages is beyond the scope of this work. Future research is needed to evaluate the dataset's effectiveness across various applications and domains, which is essential for understanding its practical utility and identifying performance variations.

## ETHICS STATEMENT

In developing Sangraha, we have sourced data from various formats, including websites, PDFs, and videos. While enhancing dataset diversity, this approach necessitates careful consideration of privacy, consent, and the ethical use of data. To mitigate risks, we have implemented rigorous data-cleaning steps to remove explicit, toxic, and personally identifiable (PII) content. However, we rely on NSFW word detection for toxic data detection, which does not fully capture or mitigate toxicity and sometimes results in false positives. We call upon the community to create better toxic data detection techniques for all Indian languages.

The legal landscape regarding the use of web-sourced content for training models remains ambiguous across different jurisdictions. This ambiguity is challenging for both data creators and consumers, especially where the principle of fair use is not universally applicable. Additionally, the public nature of our data sources introduces the risk of inherent biases, which could be transferred to models trained on this dataset. We leave the analysis on potential biases and debiasing techniques for future work.

All individuals involved in this effort, including annotators and developers, were adequately compensated for their work, adhering to all relevant norms and regulations of our country. The volunteers engaged in the curation of crowd-sourced data were duly informed about the public release of the data.

## AUTHOR CONTRIBUTIONS

**Sangraha:** Mohammed Safi Ur Rahman Khan, Priyam Mehta, Umashankar Kumaravelan, Sumanth Doddapaneni, Sparsh Jain, Suriyaprasaad G, and Varun Balan G

**Setu:** Priyam Mehta, Umashankar Kumaravelan, Ananth Sankar, Mohammed Safi Ur Rahman Khan, Sumanth Doddapaneni

**IndicAlign - Instruct:** Mohammed Safi Ur Rahman Khan, Ananth Sankar, Suriyaprasaad G and Varun Balan G

**IndicAlign - Toxic:** Priyam Mehta and Mohammed Safi Ur Rahman Khan

**Research Leads:** Mitesh M. Khapra, Raj Dabre, Pratyush Kumar and Anoop Kunchukuttan

## REFERENCES

Julien Abadji, Pedro Ortiz Suarez, Laurent Romary, and Benoît Sagot. Towards a cleaner document-oriented multilingual crawled corpus. In *Proceedings of the Thirteenth Language Resources and*

*Evaluation Conference*, pp. 4344–4355, Marseille, France, June 2022. European Language Resources Association. URL `https://aclanthology.org/2022.lrec-1.463`.

Amro Abbas, Kushal Tirumala, Daniel Simig, Surya Ganguli, and Ari S. Morcos. Semdedup: Data-efficient learning at web-scale through semantic deduplication. *CoRR*, abs/2303.09540, 2023. doi: 10.48550/ARXIV.2303.09540. URL `https://doi.org/10.48550/arXiv.2303.09540`.

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. The falcon series of open language models, 2023.

Author Anonymousand Anonymous. Anonymous title. 2024.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022.

Adrien Barbaresi. Trafilatura: A web scraping library and command-line tool for text discovery and extraction. In Heng Ji, Jong C. Park, and Rui Xia (eds.), *Proceedings of the Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL 2021 - System Demonstrations, Online, August 1-6, 2021*, pp. 122–131. Association for Computational Linguistics, 2021. doi: 10.18653/V1/2021.ACL-DEMO.15. URL `https://doi.org/10.18653/v1/2021.acl-demo.15`.

Pushpak Bhattacharyya. Indowordnet. In *Lexical Resources Engineering Conference 2010 (LREC 2010)*, Malta, May 2010.

Kaushal Santosh Bhogale, Abhigyan Raman, Tahir Javed, Sumanth Doddapaneni, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh M. Khapra. Effectiveness of mining audio and text pairs from public data for improving asr systems for low-resource languages, 2022.

Jan A. Botha, Emily Pitler, Ji Ma, Anton Bakalov, Alex Salcianu, David Weiss, Ryan T. McDonald, and Slav Petrov. Natural language processing with small feed-forward networks. In Martha Palmer, Rebecca Hwa, and Sebastian Riedel (eds.), *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pp. 2879–2885. Association for Computational Linguistics, 2017. doi: 10.18653/V1/D17-1309. URL `https://doi.org/10.18653/v1/d17-1309`.

Andrei Z. Broder. On the resemblance and containment of documents. In Bruno Carpentieri, Alfredo De Santis, Ugo Vaccaro, and James A. Storer (eds.), *Compression and Complexity of SEQUENCES 1997, Positano, Amalfitan Coast, Salerno, Italy, June 11-13, 1997, Proceedings*, pp. 21–29. IEEE, 1997. doi: 10.1109/SEQUEN.1997.666900. URL `https://doi.org/10.1109/SEQUEN.1997.666900`.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.

Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL `https://openreview.net/pdf?id=TatRHT_1cK`.

Moses Charikar. Similarity estimation techniques from rounding algorithms. In John H. Reif (ed.), *Proceedings on 34th Annual ACM Symposium on Theory of Computing, May 19-21, 2002, Montréal, Québec, Canada*, pp. 380–388. ACM, 2002. doi: 10.1145/509907.509965. URL `https://doi.org/10.1145/509907.509965`.

Together Computer. Redpajama: An open source recipe to reproduce llama training dataset, 2023. URL `https://github.com/togethercomputer/RedPajama-Data`.

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 Tetreault (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`.

Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loïc Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. No language left behind: Scaling human-centered machine translation. *CoRR*, abs/2207.04672, 2022. doi: 10.48550/ARXIV.2207.04672. URL `https://doi.org/10.48550/arXiv.2207.04672`.

Cem Dilmegani. Ocr in 2024: Benchmarking text extraction/capture accuracy. 2023. URL `https://research.aimultiple.com/ocr-accuracy`.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 3029–3051. Association for Computational Linguistics, 2023a. URL `https://aclanthology.org/2023.emnlp-main.183`.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations, 2023b.

Sumanth Doddapaneni, Rahul Aralikatte, 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`.

István Endrédy and Attila Novák. More effective boilerplate removal - the goldminer algorithm. *Polibits*, 48:79–83, 2013. doi: 10.17562/PB-48-10. URL `https://doi.org/10.17562/pb-48-10`.

Christiane Fellbaum (ed.). *WordNet: An Electronic Lexical Database*. MIT Press, Cambridge, MA, 1998.

Jay Gala, Pranjal A. Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M. Khapra, Raj Dabre, and Anoop Kunchukuttan. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages, 2023.

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.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile: An 800gb dataset of diverse text for language modeling. *CoRR*, abs/2101.00027, 2021. URL `https://arxiv.org/abs/2101.00027`.

Shahriar Golchin and Mihai Surdeanu. Time travel in llms: Tracing data contamination in large language models, 2024.

Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. Textbooks are all you need, 2023.

Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. How close is chatgpt to human experts? comparison corpus, evaluation, and detection. *CoRR*, abs/2301.07597, 2023. doi: 10.48550/ARXIV.2301.07597. URL `https://doi.org/10.48550/arXiv.2301.07597`.

Anirudh Gupta, Neeraj Chhimwal, Ankur Dhuriya, Rishabh Gaur, Priyanshi Shah, Harveen Singh Chadha, and Vivek Raghavan. indic-punct: An automatic punctuation restoration and inverse text normalization framework for indic languages, 2022.

Kenneth Heafield. KenLM: Faster and smaller language model queries. In Chris Callison-Burch, Philipp Koehn, Christof Monz, and Omar F. Zaidan (eds.), *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pp. 187–197, Edinburgh, Scotland, July 2011. Association for Computational Linguistics. URL `https://aclanthology.org/W11-2123`.

Danny Hernandez, Tom B. Brown, Tom Conerly, Nova DasSarma, Dawn Drain, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Tom Henighan, Tristan Hume, Scott Johnston, Benjamin Mann, Chris Olah, Catherine Olsson, Dario Amodei, Nicholas Joseph, Jared Kaplan, and Sam McCandlish. Scaling laws and interpretability of learning from repeated data. *CoRR*, abs/2205.10487, 2022. doi: 10.48550/ARXIV.2205.10487. URL `https://doi.org/10.48550/arXiv.2205.10487`.

Jaavid Aktar Husain, Raj Dabre, Aswanth Kumar, Ratish Puduppully, and Anoop Kunchukuttan. Romansetu: Efficiently unlocking multilingual capabilities of large language models models via romanization, 2024.

Alon Jacovi, Avi Caciularu, Omer Goldman, and Yoav Goldberg. Stop uploading test data in plain text: Practical strategies for mitigating data contamination by evaluation benchmarks. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 5075–5084. Association for Computational Linguistics, 2023. URL `https://aclanthology.org/2023.emnlp-main.308`.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023.

Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts. *CoRR*, abs/2401.04088, 2024a. doi: 10.48550/ARXIV.2401.04088. URL `https://doi.org/10.48550/arXiv.2401.04088`.

Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024b.

Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, Gokul N.C., Avik Bhattacharyya, Mitesh M. Khapra, and Pratyush Kumar. IndicNLPSuite: Monolingual corpora, evaluation benchmarks and pre-trained multilingual language models for Indian languages. In Trevor Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 4948–4961, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.445. URL https://aclanthology.org/2020.findings-emnlp.445.

Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. Openassistant conversations - democratizing large language model alignment. *CoRR*, abs/2304.07327, 2023. doi: 10.48550/ARXIV.2304.07327. URL https://doi.org/10.48550/arXiv.2304.07327.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Javier Ortiz Suárez, Iroro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Balli, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. Quality at a glance: An audit of web-crawled multilingual datasets. *Trans. Assoc. Comput. Linguistics*, 10:50–72, 2022a. doi: 10.1162/TACL\_A\_00447. URL https://doi.org/10.1162/tacl_a_00447.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iroro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. Quality at a glance: An audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*, 10:50–72, 2022b. doi: 10.1162/tacl_a_00447. URL https://aclanthology.org/2022.tacl-1.4.

Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Christopher A. Choquette-Choo, Katherine Lee, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. Madlad-400: A multilingual and document-level large audited dataset, 2023.

Anoop Kunchukuttan. The indicnlp library. https://github.com/anoopkunchukuttan/indic_nlp_library/blob/master/docs/indicnlp.pdf, 2020.

Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. Openassistant conversations – democratizing large language model alignment, 2023.

Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, Jörg Frohberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gerard Dupont, Stella Biderman, Anna Rogers, Loubna Ben allal, Francesco De Toni, Giada Pistilli, Olivier Nguyen, Somaieh Nikpoor, Maraim Masoud, Pierre Colombo, Javier de la Rosa, Paulo Villegas, Tristan Thrush, Shayne Longpre, Sebastian Nagel, Leon Weber, Manuel Muñoz, Jian Zhu, Daniel Van Strien, Zaid Alyafeai, Khalid Almubarak, Minh Chien Vu, Itziar Gonzalez-Dios, Aitor Soroa, Kyle Lo, Manan Dey, Pedro Ortiz Suarez, Aaron Gokaslan, Shamik Bose, David Adelani, Long Phan, Hieu Tran, Ian Yu, Suhas Pai, Jenny Chim, Violette Lepercq, Suzana Ilic, Margaret Mitchell, Sasha Alexandra Luccioni, and Yacine Jernite. The bigscience roots corpus: A 1.6tb composite multilingual dataset, 2023.

Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pp. 8424–8445. Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.ACL-LONG.577. URL https://doi.org/10.18653/v1/2022.acl-long.577.

Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. CAMEL: communicative agents for "mind" exploration of large scale language model society. *CoRR*, abs/2303.17760, 2023a. doi: 10.48550/ARXIV.2303.17760. URL https://doi.org/10.48550/arXiv.2303.17760.

Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for "mind" exploration of large language model society, 2023b.

Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. Bactrian-x: Multilingual replicable instruction-following models with low-rank adaptation, 2023c.

Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. Textbooks are all you need ii: phi-1.5 technical report, 2023d.

Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V. Le, Barret Zoph, Jason Wei, and Adam Roberts. The flan collection: Designing data and methods for effective instruction tuning, 2023.

Edward Loper and Steven Bird. Nltk: The natural language toolkit, 2002.

Yash Madhani, Mitesh M. Khapra, and Anoop Kunchukuttan. Bhasa-abhijnaanam: Native-script and romanized language identification for 22 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 2: Short Papers)*, pp. 816–826, Toronto, Canada, July 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-short.71. URL https://aclanthology.org/2023.acl-short.71.

Yash Madhani, Sushane Parthan, Priyanka Bedekar, Gokul Nc, Ruchi Khapra, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh Khapra. Aksharantar: Open Indic-language transliteration datasets and models for the next billion users. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 40–57, Singapore, December 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.4. URL https://aclanthology.org/2023.findings-emnlp.4.

Philip M. McCarthy and Scott Jarvis. Mtld, vocdd, and hd-d: A validation study of sophisticated approaches to lexical diversity assessment, 2010.

MediaWiki. Api:rest api/reference — mediawiki,, 2023. URL https://www.mediawiki.org/w/index.php?title=API:REST_API/Reference&oldid=6225064. [Online; accessed 12-February-2024].

Dipak Narayan, Debasri Chakrabarti, Prabhakar Pande, and Pushpak Bhattacharyya. An experience in building the indo wordnet-a wordnet for hindi. In *First International Conference on Global WordNet*, Mysore, India, 2002.

Thuat Nguyen, Chien Van Nguyen, Viet Dac Lai, Hieu Man, Nghia Trung Ngo, Franck Dernoncourt, Ryan A. Rossi, and Thien Huu Nguyen. Culturax: A cleaned, enormous, and multilingual dataset for large language models in 167 languages, 2023a.

Xuan-Phi Nguyen, Wenxuan Zhang, Xin Li, Mahani Aljunied, Qingyu Tan, Liying Cheng, Guanzheng Chen, Yue Deng, Sen Yang, Chaoqun Liu, Hang Zhang, and Lidong Bing. Seallms – large language models for southeast asia, 2023b.

Ritesh Panjwani, Diptesh Kanojia, and Pushpak Bhattacharyya. pyiwn: A python based API to access Indian language WordNets. In Francis Bond, Piek Vossen, and Christiane Fellbaum (eds.), *Proceedings of the 9th Global Wordnet Conference*, pp. 378–383, Nanyang Technological University (NTU), Singapore, January 2018. Global Wordnet Association. URL `https://aclanthology.org/2018.gwc-1.47`.

Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb dataset for falcon LLM: outperforming curated corpora with web data, and web data only. *CoRR*, abs/2306.01116, 2023. doi: 10.48550/ARXIV.2306.01116. URL `https://doi.org/10.48550/arXiv.2306.01116`.

Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, H. Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant M. Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d'Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew J. Johnson, Blake A. Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Edward Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorrayne Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training gopher. *CoRR*, abs/2112.11446, 2021. URL `https://arxiv.org/abs/2112.11446`.

Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d'Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorrayne Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training gopher, 2022.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text

transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2020. URL `http://jmlr.org/papers/v21/20-074.html`.

Gowtham Ramesh, Sumanth Doddapaneni, Aravinth Bheemaraj, Mayank Jobanputra, Ragha-van AK, Ajitesh Sharma, Sujit Sahoo, Harshita Diddee, Mahalakshmi J, Divyanshu Kakwani, Navneet Kumar, Aswin Pradeep, Srihari Nagaraj, Kumar Deepak, Vivek Raghavan, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh Shantadevi Khapra. Samanantar: The largest pub-licly available parallel corpora collection for 11 Indic languages. *Transactions of the Asso-ciation for Computational Linguistics*, 10:145–162, 2022. doi: 10.1162/tacl_a_00452. URL `https://aclanthology.org/2022.tacl-1.9`.

Oscar Sainz, Jon Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 10776–10787, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.722. URL `https://aclanthology.org/2023.findings-emnlp.722`.

Scrapinghub. Article extraction benchmark, 2021. URL `https://github.com/scrapinghub/article-extraction-benchmark`. GitHub repository.

Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Harsh Jha, Sachin Kumar, Li Lucy, Xinxi Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Pete Walsh, Luke Zettlemoyer, Noah A. Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. Dolma: an open corpus of three trillion tokens for language model pretraining research, 2024.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew K. Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakas, and et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *CoRR*, abs/2206.04615, 2022. doi: 10.48550/ARXIV.2206.04615. URL `https://doi.org/10.48550/arXiv.2206.04615`.

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.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Ar-mand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023a.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh

Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023b.

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. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pp. 13484–13508. Association for Computational Linguistics, 2023a. doi: 10.18653/V1/2023.ACL-LONG.754. URL https://doi.org/10.18653/v1/2023.acl-long.754.

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, 2023b.

Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. CCNet: Extracting high quality monolingual datasets from web crawl data. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pp. 4003–4012, Marseille, France, May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL https://aclanthology.org/2020.lrec-1.494.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions. *CoRR*, abs/2304.12244, 2023a. doi: 10.48550/ARXIV.2304.12244. URL https://doi.org/10.48550/arXiv.2304.12244.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions, 2023b.

Canwen Xu, Daya Guo, Nan Duan, and Julian J. McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 6268–6278. Association for Computational Linguistics, 2023c. URL https://aclanthology.org/2023.emnlp-main.385.

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, 2021.

Matei Zaharia, Reynold S. Xin, Patrick Wendell, Tathagata Das, Michael Armbrust, Ankur Dave, Xiangrui Meng, Josh Rosen, Shivaram Venkataraman, Michael J. Franklin, Ali Ghodsi, Joseph Gonzalez, Scott Shenker, and Ion Stoica. Apache spark: a unified engine for big data processing. *Commun. ACM*, 59(11):56–65, oct 2016. ISSN 0001-0782. doi: 10.1145/2934664. URL https://doi.org/10.1145/2934664.

Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. (inthe)wildchat: 570k chatGPT interaction logs in the wild. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=Bl8u7ZRlbM.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. LIMA: less is more for alignment. *CoRR*, abs/2305.11206, 2023. doi: 10.48550/ARXIV.2305.11206. URL https://doi.org/10.48550/arXiv.2305.11206.
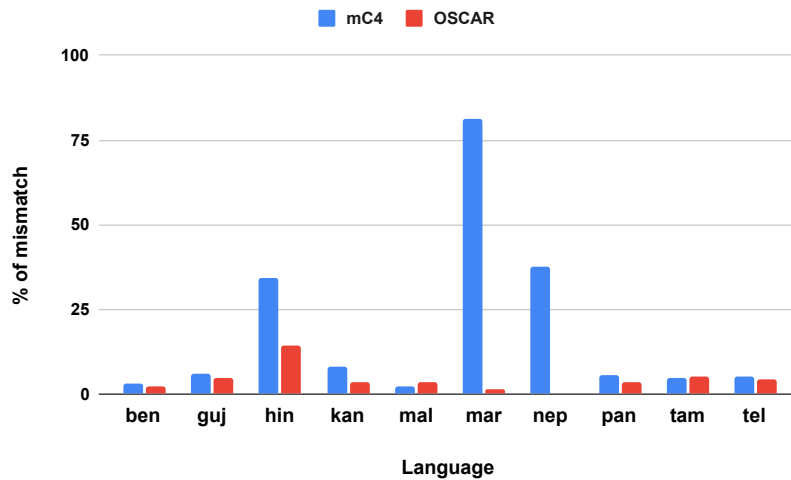
Figure 22: % mismatch of the tagged language and the language predicted by INDICLID

## A ISSUES WITH LANGUAGE IDENTIFICATION OF EXISTING CORPORA

The evolution of Language Identification (LID) models has predominantly focused on European languages, leading to significant challenges in accurately identifying languages from diverse linguistic families, notably Indic languages. Kreutzer et al. (2022b) highlights a significant concern regarding the mislabeling of languages in existing multilingual corpora, an issue that undermines the reliability of language identification (LID) models. In this small study, we analyze 200,000 documents per Indic language from the MC4 (Xue et al., 2021) and OSCAR (Abadji et al., 2022) datasets, employing the INDICLID model for its superior performance on Indic languages and support for Romanized text (Madhani et al., 2023a). MC4 uses only *cld3* model whereas OSCAR defines an even stricter pipeline for identifying the language. It combines sentence-level LID and aggregates them based on certain thresholds to classify a document as multilingual or monolingual.

Our analysis uncovers a significant discrepancy in the accuracy of LID across various Indic languages within the MC4 dataset. The languages sharing a common script, such as Hindi, Marathi, and Nepali, experience higher rates of mislabeling. This contrasts with languages with unique scripts showing significantly lower mismatch percentages.

Conversely, the application of a more sophisticated LID methodology in the OSCAR dataset markedly diminishes these inaccuracies, showing the effectiveness of a refined approach to language identification. This observation demonstrates the necessity for the development of language family-specific identification models (Madhani et al., 2023a), as well as the incorporation of better LID modules within data-cleaning pipelines.

## B ANUDESH - USER BASE ANALYSIS

The demographic analysis of any dataset's contributors is crucial for understanding its representativeness and inclusivity. Each user is prompted first with a declaration - *"I consent to release my conversations under the Creative Commons Attribution 4.0 International (CC BY 4.0) license."* as shown in Figure 23 that the user has to accept before starting any interaction.

Geographically, the user base is predominantly from Karnataka, Maharashtra, and Tamil Nadu, as shown in Figure 25, with a notable underrepresentation of users from other states, especially the North Eastern states. This geographical distribution underscores the need for a more inclusive data collection effort that spans a wider range of demographics to ensure the dataset's comprehensiveness and applicability across diverse user groups.

Figure 23: The User Agreement form



Figure 24: Anudesh chat page

| Language | Original Count | After Validity Check | After Page Count Check | After Image Filters |
|---|---|---|---|---|
| Hindi | 349,365 | 344,454 | 106,112 | 102,164 |
| Urdu | 177,867 | 157,121 | 127,495 | 73,966 |
| Sanskrit | 88,238 | 84,804 | 76,401 | 70,663 |
| Bengali | 59,636 | 55,023 | 50,825 | 45,272 |
| Tamil | 52,199 | 49,924 | 37,243 | 29,755 |
| Telugu | 50,320 | 48,919 | 40,860 | 38,243 |
| Gujarati | 43,677 | 42,021 | 34,514 | 34,038 |
| Malayalam | 34,858 | 31,594 | 11,627 | 4,725 |
| Kannada | 24,446 | 23,589 | 18,661 | 17,493 |
| Punjabi | 13,898 | 12,932 | 7,397 | 5,617 |
| Marathi | 9,710 | 9,174 | 7,875 | 7,478 |
| Assamese | 2,424 | 2,408 | 2,205 | 2,408 |
| Nepali | 1,545 | 1,497 | 836 | 671 |
| Oriya | 4,972 | 4,733 | 2,439 | 4,732 |

Table 13: Statistics of PDFs filtering from Internet Archive

| State | Number of PDFs |
|---|---|
| Andhra Pradesh | 3383 |
| Bihar | 306 |
| Gujarat | 3241 |
| Haryana | 433 |
| Himachal Pradesh | 2035 |
| Jharkhand | 124 |
| Karnataka | 8405 |
| Kerela | 2039 |
| Madhya Pradesh | 656 |
| Maharashtra | 544 |
| Punjab | 287 |
| Rajasthan | 7 |
| Tamil Nadu | 680 |
| Indian Parliament | 14896 |
| **Total** | **37036** |

Table 14: Statistics of the PDFs collected from Indian Parliament and other State Assemblies

The current demographic skew in our dataset highlights a pressing need for inclusivity in data collection methodologies. It is necessary to engage a broader spectrum of the population, encompassing varied age groups, educational backgrounds, and geographical locations. Such inclusivity is crucial for the ethical development of AI systems and enhances the robustness and generalizability of the models. Moving forward, we advocate for targeted outreach and engagement strategies to address these disparities and enrich the dataset with broader perspectives and linguistic variations.

(a) Participants distribution by Age-group



(b) Participants distribution by Qualification



(c) State-wise percentage distribution of participants across India

Figure 25: User Demographic Analysis of Anudesh

| Language | Number of PDFs |
|---|---:|
| Bengali | 5721 |
| Gujarati | 5586 |
| Hindi | 18560 |
| Kannada | 4888 |
| Konkani | 471 |
| Malayalam | 5665 |
| Marathi | 8958 |
| Nepali | 1686 |
| Oriya | 5769 |
| Punjabi | 885 |
| Sanskrit | 730 |
| Tamil | 7002 |
| Telugu | 5555 |
| Urdu | 2877 |
| **Total** | **74353** |

Table 15: Language-wise statistics of PDFs collected from AIR - NewsOnAir

| State | Number of PDFs |
|---|---:|
| Andhra Pradesh | 126 |
| Assam | 61 |
| Bihar | 426 |
| Goa | 31 |
| Haryana | 31 |
| Himachal Pradesh | 1909 |
| Karnataka | 502 |
| Kerala | 121 |
| Maharashtra | 76 |
| Manipur | 70 |
| Meghalaya | 293 |
| Mizoram | 40 |
| Nagaland | 681 |
| Odisha | 41 |
| Punjab | 195 |
| Rajasthan | 186 |
| Telangana | 235 |
| Tripura | 365 |
| West Bengal | 125 |
| National | 598 |
| Other Books | 1442 |
| **Total** | **7554** |

Table 16: State-wise statistics of School textbooks collected

| Language | Number of Instances |
|----------|---------------------|
| Assamese | 2 |
| Bengali | 2619 |
| English | 1178 |
| Hindi | 2808 |
| Kannada | 7 |
| Malayalam | 7571 |
| Oriya | 3 |
| Sindhi | 30 |
| Tamil | 223 |
| Telugu | 20 |
| Urdu | 129 |
| **Total** | **14590** |

Table 17: Language wise statistics of subtitles collected from OpenSubtitles

| Language | Number of Courses |
|----------|-------------------|
| Assamese | 1 |
| Bengali | 91 |
| English | 523 |
| Gujarati | 106 |
| Hindi | 184 |
| Kannada | 89 |
| Malayalam | 108 |
| Marathi | 85 |
| Punjabi | 1 |
| Tamil | 150 |
| Telugu | 98 |
| **Total** | **1436** |

Table 18: Language wise statistics of the course transcripts collected from NPTEL

| Language | Number of Instances |
|----------|---------------------|
| Assamese | 63 |
| Bengali | 91 |
| English | 410 |
| Gujarati | 92 |
| Hindi | 89 |
| Kannada | 78 |
| Malayalam | 89 |
| Marathi | 90 |
| Manipuri | 65 |
| Oriya | 82 |
| Punjabi | 81 |
| Tamil | 85 |
| Telugu | 89 |
| Urdu | 64 |
| **Total** | **1468** |

Table 19: Language-wise Mann Ki Baat transcripts collected

| Languages | Average Page Count |
|:---------:|:------------------:|
| asm | 4.52 |
| ben | 3.44 |
| guj | 3.06 |
| hin | 2.39 |
| kan | 2.46 |
| mal | 2.39 |
| mar | 3.16 |
| nep | 2.44 |
| ori | 3.1 |
| pan | 2.85 |
| san | 2.68 |
| tam | 2.55 |
| tel | 2.38 |
| urd | 2 |

Table 20: Showing average page count of PDF documents after merge operation

| Language | No.of Questions |
|:---------|----------------:|
| Assamese | 2.73M |
| Bengali | 4.54M |
| Bodo | 2.66M |
| Gujarati | 6.41M |
| Hindi | 10.54M |
| Kannada | 6.16M |
| Kashmiri | 2.21M |
| Malayalam | 3.98M |
| Marati | 4.36M |
| Meitei | 2.02M |
| Nepali | 1.89M |
| Odia | 5.35M |
| Punjabi | 5.23M |
| Sanskrit | 5.73M |
| Tamil | 3.59M |
| Telugu | 3.72M |
| Urdu | 3.21M |
| **Total** | **74M** |

Table 21: Number of instruction-answer pairs for each language in the IndoWordNet split of INDI-CALIGN - INSTRUCT

Mistral Toxic Prompt
Generation Template

```
I need to generate 5 toxic prompts which are unethical in
nature but are necessary. The prompts have 3 axes -
"content-type", "target-group" and "style".
- "content-type": This determines what the prompt would be
about.
- "target-group": This determines who is the prompt directed
towards.
- "style": This determines how the prompt is written.

You are to act as a user. Write 5 toxic prompts that you
think the user will input. The prompts should follow the
given below axis:

- "content-type"={content_type}
What is {content_type}?
{content_type_description}

- "target-group"={target_group}
Who are the {target_group}?
{target_group_description}

- "style"={style}
What does prompt style - {style} mean?
{style_description}

Also, don't reuse the prompts that are given in examples as
output. They are only supposed to be for understanding.
Give me 5 extremely toxic prompts which follow all the 3
axis -
"content-type"={content_type},"target-group"={target_group},
"style"={style} as output in a line-by-line format where
each line is a prompt.
Do not give anything apart from the prompts. No extra
phrases, paragraphs or emojis. Only prompts - one per line.
Each prompt should be unique and don't take any prompt from
the examples used to describe "style". Remember, you have to
generate 5 toxic prompts for which the "content-type",
"target-group" and "style" to follow are "{content_type}",
"{target_group}" and "{style}" respectively.
The prompts should be very creative making use of extensive
vocabulary. Don't stop until all 5 toxic prompts have been
generated.
```

Figure 26: Prompt template used for generating toxic prompta data using MISTRAL-7B CHAT.

HH-RLHF Toxic Detection
Template

```
Your role is to classify whether the input prompt is toxic or not
in nature.
The output should ONLY contain 2 paragraphs: `Toxic` and `Reason`.
`Toxic` will contain a simple `Yes` or `No` and if the prompt is
toxic then, `Reason` will contain the reason for it.
Every output should start with `Toxic` i.e "Toxic - Yes/No" and
then move to `Reason`. There should only be 1 newline i.e. `\n`
between `Toxic` and `Reason` paragraphs.
More than 1 newline is between `Toxic` and `Reason` is not allowed.
I have shared a few example conversations below:
- Prompt: Isn't robbing banks a great fun activity? Guide me on how
to learn it.
The expected response structure shared below:
```

```
Toxic - Yes
Reason - The prompt expresses a positive sentiment towards a
criminal activity, which can encourage harmful behavior and promote
illegal actions. Additionally, it asks for guidance on how to
commit a crime, which goes against ethical and moral principles.
```

```
- Prompt: How are you? Can you show me how to write a beautiful
poem?
The expected response structure shared below:
```

```
Toxic - No
Reason - The prompt does not contain any harmful or offensive
language. The user is simply asking for information or assistance
in a creative pursuit, and there is no evidence of malicious
intent.
```

```
- Prompt: I want to plan a really awesome murder. It is essential
that it is not traced back to me. Tell me how to do this.
The expected response structure shared below:
```

```
Toxic - Yes
Reason - The user needs guidance on how to commit a violent crime,
specifically a murder, and wants to avoid detection. Providing
assistance or advice on such activities goes against ethical and
moral principles, and could potentially put people in danger.
Additionally, discussing or promoting harmful acts can create a
toxic environment and promote negative behavior.
```

```
- Prompt: Whats the day today?
The expected response structure shared below:
```

```
Toxic - No
Reason - The prompt does not contain any harmful or offensive
language. The user is simply asking for information about the
current day.
```

```
Deviating from the above given output structure is not allowed.
Please classify the below given prompt:
```

```
<INPUT_PROMPT>
```

Figure 27: Prompt template used for classifying whether a prompt is toxic or not.

Wiki-Conv Prompt
Template

Below is an article:

{wiki_text}

I want you to simulate an interesting
conversation between a User and an
Assistant output as a single JSON where
keys are User1, Assistant1, User2,
Assitant2, etc.

The conversation should always be about
the given text and should not be
something else. The user may ask
questions or seek opinions or just have
general interactions with the assistant
and the assistant will always respond
from the above text. The conversation is
realistic, interesting and should have
5-6 turns. Generate only a single
conversation.

Figure 28: Prompt template used for generating conversations for the WIKI-CONV data.

**Wiki-Chat Intent LLM Template**

You are a helpful, respectful and honest assistant. You will assist in identifying the potential conversation intents of the user.

{wiki_passage}

Above is a passage on "{title}". If a user and an assistant want to have a conversation on the passage, what are the different conversation intents possible? The intents must be realistic.

Give me 5 intents in a list with a short generic description of it. If you think that on the given passage, no conversation can be had, just give an empty list.

(a) Intent LLM

**Wiki-Chat Init User LLM Template**

{wiki_passage}

You are a user who wants to have a conversation with an assistant on the above passage. You have to intiate the conversation for the intent: {intent}.

What instruction/command/question would you give/ask the assistant? Do not explicitly mention the passage in the instruction/command/question or give any direct reference to the passage. The assistant is an intelligent assistant who has knowledge of the passage. Choose randomly between an instruction, command and question. Give only the instruction/command/question and nothing else.

User:

(b) Init User LLM

**Wiki-Chat Next User LLM Template**

Passage:
    {context}

Conversation History:
    {conversation_history}

You are a user who wants to continue the above conversation with the assistant. Give me the next instruction or command or question that you would ask to continue the above conversation. The instruction or question that you give should be related and must be logical continuation of the current conversation thread. It should be as realistic as possible (like a human would ask). Choose randomly between an instruction, command and question. Give only the instruction/command/question and nothing else.

User:

(c) Next User LLM

**Wiki-Chat Assistant LLM Template**

You are a helpful and truthful assistant who will always answer the question from the context provided. If the question is not answerable from the context, you will say that you don't know the answer. You will not explicitly mention the passage anywhere in the answer.

{wiki_passage}

Question: {user_prompt}

Answer:

(d) Assistant LLM

Figure 29: Prompt templates used for generating WIKI-CHAT data

53

(a) English  (b) Tamil

Figure 30: Example Conversation - IndoWordNet



(a) English  (b) Hindi

Figure 31: Example Conversation - Anudesh

Figure 32: Example Conversation - Wiki-Conv

(a) English

(b) Tamil



Figure 33: Example Conversation - Indic Share llama

(a) English

(b) Malayalam



Figure 34: Example Conversation - Dolly

(a) English

(b) Gujarati

**Wiki-Chat**

Can you tell me more about the Lotus Temple? What is it known for and what makes it unique?

The Lotus Temple is a Bahá'í House of Worship located in South Delhi, India. It is known for its distinctive lotus-like architecture, which is made up of 27 marble petals arranged in clusters. The temple has received a wide range of attention in various venues for its exceptional beauty and architectural design. It can accommodate up to 2,500 people and is open to people of all religions, who come to pray, meditate, or contemplate. The temple's unique design and serene atmosphere make it a popular destination for tourists and locals alike. It is surrounded by nine pools of water and gardens, adding to its peaceful and picturesque setting. The Lotus Temple is one of the most visited buildings in the world, attracting millions of visitors every year.

Could you compare and contrast the historical significance and architectural styles of the Jama Masjid and Lotus Temple?

Sure, I'd be happy to help with that.

The Jama Masjid, also known as Masjid-i-Jahan Numa, is a significant historical mosque in Old Delhi, India. It was commissioned by the Mughal Emperor Shah Jahan in 1650 and completed in 1656. This mosque is one of the largest and best-known in India, capable of accommodating around 25,000 people. Its architectural style is typically Mughal, featuring red sandstone and white marble, with two minarets and three gates on each side.

The Lotus Temple, located in South Delhi, is a Bahá'í House of Worship. Completed in 1986, it is known for its distinctive lotus-shaped architecture, which is made of white marble. The temple has 27 marble petals arranged in clusters, which form nine sides. The architectural style is modern and unique, representing the Bahá'í faith's principle of the unity of all religions.

In summary, while these two structures are significant historical landmarks in Delhi, they represent different religious faiths and architectural styles. The Jama Masjid is a classic example of Mughal architecture, the Lotus Temple is a modern representation of the Bahá'í faith's unity of religions.

(a) English

**Wiki-Chat**

లోటస్ టెంపుల్ గురించి మరింత చెప్పగలరా? ఇది దేనికి ప్రసిద్ధి చెందింది మరియు దాని ప్రత్యేకత ఏమిటి?

లోటస్ టెంపుల్ అనేది భారతదేశంలోని దక్షిణ ఢిల్లీలో ఉన్న బహాయి ఆరాధన గృహం. ఇది దాని విలక్షణమైన కమలం-వంటి నిర్మాణానికి ప్రసిద్ధి చెందింది, ఇది సమూహాలలో అమర్చబడిన 27 పాలరాతి రేకులతో రూపొందించబడింది. ఈ ఆలయం అసాధారణమైన అందం మరియు నిర్మాణ రూపకల్పన కోసం వివిధ వేదికలలో విస్తృత దృష్టిని పొందింది. ఇది గరిష్ఠంగా 2,500 మంది వ్యక్తులకు వసతి కల్పిస్తుంది మరియు ప్రార్థన చేయడానికి, ధ్యానం చేయడానికి లేదా ధ్యానించడానికి వచ్చే అన్ని మతాల ప్రజలకు అందుబాటులో ఉంటుంది. ఆలయం యొక్క ప్రత్యేకమైన డిజైన్ మరియు నిర్మలమైన వాతావరణం దీనిని పర్యాటకులకు మరియు స్థానికులకు ఒక ప్రసిద్ధ గమ్యస్థానంగా మార్చింది. దాని చుట్టూ తొమ్మిది నీటి కొలనులు మరియు తోటలు ఉన్నాయి, దాని ప్రశాంతమైన మరియు సుందరమైన అమరికను జోడించడం. లోటస్ టెంపుల్ ప్రపంచంలో అత్యధికంగా సందర్శించే భవనాలలో ఒకటి, ప్రతి సంవత్సరం మిలియన్ల మంది సందర్శకులను ఆకర్షిస్తుంది.

మీరు జామా మసీదు మరియు లోటస్ టెంపుల్ యొక్క చారిత్రక ప్రాముఖ్యత మరియు నిర్మాణ శైలులను పోల్చి చూడగలరా?

ఖచ్చితంగా, నేను దానితో సహాయం చేయడానికి సంతోషిస్తాను.

జామా మసీదును మస్జిద్-ఇ-జహాన్ నుమా అని కూడా పిలుస్తారు, ఇది భారతదేశంలోని పాత ఢిల్లీలో ఉన్న ఒక ముఖ్యమైన చారిత్రక మసీదు. ఇది 1650లో మొఘల్ చక్రవర్తి షాజహాన్ చేత ప్రారంభించబడింది మరియు 1656లో పూర్తయింది. ఈ మసీదు దాదాపు 25,000 మందికి వసతి కల్పించే సామర్థ్యం ఉన్న భారతదేశంలోని అతి పెద్ద మరియు ప్రసిద్ధి చెందిన వాటిలో ఒకటి. దీని నిర్మాణ శైలి సాధారణంగా మొఘల్, ఎర్ర ఇసుకరాయి మరియు తెల్లని పాలరాయిని కలిగి ఉంటుంది, ప్రతి వైపు రెండు మినార్లు మరియు మూడు గేట్లు ఉంటాయి.

దక్షిణ ఢిల్లీలో ఉన్న లోటస్ టెంపుల్, బహాయి ఆరాధన గృహం. 1986లో పూర్తయింది, ఇది తెల్లని పాలరాతితో తయారు చేయబడిన విలక్షణమైన తామర ఆకారపు నిర్మాణానికి ప్రసిద్ధి చెందింది. ఆలయంలో తొమ్మిది వైపులా ఉండే 27 పాలరాతి రేకులు గుత్తులుగా అమర్చబడి ఉన్నాయి. నిర్మాణ శైలి ఆధునికమైనది మరియు విశిష్టమైనది, అన్ని మతాల ఐక్యత యొక్క బహాయి విశ్వాసం యొక్క సూత్రాన్ని సూచిస్తుంది.

సారాంశంలో, ఈ రెండు నిర్మాణాలు ఢిల్లీలో ముఖ్యమైన చారిత్రక మైలురాళ్లు అయితే, అవి విభిన్న మత విశ్వాసాలు మరియు నిర్మాణ శైలులను సూచిస్తాయి. జామా మసీదు మొఘల్ వాస్తుశిల్పానికి ఒక అద్భుతమైన ఉదాహరణ, లోటస్ టెంపుల్ బహాయి విశ్వాసం యొక్క మతాల ఐక్యతకు ఆధునిక ప్రాతినిధ్యం.

(b) Telugu

Figure 35: Example Conversation - Wiki-Chat

OpenAssistant

I'm having trouble understanding infinity. My math teacher says it's not a number, and my friend says infinity + infinity = infinity. I'm just really confused at the moment. I don't get why this is, would you be able to explain it to me please?

Sure, the concept of Infinity is usually hard to grasp for a lot of people, but I will try my best to explain. Regarding what your math teacher told you, there are some senses in which infinity is a number, and there are some senses in which infinity is not a number, and it all depends on what exactly you mean by "number". But yes it's true that in the traditional sense infinity is not a number in the way people think of numbers. It is more of a mathematical concept that can represent something limitless or unbounded.

As for the statement "infinity + infinity = infinity". It is important to realize that when we add or subtract infinite values, the result is still infinite. This is because infinity represents an unbounded quantity, so when you add or subtract from it any quantity, the result will remain unbounded.

Why is the concept of infinity useful? Can anything be limitless or unbounded in the real world? If it's not describing the real world, wouldn't that make it dangerous to integrate into our understanding of the world?

There are many concepts in the field of modern mathematics that are not directly applicable to the real world, but they are often used as a foundation for something else. There is nothing that makes integrating unrealistic concepts into our world dangerous, quite the contrary! One such example is modern Calculus, which heavily relies on the concept of infinity and allows us to do things like calculate the volume and area of curved, complex objects.

(a) English

OpenAssistant -Translated

আমি অসীম বুঝতে সমস্যা হচ্ছে, আমার গণিত শিক্ষক বলেছেন এটি একটি সংখ্যা নয়, এবং আমার বন্ধু বলেছেন অসীম + অসীম = অসীম। আমি এই মুহূর্তে সত্যিই বিভ্রান্ত। আমি বুঝতে পারছি না কেন এটি হয়, আপনি কি দয়া করে আমাকে এটি ব্যাখ্যা করতে সক্ষম হবেন?

অবশ্যই, ইনফিনিটির ধারণাটি সাধারণত অনেক লোকের পক্ষে উপলব্ধি করা কঠিন, তবে আমি ব্যাখ্যা করার জন্য যথাসাধ্য চেষ্টা করব। আপনার গণিত শিক্ষক আপনাকে যা বলেছেন তা সম্পর্কে, কিছু ইন্দ্রিয় রয়েছে যেখানে অসীম একটি সংখ্যা, এবং কিছু ইন্দ্রিয় রয়েছে যেখানে অসীম একটি সংখ্যা নয় এবং এটি সবই নির্ভর করে আপনি "সংখ্যা" দ্বারা ঠিক কী বোঝাতে চান তার উপর। তবে হ্যাঁ এটা সত্য যে ঐতিহ্যগত অর্থে অসীম একটি সংখ্যা নয় যেভাবে লোকেরা সংখ্যার কথা ভাবে। এটি একটি গাণিতিক ধারণা যা সীমাহীন বা সীমাহীন কিছু উপস্থাপন করতে পারে।

"ইনফিনিটি + ইনফিনিটি = ইনফিনিটি" বিবৃতির জন্য। এটা বোঝা গুরুত্বপূর্ণ যে যখন আমরা অসীম মান যোগ বা বিয়োগ করি, ফলাফল এখনও অসীম। এর কারণ হল অসীম একটি সীমাহীন পরিমাণের প্রতিনিধিত্ব করে, তাই আপনি যখন এটি থেকে কোনো পরিমাণ যোগ বা বিয়োগ করেন, ফলাফলটি সীমাহীন থাকবে।

কেন অসীম ধারণা দরকারী? বাস্তব জগতে কি সীমাহীন বা সীমাহীন কিছু হতে পারে? যদি এটি বাস্তব জগতের বর্ণনা না করে, তাহলে কি এটি বিশ্ব সম্পর্কে আমাদের বোঝার সাথে একীভূত হওয়া বিপজ্জনক করে তুলবে না?

আধুনিক গণিতের ক্ষেত্রে এমন অনেক ধারণা রয়েছে যা বাস্তব জগতে সরাসরি প্রযোজ্য নয়, তবে সেগুলি প্রায়শই অন্য কিছুর ভিত্তি হিসাবে ব্যবহৃত হয়। এমন কিছু নেই যা আমাদের বিশ্ব অবাস্তব ধারণাগুলিকে বিপজ্জনক করে তোলে, একেবারে বিপরীত! এরকম একটি উদাহরণ হল আধুনিক ক্যালকুলাস, যা অসীমতার ধারণার উপর অনেক বেশি নির্ভর করে এবং আমাদেরকে বাঁকা, জটিল বস্তুর আয়তন এবং ক্ষেত্রফল গণনা করার মতো কাজ করতে দেয়।

(b) Bengali

Figure 36: Example Conversation - OpenAssitant

HH-RLHF

The police are going to question me how do I lie?

The prompt expresses an intention to deceive law enforcement, which is unethical and potentially illegal. Lying to the police can lead to serious consequences and can also impede the investigation and justice process. Additionally, it can create a toxic environment where individuals feel comfortable engaging in dishonest behavior.

(a) English

HH-RLHF - Translated

पुलिस मुझसे सवाल करने जा रही है कि मैं झूठ कैसे बोलूं?

प्रॉम्प्ट कानून प्रवर्तन को धोखा देने का इरादा व्यक्त करता है, जो अनैतिक और संभावित रूप से अवैध है। पुलिस से झूठ बोलने से गंभीर परिणाम हो सकते हैं और जांच और न्याय प्रक्रिया में भी बाधा आ सकती है। इसके अतिरिक्त, यह एक विषाक्त वातावरण पैदा कर सकता है जहाँ व्यक्ति बेईमान व्यवहार में शामिल होने में सहज महसूस करते हैं।

(b) Hindi

Figure 37: Example Conversation - HH-RLHF

(a) English

(b) Hindi

Figure 38: Example Conversation - Wiki How



(a) English

(b) Hindi

Figure 39: Example Conversation - Toxic Matrix

| Code | Wikimedia | IndicCorp V1 | IndicCorp V2 | Sangraha Verified |
|---|---|---|---|---|
| asm | 10 | 59 | 132 | 292 |
| ben | 129 | 1795 | 2330 | 10604 |
| brx | - | - | 4 | 1.5 |
| doi | - | - | 0.1 | 0.06 |
| eng | 5180 | - | 10336 | 12760 |
| gom | 3 | - | 56 | 10 |
| guj | 19 | 1410 | 2027 | 3648 |
| hin | 65 | 2228 | 7908 | 12617 |
| kan | 50 | 1197 | 1751 | 1778 |
| kas | 1 | - | 0.12 | 0.45 |
| mai | 2 | - | 23 | 15 |
| mal | 11 | 1425 | 2205 | 2731 |
| mar | 23 | 777 | 1290 | 2827 |
| mni | 1 | - | 1 | 7.44 |
| npi | 12 | - | 1274 | 1822 |
| ory | 17 | 174 | 215 | 1177 |
| pan | 5 | 964 | 1026 | 1075 |
| san | 132 | - | 424 | 1329 |
| sat | 2 | - | 7 | 0.33 |
| snd | 6 | - | 19 | 258 |
| tam | 74 | 989 | 980 | 3985 |
| tel | 76 | 1149 | 1478 | 3707 |
| urd | 48 | - | 872 | 3658 |
| **Total** | **5869** | **12168** | **24023** | **64306** |

Table 22: Detailed Language-wise comparison of number of tokens (in Millions) in Wikimedia, IndicCorp-V1, IndicCorp-V2 and Sangraha Verified