

Towards Building Large Scale Datasets and State-of-the-Art Automatic Speech Translation Systems for 14 Indian Languages

Ashwin Sankar^{1,2*} Sparsh Jain^{1*} Nikhil Narasimhan¹ Devilal Choudhary^{4†}

Dhairya Suman^{3‡} Mohammed Safi Ur Rahman Khan^{1,2}

Anoop Kunchukuttan^{1,5} Mitesh M Khapra^{1,2} Raj Dabre^{1,2,6,7‡§}

¹Nilekani Centre at AI4Bharat ²Indian Institute of Technology, Madras

³Indian Institute of Technology, Delhi ⁴Delhi Technological University

⁵Microsoft ⁷Indian Institute of Technology, Bombay

 huggingface.co/BhasaAnuvaad  github.com/BhasaAnuvaad

Abstract

Speech translation for Indian languages remains a challenging task due to the scarcity of large-scale, publicly available datasets that capture the linguistic diversity and domain coverage essential for real-world applications. Existing datasets cover a fraction of Indian languages and lack the breadth needed to train robust models that generalize beyond curated benchmarks. To bridge this gap, we introduce BHASAAANUVAAD, the largest speech translation dataset for Indian languages, spanning over 44 thousand hours of audio and 17 million aligned text segments across 14 Indian languages and English. Our dataset is built through a threefold methodology: (a) aggregating high-quality existing sources, (b) large-scale web crawling to ensure linguistic and domain diversity, and (c) creating synthetic data to model real-world speech disfluencies. Leveraging BHASAAANUVAAD, we train INDIC-SEAMLESS, a state-of-the-art speech translation model for Indian languages that performs better than existing models. *Our experiments demonstrate improvements in the translation quality, setting a new standard for Indian language speech translation.* We will release all the code, data and model weights in the open-source, with permissive licenses to promote accessibility and collaboration.

1 Introduction

Automatic Speech Translation (AST) has become crucial to break language barriers and enable communication across languages and cultures. Traditionally, speech translation systems have relied on cascaded architectures, where Automatic Speech Recognition (ASR) is followed by Machine Translation (MT). However, recent advances have led to more integrated end-to-end (E2E) models (Babu

et al., 2021; Pratap et al., 2024; Radford et al., 2023; Communication et al., 2023) that directly translate speech from one language into text in another. Additionally, Audio-LLM-based systems have emerged (Wu et al., 2023; Chu et al., 2023; Fathullah et al., 2023; Gaido et al., 2024) further advancing the field. Despite these advancements, progress has been concentrated on English and other European languages, while low- and mid-resource languages, including Indian languages, have remained largely underrepresented. The primary challenge is the lack of training data, which leads to suboptimal models.

India’s linguistic diversity, with 22 officially recognized languages and numerous dialects, presents significant challenges for speech translation, including diverse phonetic systems (Mujadia and Sharma, 2023), frequent code-switching (Shankar et al., 2024), and syntactic variations. However, the primary bottleneck remains the scarcity of large-scale, high-quality datasets, aggravated by the limited web presence of Indian languages. To bridge this gap, we introduce BHASAAANUVAAD, a large-scale speech translation dataset comprising of over 44K hours of speech translation data covering 14 Indian languages alongside English, making it the most comprehensive speech translation resource for the region to date. The included languages - Assamese, Bengali, Gujarati, Hindi, Kannada, Malayalam, Marathi, Nepali, Odia, Punjabi, Sindi, Tamil, Telugu, and Urdu - span both medium and low resource languages. Our dataset construction follows a three-pronged approach to ensure high-quality and diverse data collection: **(1) Aggregating High-Quality Existing Sources:** We incorporate publicly available speech-text pairs from prior works, refining them with improved alignment techniques to enhance usability for speech translation tasks. **(2) Large-Scale Web Crawling:** To expand linguistic and domain coverage, we employ an automated web-mining pipeline that

* Equal Contribution.

† Work done during internship at AI4Bharat.

‡ Work done while at NICT, Japan.

§ Corresponding Author: raj.dabre@cse.iitm.ac.in

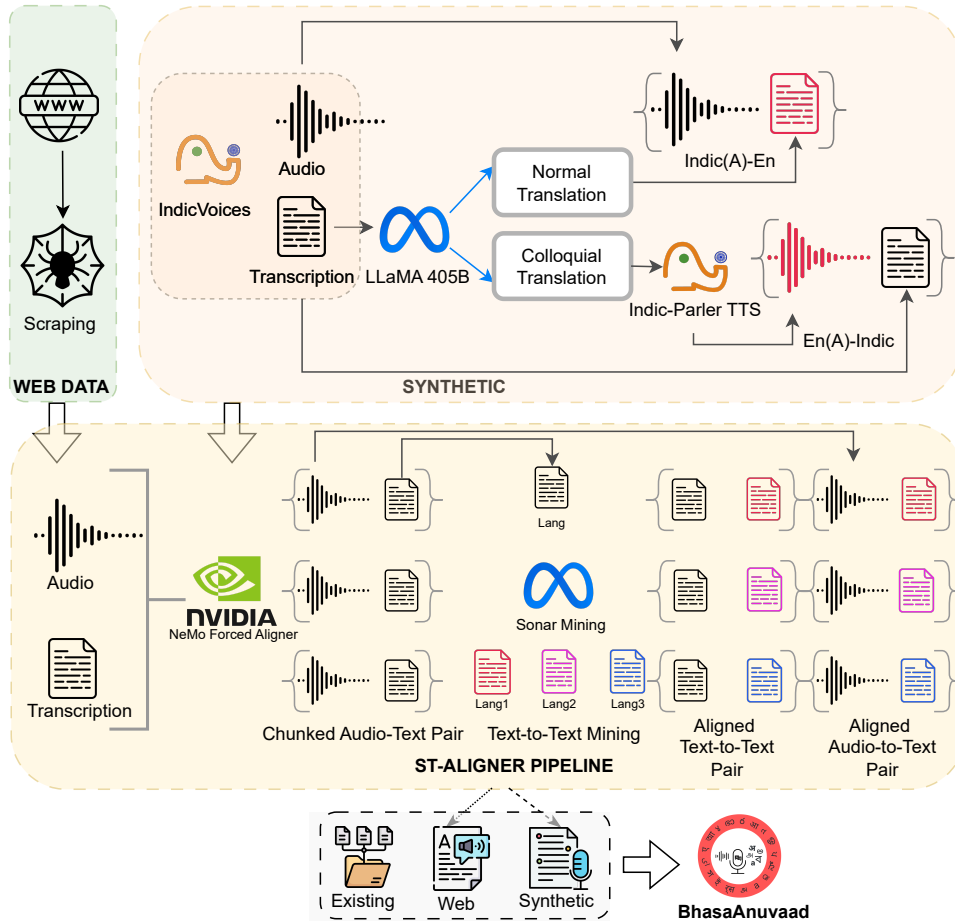


Figure 1: Overview of the creation of BhasaAnuvaad.

we call ST-ALIGNER (Figure 1) to extract speech-text pairs from diverse online sources, including government portals, educational materials, spiritual discourses, and multilingual podcasts. This enables us to capture naturally occurring speech across varied contexts, improving real-world applicability. **(3) Synthetic Data Generation:** To address gaps in speech diversity, we leverage the INDICVOICES ASR dataset and use LLAMA-3.1-405B-INSTRUCT for synthetic translations into 13 languages, while generating English speech using a TTS model (Sankar et al., 2025).

We also build BHASAANUVAAD-TEST by randomly sampling 20 minutes of data for each language pair from BHASAANUVAAD. Additionally, we introduce INDICSEAMLESS, a finetuned SEAMLESSM4T-V2-LARGE model on BHASAANUVAAD that surpasses existing AST setups - both cascaded and end-to-end - on established benchmark. We observe an overall average improvement of 10 chrF++ scores over previous approaches on BHASAANUVAAD-TEST demonstrating the efficacy of BHASAANUVAAD.

In summary, our contributions are: **(i)** the release of BHASAANUVAAD, the *largest* Indic-language speech translation corpus, comprising over 44 thousand hours of data across 14 Indian languages and English, **(ii)** a speech-translation data mining pipeline, ST-ALIGNER **(iii)** an *evaluation* of popular AST baselines on Conneau et al. (2022) and BHASAANUVAAD-TEST, a test split from BHASAANUVAAD, and **(iv)** the release of state-of-the-art speech translation model for 14 Indian languages. All code and datasets developed as part of this work will be made publicly available to support future research in AST for Indic languages.

2 Related Work

Cascaded and End-to-End AST. Cascaded Automatic Speech Translation (AST) has been explored extensively in early research (Di Gangi et al., 2019; Cheng et al., 2019; Bahar et al., 2020), connecting Automatic Speech Recognition (ASR) and Machine Translation (MT) systems in sequence. These systems suffer from cascaded error propagation (Sperber and Paulik, 2020; Beck et al.,

2019), loss of prosodic information (Bentivogli et al., 2021), and increased computational complexity (Lee et al., 2022). In contrast, end-to-end (E2E) models can avoid this (Weiss et al., 2017; Inaguma et al., 2021), but suffer from lack of good quality data (Sethiya and Maurya, 2023). With the rise of modern Large Language Models (LLMs), recent works have started exploring the feasibility of using Audio-based LLMs (or Speech LLMs), which combine Speech Foundation models (SFMs) and LLMs (Fathullah et al., 2023; Wu et al., 2023; Huang et al., 2023; Gaido et al., 2024).

Datasets for AST. The collection of speech-translation datasets for training end-to-end models has proven to be a persistent challenge. While efforts like SpeechMatrix (Duquenne et al., 2022) and SeamlessAlign (Communication et al., 2023) attempt to mine data from the web, the resulting datasets are often noisy. Most research has focused on curating data using comparable corpora from domains like Education (Ďurišková et al., 2024; Song et al., 2019), TED talks (Salesky et al., 2021), Parliament Debates (Koehn, 2005), among others. Additionally, with the rise of generative models, recent works have also started exploring the generation of large amounts of synthetic speech translation data in a flexible and cost-effective manner (Ye et al., 2023), (Bamfo Odoom et al., 2024). Although the use of synthetic data has been widely studied in text-based machine translation (Sennrich et al., 2016; Gala et al., 2023), its application in Speech translation, particularly for low-resource languages like Indic languages, remains underexplored, a gap which we fill.

Indian Languages AST. Unlike English, the shortage of parallel datasets has hindered the progress in Indian Languages (Mujadia and Sharma, 2023; Mhaskar et al., 2023). Multilingual speech translation models like MSLAM (Bapna et al., 2022), WHISPER (Radford et al., 2023), SEAMLESS (Communication et al., 2023) have made strides in addressing this gap, but their effectiveness in low-resource and real-world spontaneous speech settings remains largely unexplored. Furthermore, speech-based benchmarks for Indian languages are still underdeveloped (Sethiya and Maurya, 2023), with FLEURS (Conneau et al., 2022) being the primary option. In this paper, we push the boundaries of AST for Indian languages with the help of the large-scale dataset we construct.

3 BHASAANUVAAD

BHASAANUVAAD is the largest spoken translation dataset for any Indian language, totaling approximately 44,400 hours of speech and text data. Its construction involved a multi-step approach, as illustrated in Figure 1, combining (i) Aggregating existing datasets, (ii) Mining parallel speech and text data from comparable sources, and (iii) Synthetically generated speech/text parallel data. A language wise composition of this dataset is illustrated in Figure 2. A detailed breakdown of language and dataset-wise statistics is provided in Appendix A in Tables A1 & A2.

3.1 Aggregating Existing Datasets

We begin by aggregating several existing datasets in Indian Languages, which form the foundation of BHASAANUVAAD. These include:

Indic-TEDST: The Indic-TEDST (Sethiya et al., 2024) dataset includes the audio, transcripts, and translations of popular TED Talks in nine widely spoken Indian languages - Bengali, Gujarati, Hindi, Kannada, Malayalam, Marathi, Punjabi, Tamil, and Telugu - totaling approximately 203 hours of English-to-Indic data.

FLEURS: The FLEURS-train dataset (Conneau et al., 2022) is a multilingual speech corpus covering over 100 languages for training automatic speech recognition (ASR) and speech translation (AST) models. FLEURS was previously the largest source of data, covering 14 Indian languages, making it a valuable resource.

CVSS: The CVSS (Jia et al., 2022) corpus is a multilingual speech-to-speech translation dataset derived from CoVoST 2 (Wang et al., 2020), which is sampled from the Common Voice dataset. Spanning 21 languages, it includes only one Indian language, Tamil, with approximately 3 hours of data.

Khan Academy Corpus: The Khan Academy corpus (Ďurišková et al., 2024) provides data for the education domain in 29 languages. It includes 3 Indian languages, with 105 hours of audio.

SeamlessAlign: We recreate the SeamlessAlign data introduced by Communication et al. (2023) to get around 7500 hours of data for 5 Indic languages.

3.2 Mining from comparable sources

We explore the extraction of parallel speech-text data from various online sources containing high-quality content across multiple domains. While

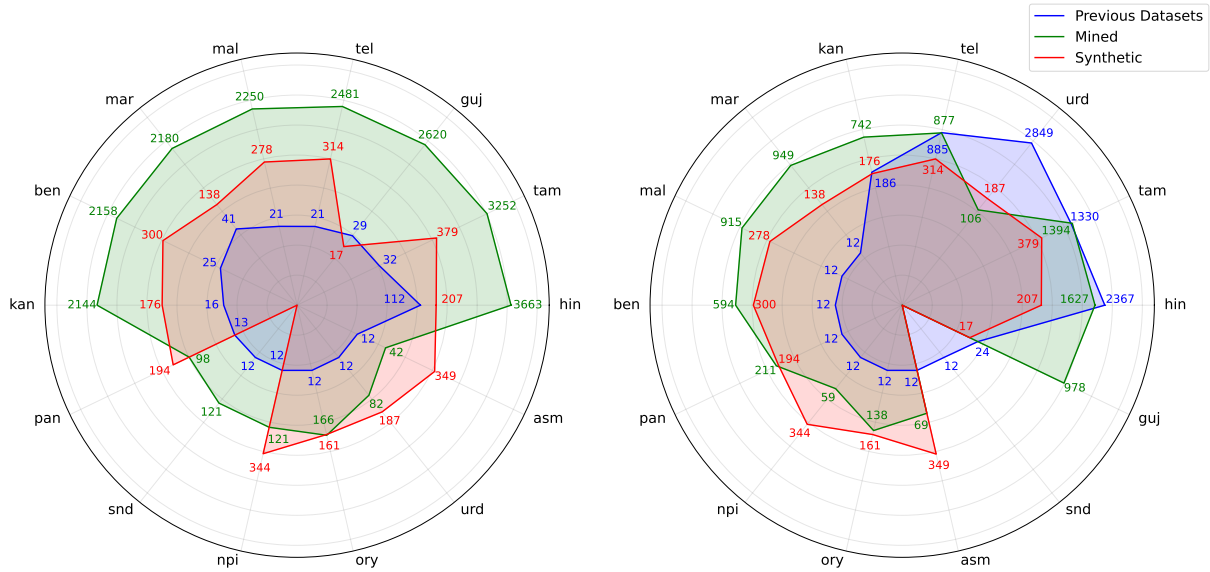


Figure 2: Distribution of dataset hours in log scale for En-Indic (left) and Indic-En (right) language pairs across three data sources: Previous Datasets (blue), Curated & Mined (green), and Synthetic (red).

many sources provide transcripts in different languages, these are often misaligned with the corresponding audio. To address this issue, we develop an automated pipeline that aligns transcripts with their corresponding source audios, to produce high-quality parallel speech-text data.

3.2.1 ST-Aligner

Our pipeline consists of four stages: normalization, punctuation restoration, audio-transcript alignment, and parallel text alignment.

1. Normalization: To ensure compatibility, all audio files are first converted to a mono channel and resampled at 16 kHz, as ASR models operate on 16 kHz audio. The corresponding transcripts are then cleaned by removing commonly found noise like extraneous symbols, extra spaces, and line breaks. We also replace the terminal punctuation marks with a non-standard character to aid sentence segmentation during alignment.

2. Punctuation Restoration: Many subtitles and transcripts lack proper punctuation, which can hinder accurate alignment. We conducted a small experiment with *indic-punct* (Gupta et al., 2022), a punctuation restoration model, and found that it did not perform well in spoken contexts. Given this limitation, we opt for the LLAMA-3.1-405B-INSTRUCT model, which demonstrated superior performance in restoring punctuation and effectively structuring the text for alignment (Figure A1).

3. Audio-Transcript Alignment: The cleaned

transcripts and audio are processed using the Nemo Forced Aligner (NFA)¹, which generates token-level, word-level, and segment-level timestamps based on CTC-based ASR models. Given the linguistic diversity of our dataset, we utilize INDIC-CONFORMER ASR models for Indian languages and Nvidia FastConformer² for English, ensuring that the alignment is optimized for the phonetic and acoustic characteristics of each language.

Long-form audio presents additional challenges, as ASR models and forced aligners struggle with excessive sequence lengths, often resulting in degraded alignment accuracy. To address this, we employ Silero Voice Activity Detection (Silero VAD) (Team, 2024), a lightweight yet highly effective neural model that detects speech boundaries and segments long recordings into smaller, more manageable chunks. This pre-segmentation step not only improves computational efficiency but also enhances alignment precision by reducing errors that arise from excessive length mismatches between audio and transcript.

4. Parallel Text Alignment: Once audio segments are paired with transcripts, we align them with their corresponding target language text. We generate sentence embeddings for both the source and target texts using SONAR (Duquenne et al., 2023). Next,

¹https://github.com/NVIDIA/NeMo/tree/main/tools/nemo_forced_aligner

²https://catalog.ngc.nvidia.com/orgs/nvidia/teams/nemo/models/stt_en_fastconformer_hybrid_large_pc

we compute cosine similarity scores between these embeddings to identify the most likely parallel sentence pairs.

3.2.2 Sources of Mining Corpora

We apply the pipeline described in Section 3.2.1 on the following sources to extract translation pairs.

Spoken Tutorial: Spoken tutorial, is an initiative under the National Mission on Education under the Ministry of Human Resource Development, Government of India to enable the translation of various vocational educational content in different languages. This dataset covers Indic-En translations with 499.04 hours of educational tutorials in 12 Indian languages.

UGCE Resources: The University Grants Commission translates a wide range of university-level courses into multiple Indian languages. Originally delivered in English (En audio), these courses have their transcripts translated into 8 additional languages, resulting in over 430 hours of aligned audio-text data for each language pair.

VaaniPedia: Vaanipedia is a comprehensive collection of religious and spiritual audio-text content in 14 Indian languages, comprising over 121 hours of aligned audio shared across each language pair.

Mann ki Baat: Mann Ki Baat is a monthly radio broadcast by the Prime Minister of India that addresses various national issues. Each episode is transcribed and manually translated into 13 Indian languages. These translations are also manually recorded, resulting in over 435 hours of speech data across all language pairs, covering both Indic-to-English and English-to-Indic directions. *Notably, Mann Ki Baat is the only source in this dataset that includes Manipuri, which is provided in Bengali script (rather than the Meitei script). While tables and figures do not explicitly highlight it, we are releasing the Manipuri data in Bengali script as sourced from Mann Ki Baat.*

WordProject: This dataset consists of Bible audiobooks sourced from the Word Project website in 12 Indian languages, totaling approximately 1,612 hours of speech data across both Indic-En and En-Indic directions.

NPTEL: The National Programme on Technology Enhanced Learning (NPTEL) is an Indian e-learning platform for university-level science, technology, engineering, and mathematics (STEM) subjects. It is the largest e-repository in the world, offering courses in engineering, basic sciences,

and selected humanities and management subjects³. These courses are primarily recorded in English, but they are transcribed and re-recorded in multiple languages manually. This is the largest component of our dataset, contributing over 22,500 hours of spoken content across 10 Indian languages. It covers a wide range of STEM subjects and includes both En-XX and XX-En translations, making it one of the most significant sources in our collection.

While these datasets are high quality, they are not diverse with majority of the content being from Educational domain. Furthermore, the audio styles they cover lack the diversity necessary for models to adapt effectively to real-world, where background noise, spontaneous speech, and varying accents are common. To overcome these limitations, we augment our dataset by using synthetic data generation as described below.

3.3 Synthetic Data Generation

To ensure that we have audio diversity across different domains, demographics, and languages, we use INDICVOICES (Javed et al., 2024) which is the largest collection of natural and spontaneous speech for Indian languages containing over 4,000 hours of transcribed speech. Its extensive language diversity makes it well-suited for training models capable of performing effectively in real-world scenarios.

Synthetic Translations: The INDICVOICES dataset consists of spontaneous speech, which includes disfluencies, repetitions, and informal constructions that pose challenges for direct text-based translation. Since INDICVOICES also provides normalized transcripts, where such artifacts are already addressed, we use these instead of verbatim transcripts for translation. This ensures cleaner and more structured inputs while retaining the essence of conversational speech.

To generate high-quality English translations, we first apply punctuation restoration and translate the unsegmented version of the transcripts using LLAMA-3.1-405B-INSTRUCT (Figure A1). By pairing the original Indic audio with these translated transcripts, we construct XX-En pairs suitable for speech translation tasks. We then utilize STALIGNER to align the translations at the segment level, preserving synchronization between the audio and textual representations.

³https://en.wikipedia.org/wiki/National_Programme_on_Technology_Enhanced_Learning

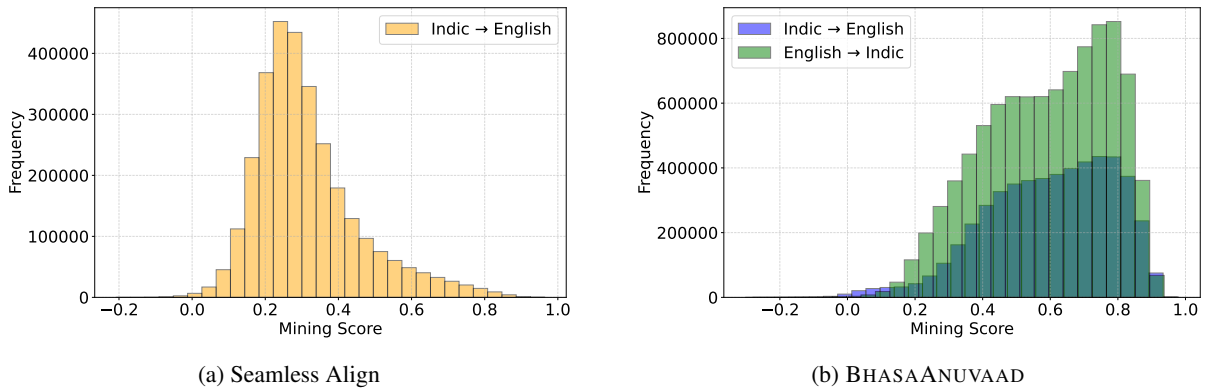


Figure 3: Distribution of SONAR mining scores for $XX \rightarrow \text{En}$ sentence pairs from the SeamlessAlign (left) and BHASAANUVAAD (right) datasets. BHASAANUVAAD consistently exhibits a higher concentration of high-quality alignments, with mining scores skewed toward the upper end of the scale. In contrast, SeamlessAlign distributions are centered around lower mining scores, indicating overall lower alignment quality. This comparison underscores the improved precision of alignments in BHASAANUVAAD .

Unlike specialized translation models such as INDICTRANS2 (Gala et al., 2023), which are optimized for formal text and often struggle with colloquial or spontaneous speech, LLAMA-3.1-405B-INSTRUCT offers a *steerable* alternative. Its instruction-following capability enables more flexible adaptation to informal and conversational registers through prompting with examples and contextual cues.

To validate this, we conducted a small-scale internal human evaluation on a subset of the INDICVOICES dataset. LLAMA-3.1-405B-INSTRUCT consistently outperformed INDICTRANS2, particularly in preserving the nuances and fluency of informal spoken language. This makes it well-suited for real-world applications, where such variability is common. In a broader multilingual evaluation, we assessed the translation quality of LLAMA-3.1-405B-INSTRUCT across 1,000 documents per language. The model achieved an average human-rated score of 8 out of 10, indicating adequacy for our use case.

Synthetic Speech Generation: To overcome the time and cost limitations of manual voice recordings, we augment our dataset synthetic speech generation. We begin with verbatim transcripts in INDICVOICES, which preserve natural speech phenomena such as disfluencies, false starts, repetitions, and pauses. These are translated into colloquial English using the LLAMA-3.1-405B-INSTRUCT model, producing outputs that reflect spontaneous conversational patterns, including hesitations and filler words (Figure A3).

To synthesize English speech with prosodic

and pronunciation features characteristic of Indian-accented speakers, we use INDIC PARLER TTS (Sankar et al., 2025), a state-of-the-art text-to-speech model fine-tuned for Indian languages. The model offers fine-grained prosodic control, enabling realistic variation in rhythm, intonation, and articulation.

By pairing the synthesized English speech with the normalized Indic transcripts from INDICVOICES, we construct a high-quality En-XX dataset. This resource enhances both linguistic diversity and domain robustness, supporting more effective training and evaluation for real-world speech translation applications.

3.4 Quality Control

To ensure the reliability of the aligned speech-text pairs, we adopt a two-fold evaluation strategy; one for audio-text alignment quality and the other for translation quality.

Alignment Quality: We quantify the accuracy of audio-transcript alignment by computing Levenshtein distance-based similarity scores between the cleaned reference transcripts and ASR-predicted outputs, following the methodology of Bhogale et al. (2022). This metric enables us to identify and filter out poorly aligned segments, thus maintaining high fidelity in timestamped pairs that are critical for training downstream speech and translation models.

Translation Quality: To evaluate the quality of the aligned parallel text, we compute cosine similarity scores between sentence embeddings of the source transcript and its corresponding translation.

These embeddings are obtained using the SONAR model (Duquenne et al., 2023), which provides robust cross-lingual sentence representations. *While the dataset is not filtered based on these metrics in the public release, we include both the alignment and translation quality scores alongside the data. This allows users to apply their own thresholds or filtering strategies depending on the specific needs of their applications.*


Comparison with Prior Web-Mined Corpora:

Figure 3 compares SONAR-based mining scores for $XX \rightarrow \text{En}$ sentence pairs from BHASAANUVAAD and the SeamlessAlign corpus. SeamlessAlign shows a concentration of low scores, indicating noisy or weak alignments, whereas BHASAANUVAAD yields a markedly higher proportion of high-quality pairs, with scores skewed toward the upper end of the distribution. Language-wise alignment score distributions are shown in Figures A4 and A5, with comprehensive plots for all languages provided in Figures A6, A7, and A8. These results underscore the effectiveness of our mining pipeline in producing cleaner and more semantically aligned sentence pairs.

4 Speech Translation Models

In this study, we primarily consider the Unified (or end-to-end) and the Cascaded Speech translation paradigms across various models. While recent advancements in Audio-LLMs (or SLLMs) introduce new paradigms by integrating Speech Foundation Models with LLMs for Automatic Speech Translation (AST), we do not include them in this study due to the lack of support for Indian languages.

4.1 Our Model

SEAMLESSM4T (SD_{BA})  We fine-tune the SEAMLESSM4T-V2-LARGE model on a filtered subset of BHASAANUVAAD (SD_{BA}; IndicSeamless), using a high-quality filtered subset to ensure effective adaptation to Indian languages. Training is conducted on $8 \times \text{A100 GPUs}$ for 430,000 steps with an effective per-device batch size of 64. We employ the Adam 8-bit optimizer with a learning rate of 1×10^{-5} and 1,000 warmup steps. To mitigate overfitting, we apply early stopping with a patience of 10.


4.1.1 Data


We construct the seed dataset by retaining only those audio-text pairs that meet the following quality thresholds: a cosine similarity-based mining

score $\sigma \geq 0.6$, and a Levenshtein distance-based alignment score $\tau \geq 0.8$, following Bhogale et al. (2022) for the alignment criterion, from the newly collected sources. These thresholds are chosen to ensure that only highly reliable and accurately aligned data contribute to model fine-tuning.


From the filtered data, we construct the evaluation set—BHASAANUVAAD-TEST (Table A3)—by randomly sampling approximately 20 minutes of speech per language for each translation direction ($\text{En} \rightarrow \text{XX}$ and $\text{XX} \rightarrow \text{En}$). This setup provides a balanced test set across all covered languages, while preserving diversity in speaker accents, linguistic complexity, and domain coverage. The remaining filtered data, excluding the test portion, is used for training the model.


4.2 Baselines

Unified Systems: In this paradigm, a single end-to-end speech translation model is used. The unified AST models considered in our study are as follows: **SEAMLESSM4T (SD)**  (Communication et al., 2023) is a foundational multilingual and multitask model supporting speech-to-speech (S2S), speech-to-text (S2T), text-to-speech (T2S), text-to-text (T2T) translation, and automatic speech recognition (ASR) for up to 100 languages. For our experiments, we use the SEAMLESS-M4T-V2-LARGE MODEL.

AZURE  is a closed-source speech translation system developed by Microsoft, providing cloud-based automatic speech recognition (ASR) and speech-to-text (S2T) translation services. While details about its architecture and training data remain proprietary, Azure’s speech translation API supports multiple languages and is widely used in commercial applications. For our study, we evaluate its performance against other baselines and SD_{BA}.

Cascaded Systems: In this paradigm, separately trained Automatic Speech Recognition (ASR) and Machine Translation (MT) models are combined in a pipeline for performing Automatic Speech Translation (AST). The different ASR and NMT models used in our experiments are listed below.

SEAMLESSM4T(SC)  (Communication et al., 2023) where we integrate the ASR and NMT capabilities of the SEAMLESSM4T model in a sequential pipeline.

WHISPER (W)  (Radford et al., 2023) is a multilingual model trained for automatic speech recognition (ASR) and speech translation (AST). For

lang	Direct						Cascaded					
	SD		SD _{BA}		AZURE		SC		W + IT2		SC + IT2	
	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA
asm	37.90	32.60	39.30	42.20	37.90	36.40	36.20	35.30	40.80	43.40	40.50	40.00
ben	47.50	31.00	48.90	43.10	47.00	36.30	46.40	36.90	51.00	42.50	50.40	38.80
guj	48.90	33.80	50.10	46.00	48.60	38.20	48.40	38.50	52.10	45.80	51.70	40.70
hin	54.20	25.00	55.10	51.90	53.10	45.40	54.20	45.60	56.00	49.90	55.70	44.80
kan	48.90	38.00	49.90	50.30	48.90	41.40	47.80	43.40	52.10	50.60	51.40	45.60
mal	49.00	39.40	51.20	50.80	50.30	43.80	47.90	44.20	54.70	49.60	53.90	47.60
mar	44.10	29.30	45.70	43.20	43.30	35.50	42.30	36.70	47.50	43.40	47.20	39.60
npi	50.00	28.00	50.40	44.40	46.80	26.70	47.20	35.60	53.70	42.90	53.10	36.30
ory	46.20	31.90	48.50	40.00	45.70	29.70	43.40	35.60	46.80	37.30	46.50	35.20
pan	50.00	27.00	51.40	34.40	48.10	28.50	48.80	28.40	51.50	34.40	50.90	29.80
snd	47.10	—	48.90	—	—	—	45.90	—	42.30	—	42.00	—
tam	50.60	33.80	52.30	40.50	50.50	42.80	49.50	41.20	53.30	48.80	53.00	44.80
tel	52.00	33.20	53.50	47.80	49.90	36.60	51.30	40.10	55.30	46.20	54.50	42.20
urd	47.20	41.10	47.30	49.70	45.80	41.10	46.60	45.50	48.80	49.60	48.80	46.10
avg.	48.10	33.20	49.46	44.75	47.38	37.11	46.85	39.33	50.42	44.84	49.97	40.90

Table 1: chrF++ (\uparrow) scores for different direct and cascaded speech translation models evaluated on the FLEURS-TEST (FL; green) and BHASAA NUVAAD-TEST (BA; yellow) datasets in $En \rightarrow XX$ direction. The colors represent a heatmap, where darker shades indicating better translation performance.

our experiments, we use the WHISPER-LARGE-V3 model which has been trained on 1 million hours of weakly labeled and 4 million hours of pseudo-labeled audio.

INDIC-CONFORMER (IC)⁴ (Bhogale et al., 2025) which is a suite of ASR models based on the Conformer architecture, models supporting all 22 scheduled Indian languages.

INDICTRANS2 (IT2) (Gala et al., 2023) which is an open-source transformer-based multilingual NMT model supporting translations across all 22 scheduled Indic languages. For our experiments, we use the INDICTRANS2-1B model.

In our experiments, we combine the INDICTRANS2-1B model and the SEAMLESS-M4T V2 LARGE translation model with various ASR models to perform automatic speech translation (AST). Specifically, the translation capabilities of INDICTRANS2-1B and SEAMLESS-M4T V2 LARGE (T2TT) are integrated with ASR outputs from systems such as WHISPER, and INDIC-CONFORMER.

4.3 Metrics

In this work, we adopt chrF++ (Popović, 2017) as our primary metric, which offers strong alignment with human judgments (Sai B et al., 2023) particularly Indian languages. We compute chrF++ scores using sacreBLEU (Post, 2018); for Indic-En⁵ evalu-

ation, we use the standard mteval-v13a tokenizer, while for En-Indic, we apply Indic-specific tokenizers from IndicNLP (Kunchukuttan, 2020) and Urduhack⁶ before scoring, ensuring linguistically informed segmentation. We also report additional metrics such as BLEU (Papineni et al., 2002) and XCOMET (Guerreiro et al., 2024) scores in Appendix Tables A4 & A5 and A6 & A7 respectively for both the directions.

5 Results

We evaluate the translation performance of various models in both $En \rightarrow XX$ and $XX \rightarrow En$ directions using chrF++ scores on the FLEURS-TEST and BHASAA NUVAAD-TEST (Table A3) datasets. Our analysis focuses on the comparative effectiveness of direct and cascaded models, highlighting the impact of fine-tuning and architectural choices.

5.1 Direct vs. Cascaded Systems

A comparison between direct and cascaded approaches on the FLEURS-TEST set reveals that their performance remains closely matched across most languages. In $En \rightarrow XX$ translation, the average chrF++ score difference between the two paradigms is marginal, suggesting that direct speech translation methods are increasingly closing the performance gap with cascaded pipelines. To analyze this further, we compare the strongest

⁵En-Indic sacreBLEU ChrF++

signature:nrefs:1|case:mixed|eff:yes|nc:6|nw:2|space:no|version:2.3.1

⁶<https://github.com/urduhack/urduhack>

lang	Direct						Cascaded					
	SD		SD _{BA}		AZURE		SC		IC + IT2		SC + IT2	
	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA
asm	44.00	46.90	49.70	53.80	38.20	41.20	46.30	49.40	47.20	51.50	47.10	50.70
ben	48.50	45.20	55.00	57.40	46.10	40.50	51.20	47.20	52.50	49.50	53.30	48.80
guj	52.80	51.20	60.20	60.00	45.40	32.70	55.90	54.20	57.10	54.50	57.50	54.40
hin	50.60	48.60	59.10	59.40	51.40	47.50	53.20	51.20	56.50	55.80	56.50	53.70
kan	47.60	47.60	54.80	56.90	47.10	38.10	50.10	50.80	52.00	52.50	52.00	52.20
mal	48.40	43.50	55.90	54.00	46.00	35.60	52.00	46.30	54.20	49.50	54.10	48.70
mar	48.40	46.20	55.30	56.20	43.80	37.70	52.20	50.10	54.30	51.10	53.50	50.20
npi	50.60	35.70	56.60	51.90	50.10	31.20	54.30	38.70	55.30	46.50	56.40	40.90
ory	48.60	50.10	54.90	57.70	48.20	44.00	52.80	53.50	55.10	52.10	54.20	52.90
pan	50.30	50.00	56.40	57.60	44.00	44.70	53.00	55.00	55.10	53.60	54.70	55.40
snd	31.10	—	43.40	—	—	—	36.20	—	39.30	—	38.50	—
tam	44.80	41.60	51.00	52.20	45.40	37.30	48.10	44.30	50.70	47.30	50.30	46.80
tel	47.50	44.50	54.30	55.00	48.30	39.40	51.50	47.50	55.60	51.20	54.40	50.10
urd	47.60	48.90	54.40	61.10	47.60	42.80	50.30	51.20	51.30	52.70	53.10	53.00
avg.	47.20	46.15	54.36	56.40	46.28	39.44	50.51	49.18	52.59	51.37	52.54	50.60

Table 2: chrF++ (\uparrow) scores for different direct and cascaded speech translation models evaluated on the FLEURS-TEST (FL; green) and BHASAA NUVAAD-TEST (BA; yellow) datasets in $XX \rightarrow En$ direction. The colors represent a heatmap, where darker shades indicating better translation performance.

direct model, SD_{BA}, with the best-performing cascaded setups: W + IT2 for $En \rightarrow XX$ and IC + IT2 for $XX \rightarrow En$.

For $En \rightarrow XX$ (Table 1), SD_{BA} and W + IT2 achieve comparable results across most languages, underscoring the efficacy of both direct speech translation and cascaded modeling. However, an exception arises in Tamil in BHASAA NUVAAD-TEST, where W + IT2 significantly outperforms SD_{BA}. This deviation warrants further investigation into language-specific factors that may influence direct model performance.

In the $XX \rightarrow En$ direction, SD_{BA} surpasses IC + IT2 across multiple languages, demonstrating strong generalization capabilities. Notably, on the BHASAA NUVAAD-TEST set, SD_{BA} achieves significantly higher CHR F++ scores for Nepali, Bengali, and Urdu, reinforcing its ability to handle diverse linguistic structures effectively. These results suggest that fine-tuned direct models can be competitive or even superior to cascaded systems, where cascading errors from ASR and MT modules compound and degrade translation quality.

5.2 Performance Gains from Fine-Tuning

Fine-tuning on the BhasaAnuvaad dataset (BA) leads to consistent improvements across both translation directions ($En \rightarrow XX$ and $XX \rightarrow En$). SD_{BA} achieves substantial gains over SD, particularly on the BA test set, where it outperforms the baseline by at least 10 points on average.

On the FLEURS-TEST set, SD_{BA} consistently outperforms SD across all languages, achieving an average improvement of over 7 points. These results highlight the effectiveness of domain-specific fine-tuning in enhancing direct speech translation performance, particularly in real-world, low-resource conditions.

6 Conclusion

In this work, we present BHASAA NUVAAD, the largest publicly available Indian language speech translation dataset, comprising 44,000+ hours of speech and 17 million aligned text segments across 14 Indian languages and English. Our dataset is built using a threefold approach: (i) aggregating high-quality sources, (ii) large-scale web crawling, and (iii) synthetic augmentation to capture real-world disfluencies. Using BHASAA NUVAAD, we train INDICSEAMLESS, a state-of-the-art speech translation model for Indian languages. Our evaluations show that fine-tuned direct speech translation models (SD_{BA}) match or surpass strong cascaded systems, particularly in low-resource settings where cascading errors degrade performance. By releasing all data, models, and code under open-source licenses, we aim to advance research in Indian language speech translation.

7 Limitations

Despite the scale and diversity of BHASAA NUVAAD, several limitations remain, addressing

which will be essential to making speech translation systems for Indian languages more robust, reliable, and truly representative of real-world usage.

First, while our dataset integrates a range of sources, spontaneous speech remains underrepresented. Real-world conversational speech often includes disfluencies, hesitations, code-switching, and speaker variability that are not fully captured in our current dataset.

Second, domain coverage remains a challenge. While BHASANUVAAD aggregates speech from multiple sources including educational content, government archives, and podcasts certain critical domains such as informal dialogues, social media discourse, medical interactions, and legal proceedings remain underrepresented.

Third, benchmarking Indian language AST remains an open challenge. Most prior evaluations rely on FLEURS, which consists of acted-out, read speech rather than spontaneous speech. While useful for measuring system performance under controlled conditions, FLEURS fails to reflect the complexities of natural, conversational speech, where hesitations, overlaps, and real-world noise significantly impact translation quality. Our experiments indicate that models achieving strong results on FLEURS do not necessarily maintain similar performance on more challenging, real-world speech scenarios.

8 Ethics

The code and datasets created in this work will be made available under permissible licenses. Generative AI systems were only used for assistance purely with the language of the paper, eg: paraphrasing, spell-check, polishing the author’s original content and for writing boiler plate code.

9 Acknowledgements

We gratefully acknowledge the generous support and funding provided by Digital India Bhashini, the Centre for Development of Advanced Computing (C-DAC) Pune, Yotta, EkStep Foundation, and Nilekani Philanthropies. We also thank Pranjal Agadh Chitale for his insightful discussions and valuable feedback throughout the development of this work.

References

- Arun Babu, Changan Wang, Andros Tjandra, Kushal Lakhota, Qiantong Xu, Naman Goyal, Kritika Singh, Patrick von Platen, Yatharth Saraf, Juan Pino, Alexei Baevski, Alexis Conneau, and Michael Auli. 2021. Xls-r: Self-supervised cross-lingual speech representation learning at scale. *arXiv preprint arXiv: 2111.09296*.
- Parnia Bahar, Patrick Wilken, Tamer Alkhouli, Andreas Guta, Pavel Golik, Evgeny Matusov, and Christian Herold. 2020. [Start-before-end and end-to-end: Neural speech translation by AppTek and RWTH Aachen University](#). In *Proceedings of the 17th International Conference on Spoken Language Translation*, pages 44–54, Online. Association for Computational Linguistics.
- Bismarck Bamfo Odoom, Nathaniel Robinson, Elijah Rippeth, Luis Tavarez-Arce, Kenton Murray, Matthew Wiesner, Paul McNamee, Philipp Koehn, and Kevin Duh. 2024. [Can synthetic speech improve end-to-end conversational speech translation?](#) In *Proceedings of the 16th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 167–177, Chicago, USA. Association for Machine Translation in the Americas.
- Ankur Bapna, Colin Cherry, Yu Zhang, Ye Jia, Melvin Johnson, Yong Cheng, Simran Khanuja, Jason Riesa, and Alexis Conneau. 2022. [mslm: Massively multilingual joint pre-training for speech and text](#). *arXiv preprint arXiv: 2202.01374*.
- Daniel Beck, Trevor Cohn, and Gholamreza Haffari. 2019. [Neural speech translation using lattice transformations and graph networks](#). In *Proceedings of the Thirteenth Workshop on Graph-Based Methods for Natural Language Processing (TextGraphs-13)*, pages 26–31, Hong Kong. Association for Computational Linguistics.
- Luisa Bentivogli, Mauro Cettolo, Marco Gaido, Alina Karakanta, Alberto Martinelli, Matteo Negri, and Marco Turchi. 2021. [Cascade versus direct speech translation: Do the differences still make a difference?](#) *arXiv preprint arXiv:2106.01045*.
- Kaushal Santosh Bhogale, Deovrat Mehendale, Tahir Javed, Devbrat Anuragi, Sakshi Joshi, Sai Sundaresan, Aparna Ananthanarayanan, Sharmistha Dey, Sathish Kumar Reddy G, Anusha Srinivasan, Abhigyan Raman, Pratyush Kumar, and Mitesh M. Khapra. 2025. [Towards bringing parity in pretraining datasets for low-resource indian languages](#). In *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5.
- Kaushal Santosh Bhogale, Abhigyan Raman, Tahir Javed, Sumanth Doddapaneni, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh M. Khapra. 2022. [Effectiveness of mining audio and text pairs from public data for improving asr systems for low-resource languages](#). *arXiv preprint*.

- Qiao Cheng, Meiyuan Fan, Yaqian Han, Jin Huang, and Yitao Duan. 2019. [Breaking the data barrier: Towards robust speech translation via adversarial stability training](#). In *Proceedings of the 16th International Conference on Spoken Language Translation*, Hong Kong. Association for Computational Linguistics.
- Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. 2023. [Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models](#). *arXiv preprint arXiv: 2311.07919*.
- Seamless Communication, Loïc Barrault, Yu-An Chung, Mariano Cora Meglioli, David Dale, Ning Dong, Paul-Ambroise Duquenne, Hady Elsahar, Hongyu Gong, Kevin Heffernan, John Hoffman, Christopher Klaiber, Pengwei Li, Daniel Licht, Jean Maillard, Alice Rakotoarison, Kaushik Ram Sadagopan, Guillaume Wenzek, Ethan Ye, Bapi Akula, Peng-Jen Chen, Naji El Hachem, Brian Ellis, Gabriel Mejia Gonzalez, Justin Haaheim, Prangthip Hansanti, Russ Howes, Bernie Huang, Min-Jae Hwang, Hirofumi Inaguma, Somya Jain, Elahe Kalbassi, Amanda Kallet, Iliia Kulikov, Janice Lam, Daniel Li, Xutai Ma, Ruslan Mavlyutov, Benjamin Peloquin, Mohamed Ramadan, Abinash Ramakrishnan, Anna Sun, Kevin Tran, Tuan Tran, Igor Tufanov, Vish Vogeti, Carleigh Wood, Yilin Yang, Bokai Yu, Pierre Andrews, Can Balioglu, Marta R. Costa-jussà, Onur Celebi, Maha Elbayad, Cynthia Gao, Francisco Guzmán, Justine Kao, Ann Lee, Alexandre Mourachko, Juan Pino, Sravya Popuri, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, Paden Tomasello, Changan Wang, Jeff Wang, and Skyler Wang. 2023. [SeamlessM4t: Massively multilingual & multimodal machine translation](#). *arXiv preprint arXiv: 2308.11596*.
- Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2022. [FLEURS: few-shot learning evaluation of universal representations of speech](#). In *IEEE Spoken Language Technology Workshop, SLT 2022, Doha, Qatar, January 9-12, 2023*, pages 798–805. IEEE.
- Matti Di Gangi, Robert Enyedi, Alessandra Brusadin, and Marcello Federico. 2019. [Robust neural machine translation for clean and noisy speech transcripts](#). In *Proceedings of the 16th International Conference on Spoken Language Translation*, Hong Kong. Association for Computational Linguistics.
- Paul-Ambroise Duquenne, Hongyu Gong, Ning Dong, Jingfei Du, Ann Lee, Vedanuj Goswami, Changan Wang, Juan Pino, Benoît Sagot, and Holger Schwenk. 2022. [Speechmatrix: A large-scale mined corpus of multilingual speech-to-speech translations](#). *ARXIV*.
- Paul-Ambroise Duquenne, Holger Schwenk, and Benoît Sagot. 2023. [SONAR: sentence-level multimodal and language-agnostic representations](#). *CoRR*, abs/2308.11466.
- Dominika Ďurišková, Daniela Jurášová, Matúš Žilinc, Eduard Šubert, and Ondřej Bojar. 2024. [Khan academy corpus: A multilingual corpus of khan academy lectures](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9743–9752, Torino, Italia. ELRA and ICCL.
- Yassir Fathullah, Chunyang Wu, Egor Lakomkin, Ke Li, Junteng Jia, Yuan Shanguan, Jay Mahadeokar, Ozlem Kalinli, Christian Fuegen, and Mike Seltzer. 2023. [Audiochatllama: Towards general-purpose speech abilities for llms](#). *arXiv preprint arXiv: 2311.06753*.
- Marco Gaido, Sara Papi, Matteo Negri, and Luisa Bentivogli. 2024. [Speech translation with speech foundation models and large language models: What is there and what is missing?](#) *Preprint*, arXiv:2402.12025.
- Jay P. Gala, Pranjal A. Chitale, AK Raghavan, Varun Gumma, Sumanth Doddapaneni, M. AswathKumar, J. Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M. Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. [Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages](#). *Trans. Mach. Learn. Res.*
- Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2024. [xcomet: Transparent machine translation evaluation through fine-grained error detection](#). *Transactions of the Association for Computational Linguistics*, 12:979–995.
- Anirudh Gupta, Neeraj Chhimwal, Ankur Dhuriya, Rishabh Gaur, Priyanshi Shah, Harveen Singh Chadha, and Vivek Raghavan. 2022. [indic-punct: An automatic punctuation restoration and inverse text normalization framework for indic languages](#). *Preprint*, arXiv:2203.16825.
- Zhichao Huang, Rong Ye, Tom Ko, Qianqian Dong, Shanbo Cheng, Mingxuan Wang, and Hang Li. 2023. [Speech translation with large language models: An industrial practice](#). *arXiv preprint arXiv: 2312.13585*.
- Hirofumi Inaguma, Yosuke Higuchi, Kevin Duh, Tatsuya Kawahara, and Shinji Watanabe. 2021. [Orthros: non-autoregressive end-to-end speech translation with dual-decoder](#). In *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7503–7507.
- Tahir Javed, Janki Nawale, Eldho George, Sakshi Joshi, Kaushal Bhogale, Deovrat Mehendale, Ishvinder Sethi, Aparna Ananthanarayanan, Hafsah Faquih, Pratiti Palit, Sneha Ravishankar, Saranya Sukumaran, Tripura Panchagnula, Sunjay Murali, Kunal Gandhi, Ambujavalli R, Manickam M, C Vajjayanthi, Krishnan Karunganni, Pratyush Kumar, and Mitesh Khapra. 2024. [IndicVoices: Towards building an inclusive multilingual speech dataset for Indian languages](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 10740–10782.

- Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Ye Jia, Michelle Tadmor Ramanovich, Quan Wang, and Heiga Zen. 2022. [CVSS corpus and massively multilingual speech-to-speech translation](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6691–6703, Marseille, France. European Language Resources Association.
- Philipp Koehn. 2005. [Europarl: A parallel corpus for statistical machine translation](#). In *Proceedings of Machine Translation Summit X: Papers*, pages 79–86, Phuket, Thailand.
- Anoop Kunchukuttan. 2020. The IndicNLP Library. https://github.com/anoopkunchukuttan/indic_nlp_library/blob/master/docs/indicnlp.pdf.
- Ann Lee, Peng-Jen Chen, Changan Wang, Jiatao Gu, Sravya Popuri, Xutai Ma, Adam Polyak, Yossi Adi, Qing He, Yun Tang, Juan Pino, and Wei-Ning Hsu. 2022. [Direct speech-to-speech translation with discrete units](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3327–3339, Dublin, Ireland. Association for Computational Linguistics.
- Shivam Mhaskar, Vineet Bhat, Akshay Batheja, Sourabh Deoghare, Paramveer Choudhary, and Pushpak Bhat-tacharyya. 2023. [Vakta-setu: A speech-to-speech machine translation service in select indic languages](#). *arXiv preprint arXiv: 2305.12518*.
- Vandan Mujadia and Dipti Sharma. 2023. [Towards speech to speech machine translation focusing on Indian languages](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 161–168, Dubrovnik, Croatia. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Maja Popović. 2017. [chr++: words helping character n-grams](#). In *Proceedings of the Second Conference on Machine Translation*, pages 612–618. Association for Computational Linguistics.
- Matt Post. 2018. [A call for clarity in reporting BLEU scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.
- Vineel Pratap, Andros Tjandra, Bowen Shi, Paden Tomasello, Arun Babu, Sayani Kundu, Ali Elkahky, Zhaoheng Ni, Apoorv Vyas, Maryam Fazel-Zarandi, et al. 2024. [Scaling speech technology to 1,000+](#) languages. *Journal of Machine Learning Research*, 25(97):1–52.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. [Robust speech recognition via large-scale weak supervision](#). In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 28492–28518. PMLR.
- Ananya Sai B, Tanay Dixit, Vignesh Nagarajan, Anoop Kunchukuttan, Pratyush Kumar, Mitesh M. Khapra, and Raj Dabre. 2023. [IndicMT eval: A dataset to meta-evaluate machine translation metrics for Indian languages](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14210–14228, Toronto, Canada. Association for Computational Linguistics.
- Elizabeth Salesky, Matthew Wiesner, Jacob Bremerman, R. Cattoni, Matteo Negri, Marco Turchi, Douglas W. Oard, and Matt Post. 2021. [The multilingual tedx corpus for speech recognition and translation](#). *INTERSPEECH*.
- Ashwin Sankar, Yoach Lacombe, Sherry Thomas, Praveen Srinivasa Varadhan, Sanchit Gandhi, and Mitesh M Khapra. 2025. [Rasmalai: Resources for adaptive speech modeling in indian languages with accents and intonations](#). In *Interspeech 2025*.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Improving neural machine translation models with monolingual data](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.
- Nivedita Sethiya and Chandresh Kumar Maurya. 2023. [End-to-end speech-to-text translation: A survey](#). *arXiv preprint arXiv: 2312.01053*.
- Nivedita Sethiya, Saanvi Nair, and Chandresh Maurya. 2024. [Indic-TEDST: Datasets and baselines for low-resource speech to text translation](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9019–9024, Torino, Italia. ELRA and ICCL.
- Bhavani Shankar, Preethi Jyothi, and Pushpak Bhat-tacharyya. 2024. [Costa: Code-switched speech translation using aligned speech-text interleaving](#). *arXiv preprint arXiv: 2406.10993*.
- Haiyue Song, Raj Dabre, Atsushi Fujita, and S. Kurohashi. 2019. [Coursera corpus mining and multistage fine-tuning for improving lectures translation](#). *International Conference on Language Resources and Evaluation*.
- Matthias Sperber and Matthias Paulik. 2020. [Speech translation and the end-to-end promise: Taking stock](#)

of where we are. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7409–7421, Online. Association for Computational Linguistics.

Silero Team. 2024. Silero vad: pre-trained enterprise-grade voice activity detector (vad), number detector and language classifier. <https://github.com/snakers4/silero-vad>.

Changhan Wang, Anne Wu, and Juan Pino. 2020. *Covost 2 and massively multilingual speech-to-text translation*. Preprint, arXiv:2007.10310.

Ron J. Weiss, Jan Chorowski, Navdeep Jaitly, Yonghui Wu, and Zhifeng Chen. 2017. *Sequence-to-sequence models can directly translate foreign speech*. In *Interspeech 2017*, pages 2625–2629.

Jian Wu, Yashesh Gaur, Zhuo Chen, Long Zhou, Yilun Zhu, Tianrui Wang, Jinyu Li, Shujie Liu, Bo Ren, Linqun Liu, and Yu Wu. 2023. *On decoder-only architecture for speech-to-text and large language model integration*. *Automatic Speech Recognition & Understanding*.

Rong Ye, Chengqi Zhao, Tom Ko, Chutong Meng, Tao Wang, Mingxuan Wang, and Jun Cao. 2023. *Gigast: A 10,000-hour pseudo speech translation corpus*. Preprint, arXiv:2204.03939.

Appendix

A Data Statistics

This section outlines the datasets used for training and evaluation, covering diverse Indic languages and both En-Indic and Indic-En directions. We begin with the test set composition, followed by training data statistics.

A.1 BHASAANUVAAD-TEST

The BHASAANUVAAD-TEST set (Table A3) contains 26.22 hours of En-Indic and 28.18 hours of Indic-En speech translation data, covering 13 Indian languages. Most languages have balanced data in both directions, with minor variations. Notably, mni and npj have smaller amounts of test data, with mni lacking Indic-En coverage. This test set provides a reliable benchmark for evaluating spoken translation performance across diverse Indic languages.

A.2 Training Data

En-Indic: The dataset combines 202.9 hours (91K utterances) of existing En-Indic speech data, 3453.8 hours (2619.3K utterances) of mined data, and 2694.1 hours (269K utterances) of synthetic data. Mined sources like NPTEL contribute the

largest share with 15940.9 hours and 6344.7K utterances, while synthetic data from IndicVoices helps balance low-resource languages. This diverse collection ensures comprehensive coverage for spoken translation across multiple Indic languages.

Indic-En The dataset includes 108 hours (90K utterances) of existing Indic-En speech data, 1076.7 hours (308.6K utterances) of mined data, and 3988.2 hours (288K utterances) of synthetic data. Among the mined sources, NPTEL contributes the largest portion with 6626.5 hours and 2603.3K utterances. The synthetic IndicVoices dataset further strengthens the collection with 3988.2 hours and 288K utterances, ensuring better representation across low-resource Indic languages. This combination provides a balanced and diverse dataset for advancing Indic-English spoken translation.

B Prompt Design for Punctuation Restoration

To address the lack of punctuation in automatic speech recognition (ASR) outputs and enhance downstream readability, we design a structured prompt for **punctuation restoration** using LLAMA-3.1-405B-INSTRUCT. As shown in Figure A1, the prompt guides the model to insert only punctuation marks into a given raw text while *strictly preserving the original word order and structure*.

The prompt emphasizes four key principles:

1. **Accuracy:** Punctuation must conform to the grammatical rules of the target language.
2. **Readability:** Insert appropriate marks (e.g., commas, periods, question marks) to enhance sentence clarity.
3. **Consistency:** Follow the punctuation style observed in reference texts or previous examples.
4. **Structure Preservation:** Do not add, delete, or reorder words—only punctuation should be inserted.

To support multilingual usage, we provide language-specific reference metadata, including the expected sentence terminator. The model is instructed to return output in a clean JSON format, making it suitable for direct integration into data processing pipelines.

lang	Existing				Mined								Synthetic			
	Indic-TedST		Fleurs		VaaniPedia		WordProject		UGCE Resources		Mann Ki Baat		NPTEL		IndicVoices	
	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.
asm	-	-	12.0	1.5	3.4	1.6	-	-	-	-	28.3	15.7	10.6	5.3	348.9	23.5
ben	13.1	6.5	12.0	1.5	121.0	52.5	50.7	23.9	431.7	329.6	38.3	21.4	1516.6	618.8	300.5	30.3
guj	17.2	1.5	12.0	1.5	121.0	52.5	50.7	23.9	431.7	329.4	39.1	22.0	1977.9	784.9	16.7	2.8
hin	100.4	52.5	12.0	1.5	121.0	52.5	50.7	23.9	431.7	329.2	37.7	21.3	3021.6	1176.8	206.6	20.4
kan	3.7	1.9	12.0	1.5	23.2	9.9	50.7	23.9	431.7	328.6	36.9	20.8	1601.4	638.0	175.6	11.7
mal	9.2	3.2	12.0	1.5	0.5	52.5	50.7	23.9	431.7	329.6	38.3	21.3	1729.1	699.3	278.1	22.9
mar	29.0	12.7	12.0	1.5	56.5	24.4	50.7	23.9	431.7	322.4	39.0	22.0	1601.7	626.3	138.0	16.2
npi	-	-	12.0	1.5	121.0	52.5	-	-	-	-	-	-	-	-	343.9	31.6
ory	-	-	12.0	1.5	79.0	34.0	49.2	23.1	-	-	36.8	20.5	0.8	0.3	161.1	19.1
pan	1.2	0.4	12.0	1.5	0.5	0.3	50.7	23.9	-	-	37.0	20.5	9.6	5.8	194.1	18.9
snd	-	-	12.0	1.5	121.0	52.5	-	-	-	-	-	-	-	-	-	-
tam	20.1	10.1	12.0	1.5	121.0	52.5	30.5	14.1	431.7	327.9	37.4	20.9	2631.9	1048.1	378.6	33.8
tel	9.1	2.2	12.0	1.5	121.0	52.5	50.7	23.9	431.7	322.6	37.8	21.4	1839.9	741.2	314.3	21.6
urd	-	-	12.0	1.5	2.3	1.1	50.7	23.9	-	-	29.0	16.0	-	-	186.6	19.1
Total	202.9	91.0	168.0	21	891.1	439.0	536.2	252.3	3453.8	2619.3	435.7	244.0	15940.9	6344.7	2694.1	269.0

Table A1: Overview of En-Indic spoken translation dataset across various Indic languages. The table details the hours of audio and the number of utterances collected from multiple sources. Number of utterances are in thousands (K)

lang	Existing				Mined								Synthetic			
	CVSS & Khan Academy		SeamlessAlign		Fleurs		WordProject		Spokentutorial		Mann Ki Baat		NPTEL		IndicVoices	
	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.	# hours	# uttr.
asm	-	-	-	-	12.0	1.5	-	-	36.9	23.2	23.0	13.9	9.3	5.7	443.6	26.0
ben	-	-	-	-	12.0	1.5	117.6	28.9	49.9	31.1	38.5	20.7	388.4	149.1	450.6	33.2
guj	12.0	10.0	-	-	12.0	1.5	82.1	29.0	63.6	37.6	37.4	21.3	795.0	300.9	30.2	3.0
hin	-	-	2355.4	1021.5	12.0	1.5	88.5	29.5	63.6	37.9	40.9	21.3	1433.9	477.5	282.9	21.6
kan	-	-	174.2	68.6	12.0	1.5	108.7	29.5	25.0	15.4	38.5	21.1	569.9	219.3	199.5	12.2
mal	-	-	-	-	12.0	1.5	134.4	55.5	24.4	15.4	35.6	24.1	721.0	320.2	378.4	24.9
mar	-	-	-	-	12.0	1.5	98.8	29.2	68.9	41.3	42.1	22.7	738.7	289.7	201.7	17.5
npi	-	-	-	-	12.0	1.5	-	-	58.8	35.4	-	-	-	-	479.1	34.0
ory	-	-	-	-	12.0	1.5	97.1	28.5	6.5	3.5	34.0	20.0	-	-	219.6	20.0
pan	-	-	-	-	12.0	1.5	113.2	1.1	19.8	11.4	78.2	8.2	-	-	231.4	19.9
snd	-	-	-	-	12.0	1.5	-	-	-	-	-	-	-	-	-	-
tam	96.0	80.0	1222.4	478.8	12.0	1.5	51.5	18.9	69.1	39.9	41.2	19.9	1231.8	538.8	454.5	35.5
tel	-	-	873.3	329.5	12.0	1.5	90.4	29.2	12.5	6.5	35.3	26.3	738.4	302.1	377.5	22.8
urd	-	-	2837.1	1107.7	12.0	1.5	94.3	29.4	-	-	11.9	7.0	-	-	239.2	17.4
Total	108.0	90.0	7462.3	3006.0	156.0	21	1076.7	308.6	499.0	298.7	456.4	226.5	6626.5	2603.3	3988.2	288.0

Table A2: Overview of Indic-En spoken translation dataset across various Indic languages. The table details the hours of audio and the number of utterances collected from multiple sources. Number of utterances are in thousands (K)

<i>lang</i>	En-XX	XX-En
<i>asm</i>	1.51	2.01
<i>ben</i>	2.51	2.50
<i>guj</i>	2.51	2.51
<i>hin</i>	2.50	2.50
<i>kan</i>	2.51	2.50
<i>mal</i>	2.13	2.50
<i>mar</i>	2.51	2.51
<i>npi</i>	0.50	0.64
<i>ory</i>	1.50	2.01
<i>pan</i>	1.02	1.98
<i>tam</i>	2.51	2.51
<i>tel</i>	2.51	2.51
<i>urd</i>	1.50	1.50
<i>tot.</i>	26.22	28.18

Table A3: Statistics of data obtained across each language in both Indic-En and En-Indic directions for BHASAAANUVAAD-TEST set. For each language and in both directions, we show the number of hours.

This prompt is especially beneficial for restoring structure in spontaneous speech transcripts, facilitating improved performance in downstream tasks such as machine translation, summarization, and language modeling where well-punctuated input is essential.

C Prompts Used for Synthetic Data Generation:

Prompt for Normal Translation (Figure A2) is structured to generate neutral, grammatically correct translations faithful to the original content. It emphasizes three key principles: (1) accuracy, ensuring no information is lost, (2) natural tone, and (3) meaningful restructuring, which encourages fluent and idiomatic English rather than literal renderings. This prompt is used to simulate conventional text translation, making the output suitable for benchmark comparison with traditional MT systems.

In contrast, **Prompt for Colloquial Translation** (Figure A3) is tailored for generating synthetic translations that mimic spontaneous, conversational speech. It guides the model to incorporate disfluencies, fillers, contractions, and informal vocabulary, reflecting the characteristics of real-world spoken language. The prompt explicitly encourages natural phrasing and flexibility in sentence structure while preserving the core meaning. This helps in constructing datasets that better model speech patterns and are particularly valuable for

Prompt for Punctuation Restoration

You are an expert in inserting punctuation. Please help in adding punctuation to the following text while strictly preserving the original words and structure.

Enhance readability by inserting only punctuation marks. Do not modify, add, or remove any words.

Follow these guidelines:\n\n

1. **Accuracy:** Ensure punctuation is applied correctly based on the language's grammatical rules.
2. **Readability:** Improve sentence clarity by inserting appropriate punctuation marks (commas, periods, question marks, etc.).
3. **Consistency:** Follow the punctuation style observed in the provided reference text.
4. **Preservation of Structure:** Do not alter word order or introduce new elements—only punctuation should be adjusted.

Reference Information:
- Language of the text: {lang}
- Sentence terminator for {lang}: {terminator}

Output Format:
Provide only the punctuated text in JSON format with the structure:

```
json
{ "punctuated_text": "Your punctuated text here" }
```

Figure A1: Prompt For Punctuation Restoration

training or evaluating speech-to-speech translation systems in informal contexts.

D Analysis of Mining Score Distributions

A histogram depicting the distribution of mining scores for each language pair in the dataset is illustrated in Figure A4, Figure A5, Figure A6, Figure A7, and Figure A8. These visualizations highlight a significant challenge: the SONAR model's performance is considerably constrained when applied to **low-resource languages**.

For instance, the Assamese-English language pair (Figure A8a) exhibits a distribution of mining scores that is heavily skewed towards lower values, with a distinct peak frequency around a mere 0.2. Similarly, the English-Sindhi pair, as depicted in Figure A7e, also demonstrates this limitation, with its most prominent frequency peak for mining scores occurring around 0.2 to 0.3, despite some presence of higher scores. This performance with Sindhi, a low-resource language, strongly suggests that the SONAR model struggles to generate

Prompt for Normal Translation

You are machine translation model. Please help in translating the following English text to tgt_lang.

"Translate the following text from English to tgt_lang.

Follow these guidelines:\n\n"

"1. **Accuracy:** Translate the text as accurately as possible, ensuring no information is missed.\n"

"2. **Tone and Naturalness:** Maintain the original tone of the text and ensure the translation sounds natural in English.\n"

"3. **Meaningful Restructuring:** Do not translate word-for-word. Restructure sentences as needed to ensure clarity and fluency in English.\n"

"Give only the translation and nothing else. Give the output as a JSON with the structure "

Figure A2: Prompt For Normal Translation

Prompt for Colloquial Translation

You are colloquial machine translation model. Please help in translating the following English text to {tgt_lang}.

I want to generate some synthetic speech data. For this, I will give you a verbatim text in {tgt_lang}, and your job is to colloquially translate it to English while making it sound natural, as if it were spoken in casual conversation.

"The translation should maintain the informal tone and spoken quirks of the original Hindi, using contractions, fillers, or casual "

"phrasing that a native English speaker might use. However, please ensure that the meaning stays accurate to the original text.\n\n"

"Avoid making the English sound overly formal or stilted. The translated text should feel spontaneous and conversational, like real speech. "

"Just give the text in English and nothing else. Below are some additional suggestions: \n\n"

"1. Use contractions (e.g., 'I'm' instead of 'I am', 'you're' instead of 'you are') to make the text sound more conversational. \n"

"2. Include natural fillers (e.g., 'you know', 'like', 'well', 'I mean') where appropriate, to reflect spoken language.\n"

"3. Use colloquial vocabulary. Avoid formal or academic words; opt for simpler, everyday language.\n"

"4. Retain casual expressions that might exist in the original text, replacing them with equivalent English expressions (e.g., 'Yaar' can be translated as 'buddy' or 'man').\n"

"5. Rearrange sentence structure if needed to reflect how a native English speaker might say it naturally.\n"

"6. Ensure smooth flow between sentences, with transitions that mimic spontaneous speech (e.g., 'So', 'Anyway', 'By the way')."

Figure A3: Prompt For Colloquial Translation

consistently effective cross-lingual embeddings for under-represented languages.

In contrast, the model demonstrates commendable efficacy with higher-resource Indian languages. For language pairs like Bengali-English (Figure A8b), Gujarati-English (Figure A8c), Malayalam-English (Figure A8d), and Marathi-English (Figure A8e), the mining scores predominantly cluster at higher values, typically peaking between 0.6 and 0.8. This indicates the SONAR model's capability to capture semantic similarities effectively when ample linguistic resources are available.

These divergent outcomes underscore a critical need: current multilingual embedding models like SONAR are not robust enough for low-resource languages. To address this, two major steps are imperative: (1) the adoption of **language-specific thresholds** when filtering mined data, rather than applying uniform cutoffs across all language pairs, and (2) the development of **improved multilingual embedding models** that can provide more accurate and reliable cross-lingual representations, particularly for under-represented languages like Assamese and Sindhi. Without such targeted

improvements, data mining efforts risk perpetuating the under-representation and degraded performance of these languages in downstream tasks.

<i>lang</i>	Direct						Cascaded					
	SD		SD _{BA}		AZURE		SC		IC + IT2		SC + IT2	
	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA
<i>asm</i>	0.77	0.76	0.75	0.83	0.59	0.70	0.58	0.60	0.54	0.64	0.58	0.78
<i>ben</i>	0.86	0.82	0.83	0.88	0.74	0.68	0.80	0.78	0.77	0.75	0.80	0.82
<i>guj</i>	0.86	0.84	0.85	0.88	0.66	0.59	0.88	0.85	0.87	0.83	0.88	0.86
<i>hin</i>	0.87	0.85	0.85	0.88	0.80	0.78	0.73	0.77	0.71	0.70	0.73	0.85
<i>kan</i>	0.86	0.84	0.84	0.89	0.72	0.69	0.77	0.77	0.76	0.77	0.77	0.85
<i>mal</i>	0.86	0.80	0.84	0.87	0.71	0.66	0.86	0.77	0.86	0.80	0.87	0.81
<i>mar</i>	0.82	0.82	0.80	0.87	0.66	0.69	0.67	0.66	0.65	0.66	0.67	0.84
<i>npi</i>	0.83	0.65	0.81	0.83	0.66	0.55	0.75	0.61	0.71	0.67	0.75	0.67
<i>ory</i>	0.82	0.85	0.82	0.89	0.73	0.75	0.80	0.78	0.81	0.81	0.80	0.87
<i>pan</i>	0.79	0.71	0.80	0.79	0.59	0.68	0.77	0.72	0.78	0.74	0.77	0.72
<i>snd</i>	0.43	-	0.60	-	-	-	0.39	-	0.21	-	0.39	-
<i>tam</i>	0.81	0.74	0.78	0.85	0.71	0.67	0.69	0.67	0.69	0.68	0.69	0.77
<i>tel</i>	0.83	0.77	0.80	0.86	0.70	0.65	0.73	0.71	0.73	0.71	0.73	0.81
<i>urd</i>	0.82	0.82	0.82	0.86	0.70	0.67	0.73	0.77	0.67	0.70	0.73	0.83
<i>total</i>	0.80	0.79	0.80	0.86	0.69	0.67	0.73	0.73	0.70	0.73	0.73	0.81

Table A4: COMET (\uparrow) scores for different direct and cascaded speech translation models evaluated on the FLEURS-TEST (FL; green) and BHASAAANUVAAD-TEST (BA; yellow) datasets in $XX \rightarrow En$ direction. The colors represent a heatmap, where darker shades indicating better translation performance.

<i>lang</i>	Direct						Cascaded					
	SD		SD _{BA}		AZURE		SC		W + IT2		SC + IT2	
	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA
<i>asm</i>	0.59	0.49	0.56	0.62	0.49	0.50	0.42	0.31	0.42	0.35	0.42	0.56
<i>ben</i>	0.77	0.59	0.74	0.73	0.69	0.61	0.77	0.46	0.78	0.50	0.77	0.67
<i>guj</i>	0.86	0.65	0.86	0.79	0.80	0.65	0.79	0.46	0.80	0.52	0.79	0.71
<i>hin</i>	0.74	0.67	0.74	0.75	0.64	0.65	0.86	0.64	0.85	0.68	0.86	0.71
<i>kan</i>	0.72	0.57	0.72	0.72	0.66	0.60	0.70	0.45	0.69	0.50	0.70	0.64
<i>mal</i>	0.85	0.77	0.84	0.84	0.79	0.76	0.82	0.63	0.82	0.64	0.82	0.81
<i>mar</i>	0.66	0.48	0.64	0.62	0.60	0.51	0.72	0.42	0.71	0.46	0.72	0.58
<i>npi</i>	0.83	0.48	0.84	0.71	0.73	0.51	0.79	0.32	0.8	0.38	0.79	0.59
<i>ory</i>	0.83	0.58	0.80	0.73	0.74	0.56	0.66	0.33	0.67	0.37	0.66	0.66
<i>pan</i>	0.82	0.49	0.82	0.57	0.76	0.64	0.74	0.29	0.74	0.39	0.74	0.50
<i>snd</i>	0.69	-	0.72	-	-	-	0.65	-	0.66	0.32	0.65	-
<i>tam</i>	0.68	0.55	0.67	0.67	0.62	0.59	0.71	0.46	0.72	0.50	0.71	0.63
<i>tel</i>	0.75	0.56	0.74	0.71	0.71	0.62	0.73	0.45	0.74	0.49	0.73	0.62
<i>urd</i>	0.73	0.60	0.73	0.72	0.66	0.59	0.78	0.49	0.78	0.54	0.78	0.65
<i>total</i>	0.75	0.58	0.74	0.71	0.68	0.60	0.72	0.44	0.73	0.47	0.72	0.64

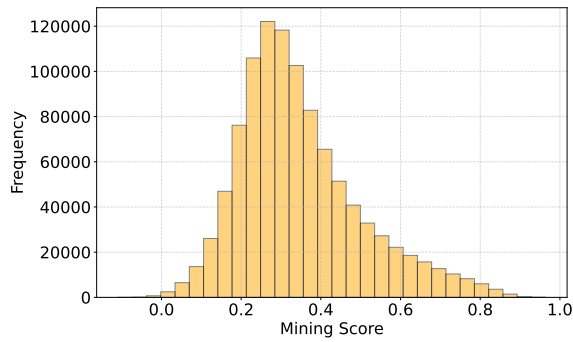
Table A5: COMET (\uparrow) scores for different direct and cascaded speech translation models evaluated on the FLEURS-TEST (FL; green) and BHASAAANUVAAD-TEST (BA; yellow) datasets in $En \rightarrow XX$ direction. The colors represent a heatmap, where darker shades indicating better translation performance.

lang	Direct						Cascaded					
	SD		SD _{BA}		AZURE		SC		W + IT2		SC + IT2	
	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA
asm	8.20	8.70	9.50	15.50	8.20	12.40	7.70	9.10	9.40	14.60	9.20	9.10
ben	14.80	8.30	17.20	16.30	15.80	11.90	14.20	10.90	18.00	14.80	17.60	10.90
guj	16.80	11.80	19.50	21.00	16.90	14.70	16.40	14.70	18.90	20.00	18.50	14.70
hin	28.70	11.50	31.60	29.50	29.40	25.20	28.70	23.30	30.30	27.50	30.30	23.30
kan	11.90	9.60	14.80	18.10	13.90	12.50	11.40	11.90	15.00	16.50	14.50	11.90
mal	9.50	8.40	13.00	17.80	12.50	13.50	9.60	11.30	14.20	16.10	13.00	11.30
mar	11.10	7.60	12.30	15.50	10.70	10.80	10.40	10.90	13.30	14.60	12.80	10.90
npi	14.10	6.90	16.00	15.00	12.70	5.80	13.00	9.70	17.3	13.30	16.90	9.70
ory	13.30	7.50	15.30	12.50	14.30	6.10	11.40	8.80	12.20	9.30	11.90	8.80
pan	24.00	11.60	26.10	17.30	22.90	10.20	22.80	12.60	25.70	16.20	25.20	12.60
snd	20.50	—	23.60	—	—	—	19.80	—	16.80	—	16.20	—
tam	18.30	6.30	14.90	16.60	12.90	13.80	11.00	9.80	13.60	16.60	13.30	9.80
tel	15.80	7.60	18.70	17.00	15.20	10.60	15.30	9.80	18.80	15.40	17.90	9.80
urd	21.80	19.70	23.00	24.90	22.10	20.50	21.20	22.50	23.60	27.30	23.90	22.50
total	16.34	9.65	18.25	18.23	15.96	12.92	15.21	12.72	17.65	17.09	17.23	12.72

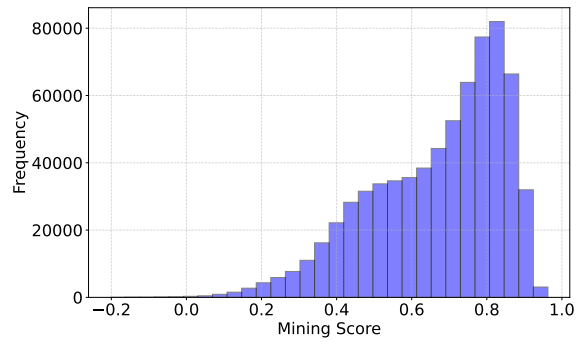
Table A6: BLEU (\uparrow) scores for different direct and cascaded speech translation models evaluated on the FLEURS-TEST (FL; green) and BHASANUVAAD-TEST (BA; yellow) datasets in $En \rightarrow XX$ direction. The colors represent a heatmap, where darker shades indicating better translation performance.

lang	Direct						Cascaded					
	SD		SD _{BA}		AZURE		SC		W + IT2		SC + IT2	
	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA	FL	BA
asm	14.70	21.40	22.40	29.50	12.80	16.70	17.00	23.60	17.90	26.40	16.90	23.60
ben	18.40	21.80	27.40	35.20	19.40	18.30	21.00	23.30	22.20	29.80	22.00	23.30
guj	22.80	27.20	33.50	37.70	19.40	12.60	25.50	29.50	27.70	34.30	26.70	29.50
hin	20.30	25.40	33.40	37.50	23.40	24.70	23.70	27.80	26.70	35.20	26.90	27.80
kan	17.30	21.10	27.60	32.70	17.50	15.40	19.80	25.70	21.80	31.30	20.50	25.70
mal	18.20	18.80	28.70	31.20	18.70	12.50	22.10	21.50	24.60	26.80	23.80	21.50
mar	18.30	21.20	28.30	33.40	17.80	14.10	21.80	24.70	24.50	29.90	22.00	24.70
npi	21.40	14.80	30.50	29.90	20.20	10.90	24.90	16.20	26.50	22.20	27.00	16.20
ory	18.30	25.20	27.20	34.70	19.90	19.90	23.20	29.70	25.90	33.40	24.00	29.70
pan	20.30	24.10	30.30	29.20	15.60	21.60	23.70	28.80	25.70	24.90	25.70	28.80
snd	6.30	—	15.30	—	—	—	8.00	—	12.70	—	10.30	—
tam	15.50	17.80	24.20	29.50	17.50	15.50	18.10	21.00	20.90	26.00	19.70	21.00
tel	18.10	19.90	27.40	31.30	19.60	15.80	22.80	22.90	26.90	29.10	25.10	22.90
urd	17.50	27.00	27.30	41.00	18.50	21.70	20.70	29.30	22.50	35.60	23.80	29.30
total	17.67	21.98	27.39	33.29	18.48	16.90	20.88	24.92	23.32	29.61	22.46	24.92

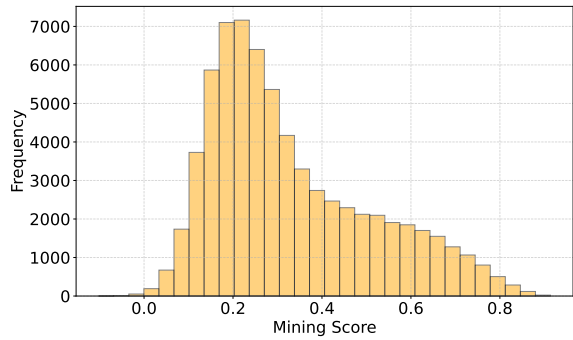
Table A7: BLEU (\uparrow) scores for different direct and cascaded speech translation models evaluated on the FLEURS-TEST (FL; green) and BHASANUVAAD-TEST (BA; yellow) datasets in $XX \rightarrow En$ direction. The colors represent a heatmap, where darker shades indicating better translation performance.



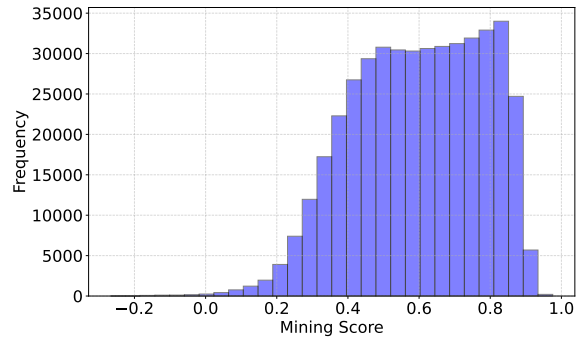
(a) Seamless Align - Hindi



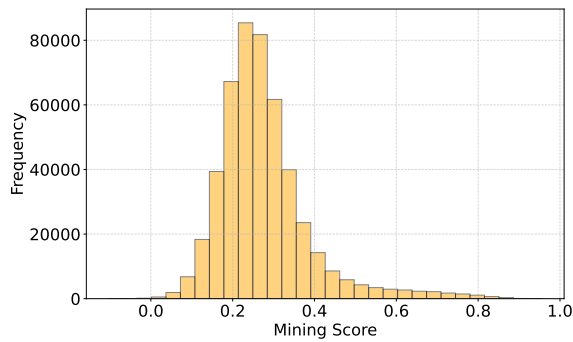
(b) BHASAANUVAAD- Hindi



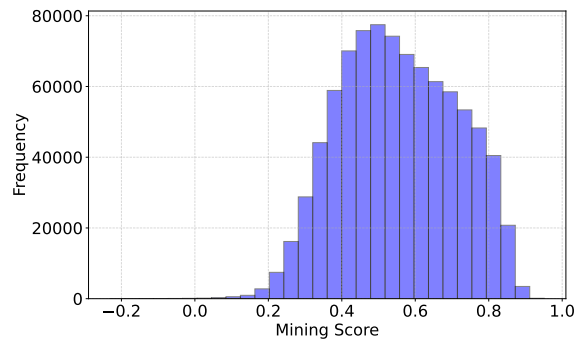
(c) Seamless Align - Kannada



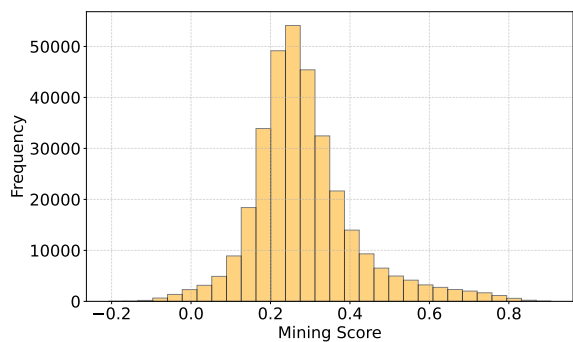
(d) BHASAANUVAAD- Kannada



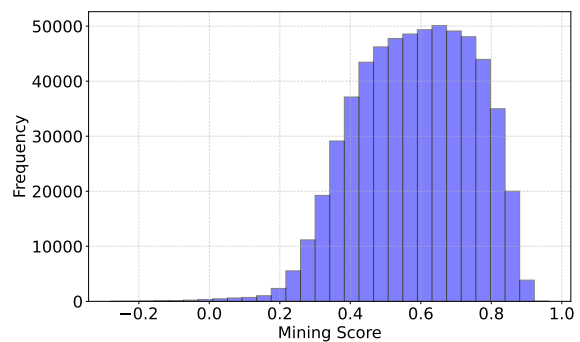
(e) Seamless Align - Tamil



(f) BHASAANUVAAD- Tamil

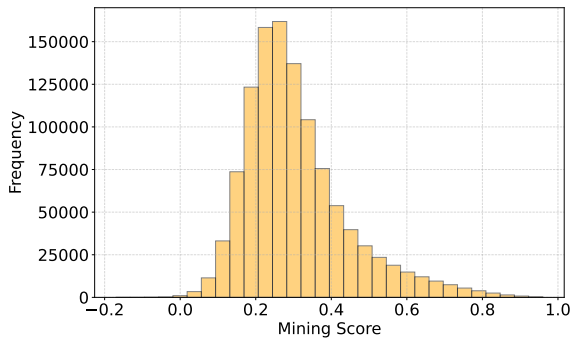


(g) Seamless Align - Telugu

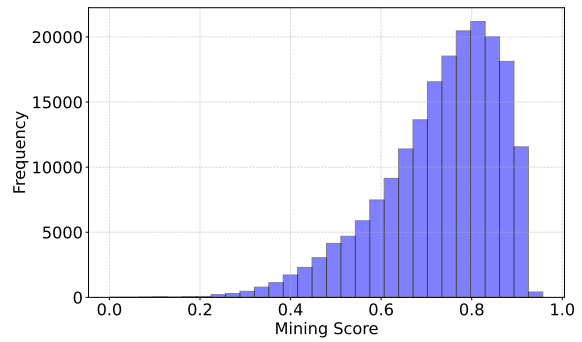


(h) BHASAANUVAAD- Telugu

Figure A4: Distribution of SONAR mining scores for $XX \rightarrow \text{En}$ sentence pairs from the SeamlessAlign (left) and BHASAANUVAAD (right) datasets across four languages: Hindi, Kannada, Tamil, and Telugu. BHASAANUVAAD consistently exhibits a higher concentration of high-quality alignments, with mining scores skewed toward the upper end of the scale. In contrast, SeamlessAlign distributions are centered around lower mining scores, indicating overall lower alignment quality. This comparison underscores the improved precision of alignments in BHASAANUVAAD.

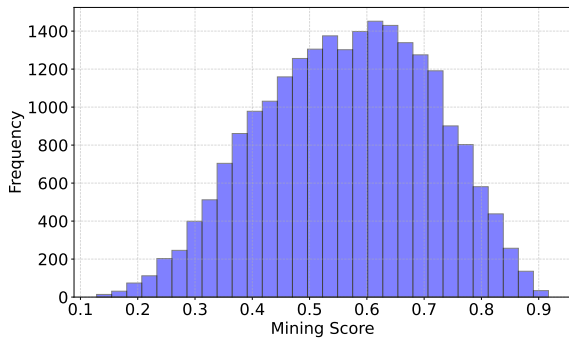


(a) Seamless Align - Urdu

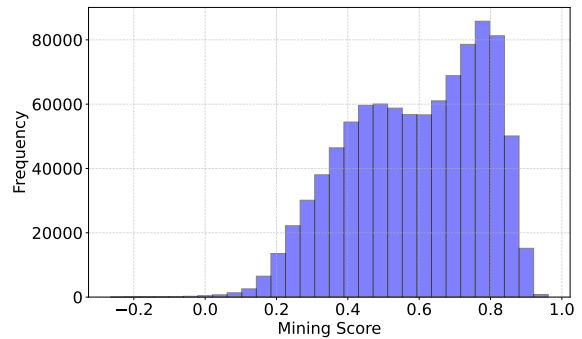


(b) BHASAANUVAAD - Urdu

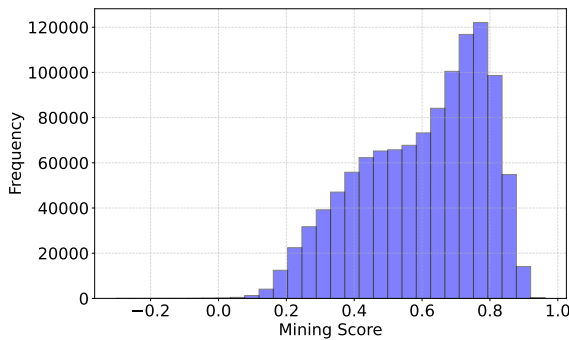
Figure A5: Distribution of SONAR mining scores for $XX \rightarrow \text{En}$ sentence pairs from the SeamlessAlign (left) and BHASAANUVAAD (right) datasets across four languages: Hindi, Kannada, Tamil, and Telugu. BHASAANUVAAD consistently exhibits a higher concentration of high-quality alignments, with mining scores skewed toward the upper end of the scale. In contrast, SeamlessAlign distributions are centered around lower mining scores, indicating overall lower alignment quality.



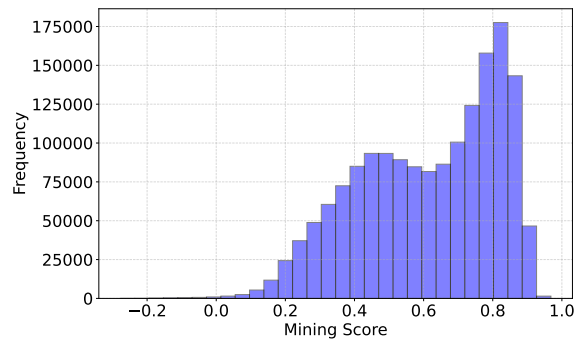
(a) English - Assamese



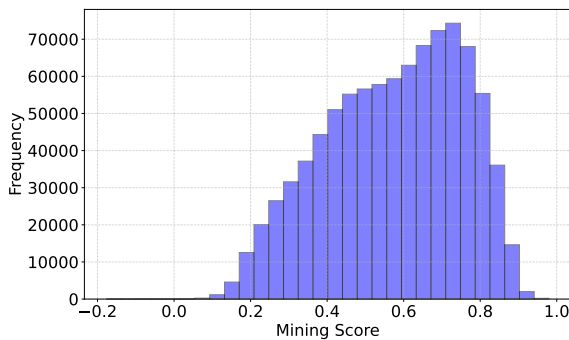
(b) English - Bengali



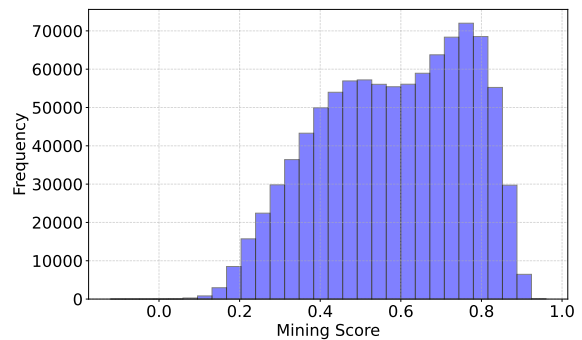
(c) English - Gujarati



(d) English - Hindi

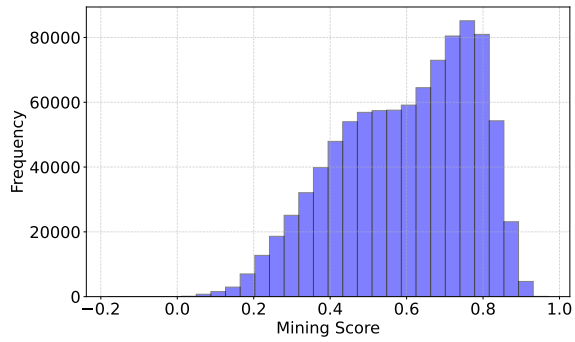


(e) English - Kannada

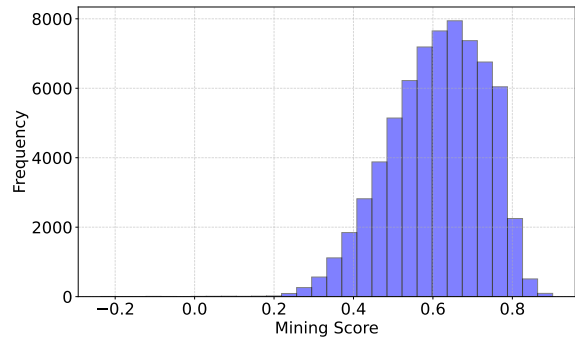


(f) English - Malayalam

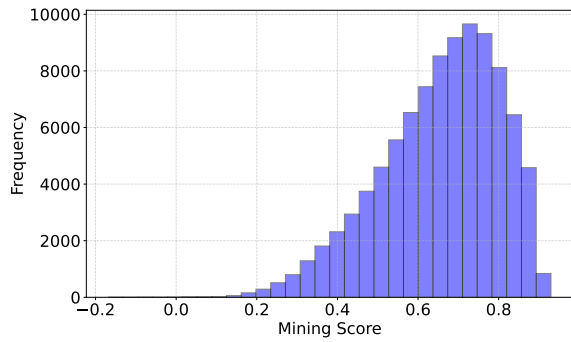
Figure A6: Histogram of SONAR mining scores for $\text{En} \rightarrow \text{XX}$ language pairs in the BHASAANUVAAD dataset. The distributions indicate a strong skew toward high-quality alignments (scores above 0.6) across most language pairs, reflecting the effectiveness of the mining and filtering pipeline used in BHASAANUVAAD.



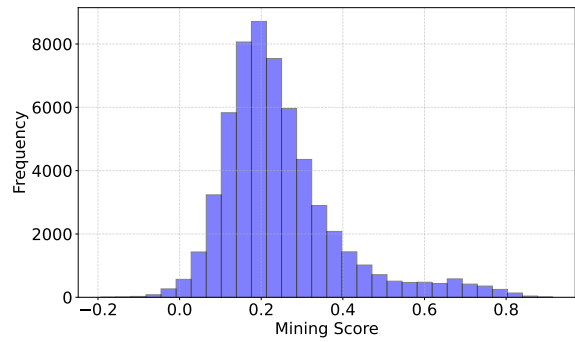
(a) English - Marathi



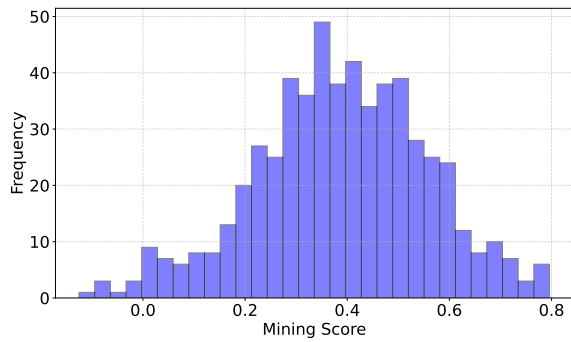
(b) English - Nepali



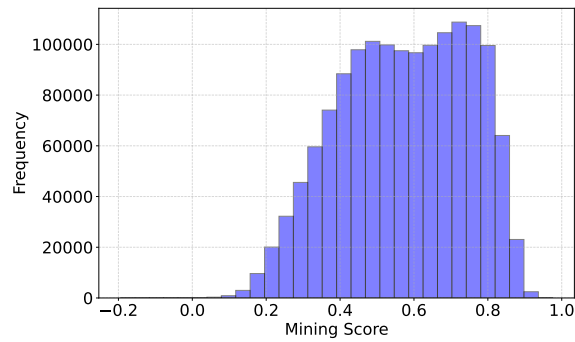
(c) English - Odia



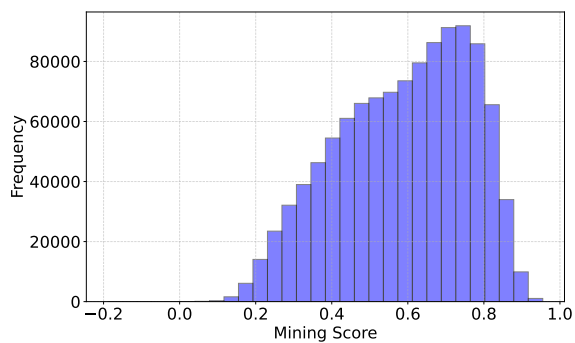
(d) English - Punjabi



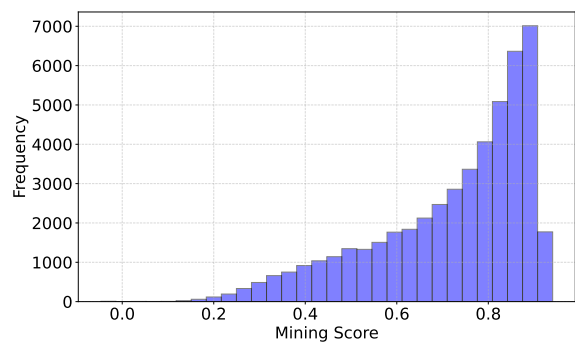
(e) English - Sindhi



(f) English - Tamil

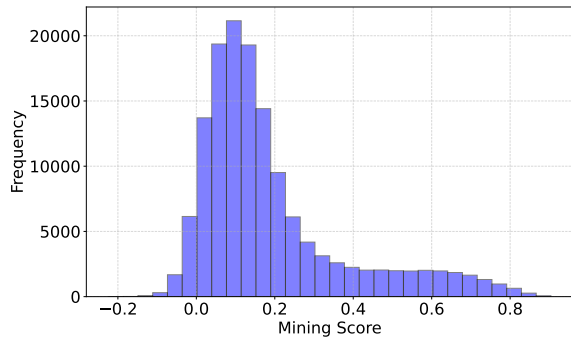


(g) English - Telugu

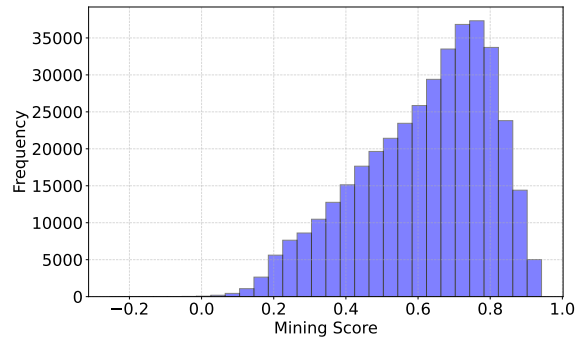


(h) English - Urdu

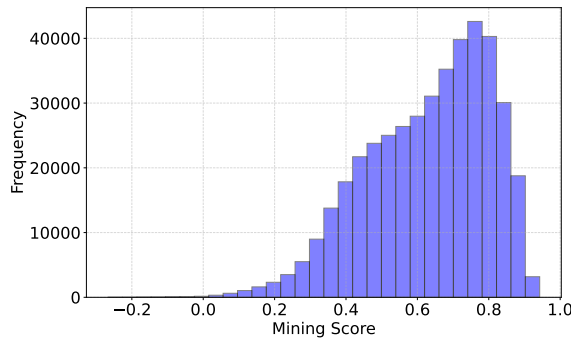
Figure A7: Histogram of SONAR mining scores for En \rightarrow XX language pairs in the BHASAANUVAAD dataset. The distributions indicate a strong skew toward high-quality alignments (scores above 0.6) across most language pairs, reflecting the effectiveness of the mining and filtering pipeline used in BHASAANUVAAD (contd).



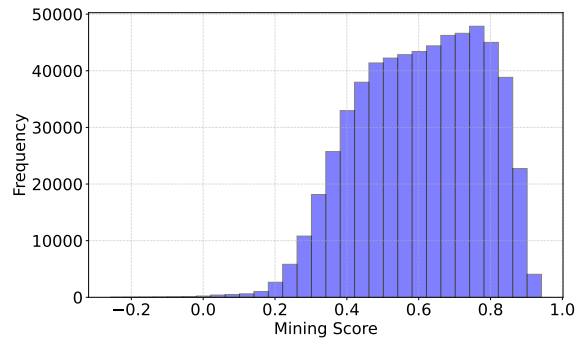
(a) Assamese - English



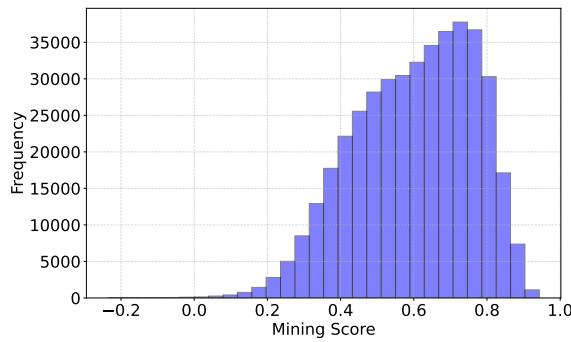
(b) Bengali - English



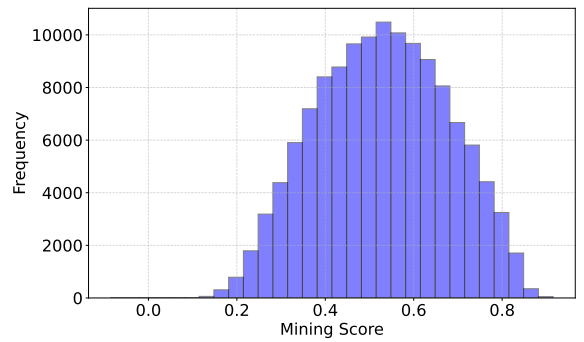
(c) Gujarati - English



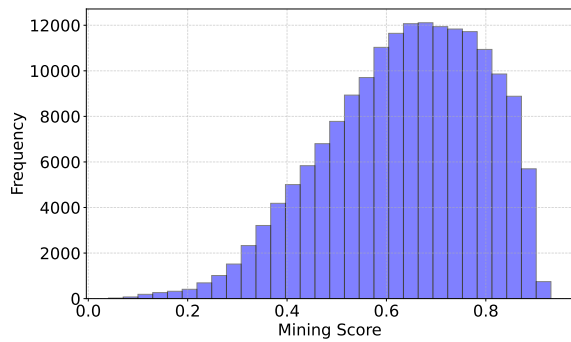
(d) Malayalam - English



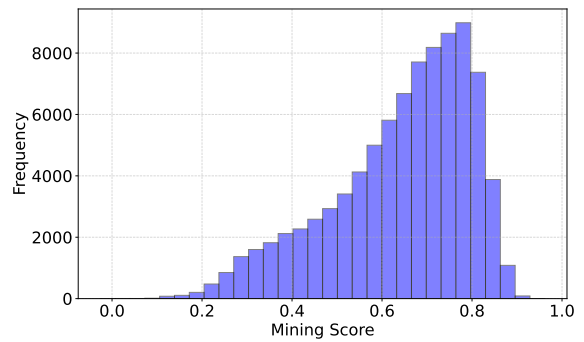
(e) Marathi - English



(f) Nepali - English



(g) Odia - English



(h) Punjabi - English

Figure A8: Histogram of SONAR mining scores for $XX \rightarrow \text{Any language pairs}$ in the BHASAANUVAAD dataset. The distributions indicate a strong skew toward high-quality alignments (scores above 0.6) across most language pairs, reflecting the effectiveness of the mining and filtering pipeline used in BHASAANUVAAD.