

# Are ASR Foundation Models Generalized Enough to Capture Regional Dialects in Low-Resource Languages?

Tawsif Tashwar Dipto<sup>1</sup>, Azmol Hossain<sup>3</sup>, Rubayet Sabbir Faruque<sup>2</sup>, Md. Rezuwan Hassan<sup>2</sup>  
Kanij Fatema<sup>2</sup>, Tanmoy Shome<sup>2</sup>, Ruwad Naswan<sup>3</sup>, Md. Foriduzzaman Zihad<sup>3</sup>, Mohaymen Ul Anam<sup>1</sup>  
Nazia Tasnim<sup>3</sup>, Hasan Mahmud<sup>1</sup>, Md. Kamrul Hasan<sup>1</sup>, Md. Mehedi Hasan Shawon<sup>2</sup>  
Farig Sadeque<sup>2</sup>, Tahsin Reasat<sup>3</sup>

<sup>1</sup>Islamic University of Technology    <sup>2</sup>Brac University    <sup>3</sup>Bengali.AI

## Bengali Dialectal ASR: Why It Matters

- Bengali has 270+ million speakers worldwide - 3rd most spoken native language in the world
- Yet current ASR systems (e.g., Whisper) are trained almost exclusively on Standard Colloquial Bengali (SCB)
- Real-life Bengali speech is highly dialectal which lives in a separate feature space than SCB
- Neglecting the dialectal variations in phonology, morphology, syntax and prosody characteristic of authentic speech results in a severe performance degradation for foundation models

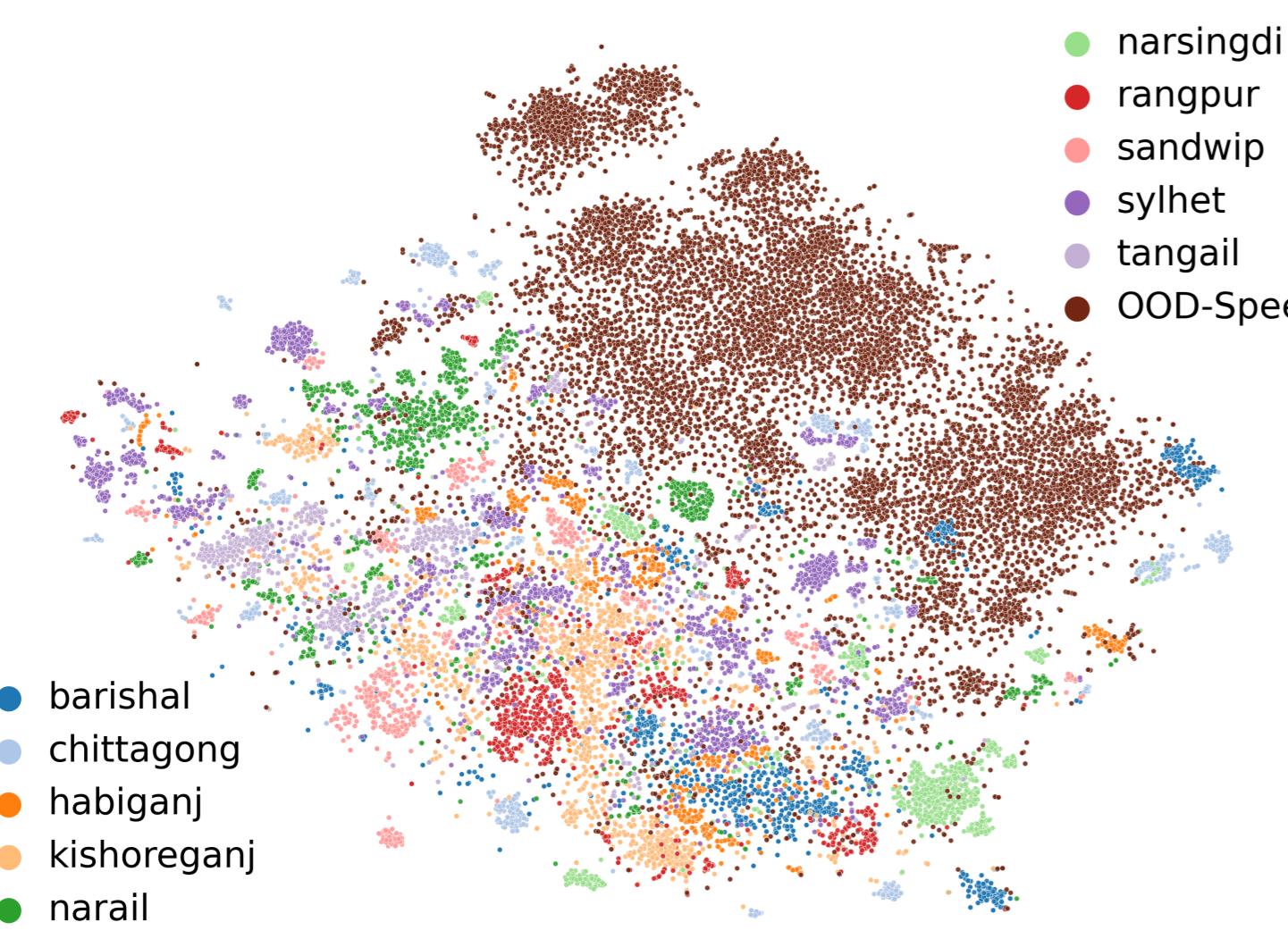


Figure 1: t-SNE of GeMAPS acoustic features shows a clear distribution shift between Standard SCB (brown cluster) and Regional Dialects.

## Linguistic Diversity

Our analysis reveals why standard models fail. Bengali dialects are not just “accents” They diverge in every linguistic dimension. We observe massive shifts from the standard dialect:

- Phonology: Aspirated shifts (e.g., Sylheti ক → খ) and heavy nasalization in Chittagong/Sandvip.
- Morphology: Unique clitics (Singular: টা→তা) and verb conjugations (খেতে পারে → খাইতারে).
- Syntax: Subject-Object-Verb deviations. Negation placement shifts (e.g., Chittagonian: যাইনা → ন যাই).
- Lexicon: High Out-of-Vocabulary (OOV) rates. Up to 62% of words in Sandvip are unique compared to SCB.

Region	Regional Sentence	Standard Bengali	English Translation
Rangpur	মোবাইল একনামে রেকর্ড করো।	[mobail' ekan'm rekor'd korbo]	I will record on mobile now.
Kishoreganj	এনো কি লিখছে দানো।	[eno' ki lik'che dsho]	See what is written here.
Narail	তালি তো বেগে দেবে।	[tal'i to beg'e gec'e]	Then it has increased.
Chittagong	কিনু অহিতো ন।	[kin'u ahito no]	Nothing will happen.
Sandvip	যেতে আসতে ন।	[jet'e asat'e no]	He won't come.
Sylhet	আমার জন্ম দেয়া থিএও।	[amer' jnnoq doe k'orjo]	Pray for me too.
Habiganj	ইস্কুলের বিজ্ঞ বাইচার?	[iskuler' vij'gn b'vicher]	What to say at school?
Narsingdi	উলটা পাটা কথা কল।	[ul'ta pata' ktha' kol]	Talk back and forth.
Tangail	এই ছবিটা দেখা নাই।	[ei' q'obiti da'xa' na'i]	This image has not been viewed.
Barishal	এনো বাঢ়ি গালে গাগ আয়।	[eno' bo'gi gal'e ga'g ay]	Now I get angry when I go home.

Table 1: Regional dialectal variations in Bengali.

## Contributions

- Ben-10 Dataset: 78 hours, 16,690 clips, 394 speakers across 10 districts.
- High Diversity: Average OOV rate of 59% compared to SCB.
- Linguistic Analysis: Comprehensive mapping of phonological and prosodic shifts.
- Benchmarking: Evaluation of 7 modern ASR models including Whisper, Conformer, and Wav2Vec2.

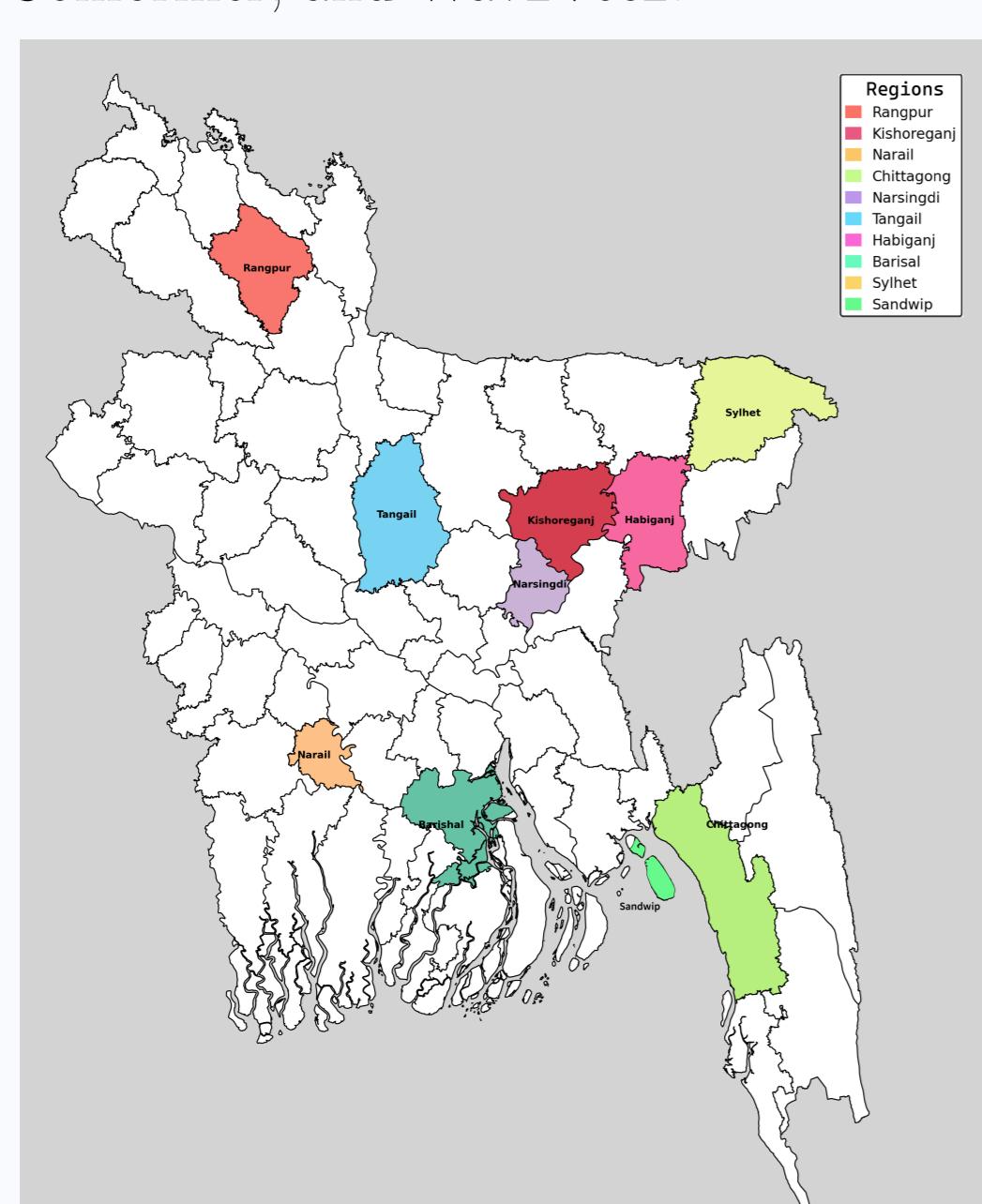


Figure 2: Mapped dialect regions with reference points

## Dataset Overview

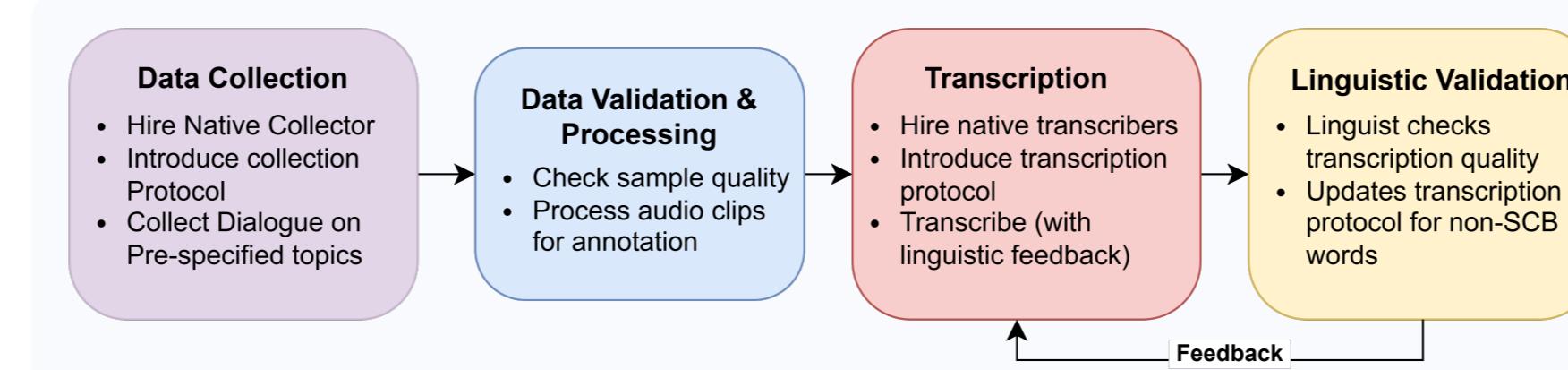


Figure 3: Dataset Creation Workflow: Collection → Validation → Transcription → Linguistic Review.

Regions: Rangpur, Kishoreganj, Narail, Chittagong, Narsingdi, Tangail, Sylhet, Habiganj, Barishal, Sandvip.

- Volume: 78 hours of spontaneous speech.
- Content: 155 conversation topics (family, politics, sports).
- Vocabulary: 62,762 unique words.

Districts	Sample Counts	Duration [H:M]	OOV%	WPM	Contributor
	Train Valid Test	Train Valid Test	Train Valid Test	(M/F/B)	
Rangpur	1,037 131 130	4:48 0:36 0:35	47.86 37.64 37.99	134.38	27 (19/8/0)
Kishoreganj	1,638 204 206	7:42 0:55 0:58	62.33 53.57 47.14	117.96	47 (25/18/4)
Narail	1,488 183 188	6:52 0:52 0:51	54.74 41.29 43.42	136.81	37 (21/12/4)
Chittagong	1,406 174 177	6:35 0:48 0:47	61.13 56.28 63.33	134.42	41 (15/22/4)
Narsingdi	1,098 136 137	5:04 0:38 0:37	52.24 39.18 37.59	148.53	26 (9/16/1)
Tangail	987 131 132	4:54 0:36 0:35	43.04 24.72 23.06	141.67	36 (18/11/7)
Habiganj	940 117 113	4:20 0:32 0:33	56.55 57.90 54.28	123.47	34 (19/15/0)
Barishal	796 105 105	3:45 0:30 0:30	48.29 43.62 44.86	123.79	26 (6/7/13)
Sylhet	2,903 356 362	13:34 1:50 1:41	63.24 50.62 50.27	126.5	94 (62/30/2)
Sandvip	1,049 129 132	4:48 0:36 0:37	61.91 51.61 52.77	144.12	26 (15/9/2)
Total	13,342 1,666 1,689	6:27 7:58 7:50	59.72 58.46 58.46	131.38	394 (209/148/37)

Table 2: Ben-10 dataset statistics. OOV → words unique to the district that are Out Of Vocabulary in comparison to SCB. M → Male, F → Female, B → Both.

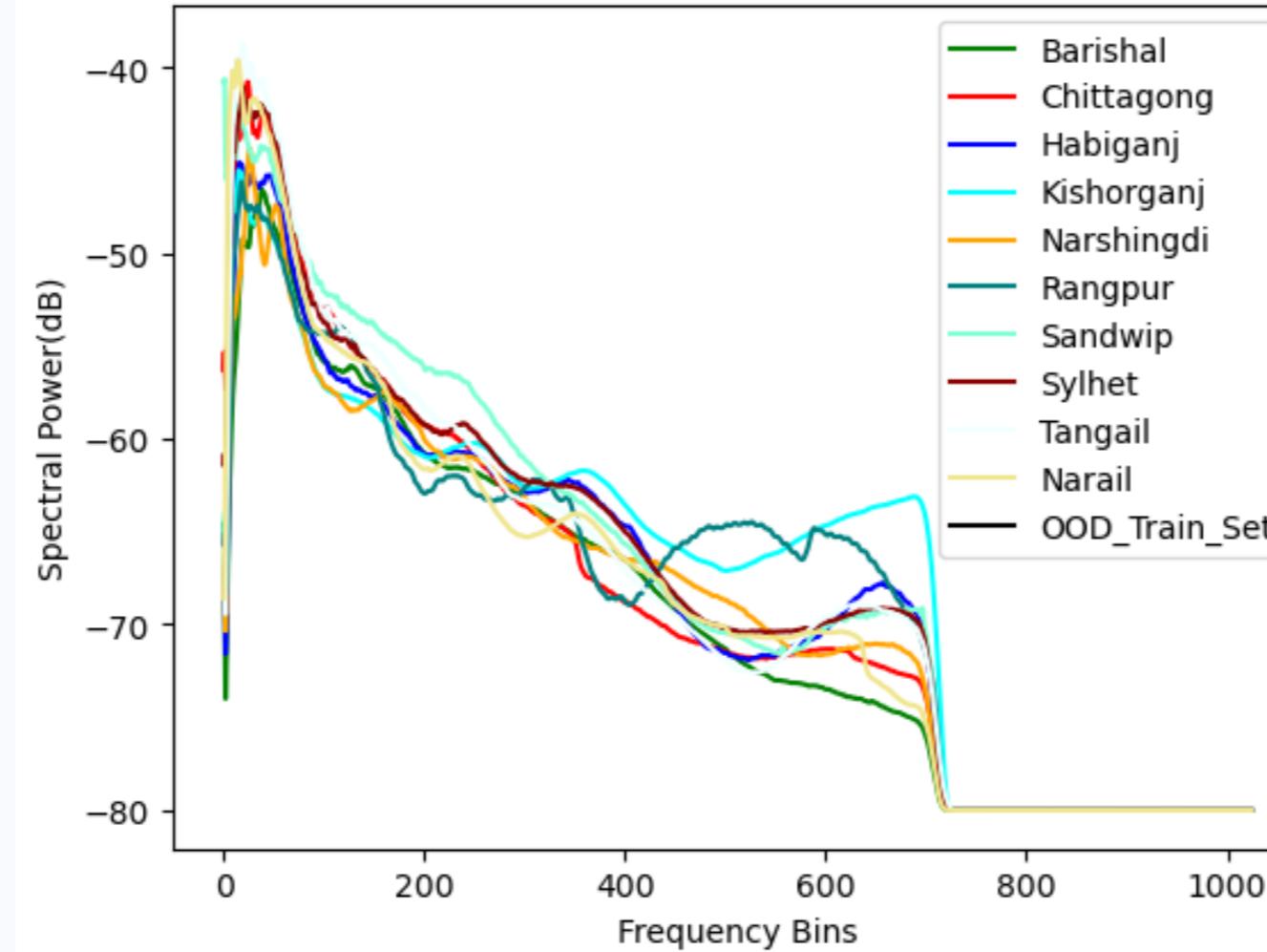


Figure 4: Long-term Spectral Average of Recordings from different dialects

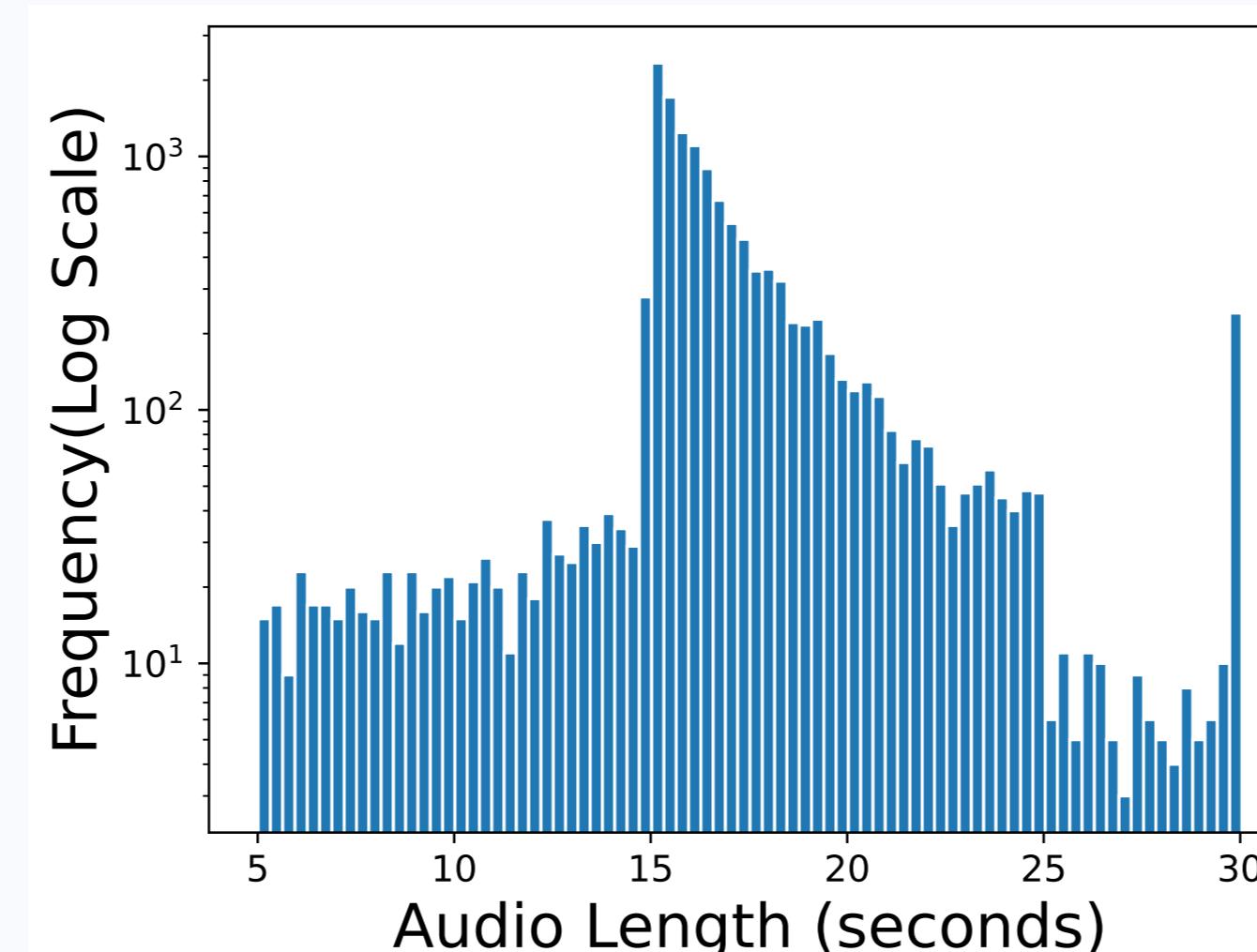
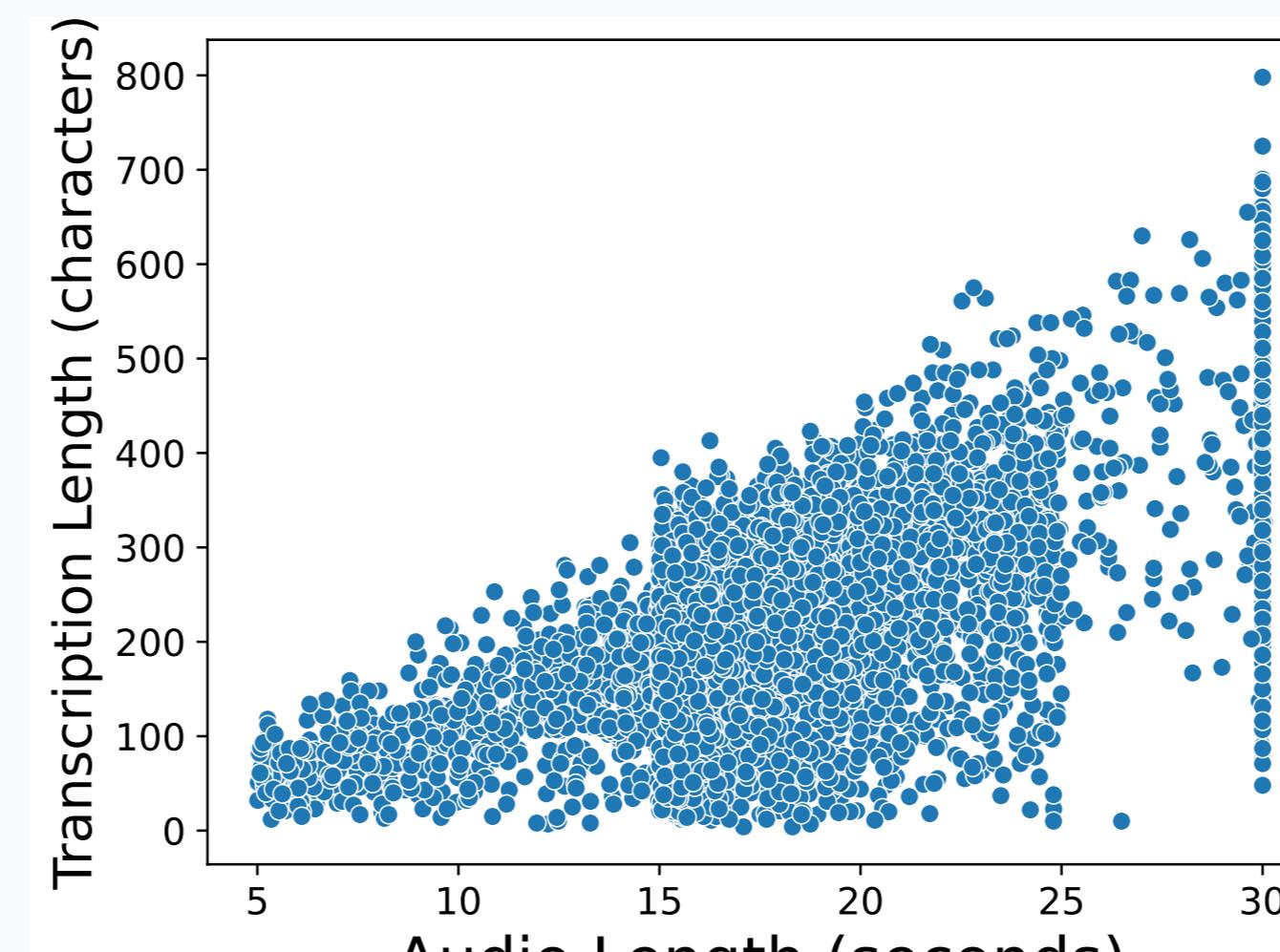


Figure 5: Recording length distribution in Training data



## Methodology

### Evaluated Models:

- Foundational: Whisper Large V3, Google Conformer ASR, Hishab Conformer.
- Fine-Tuned: Wav2vec2 (SCB vs Ben-10), Tugstugi (Whisper fine-tuned).

Metrics: WER (Word Error Rate) & CER (Character Error Rate)

## Benchmark Results

- Fine-tuning on Ben-10 significantly outperformed standard pre-trained models, achieving substantial improvements in dialectal ASR accuracy across all regions
- Zero-shot performance remains poor, highlighting the critical need for dialect-specific training data

Model	Avg WER ↓	Avg CER ↓
Whisper-Large-V3 (Zero-Shot)	1.13	0.81
Google ASR (API)	1.03	1.03
Hishab Conformer	0.87	0.56
Wav2vec2 (SCB Trained)	0.89	0.62
Tugstugi Whisper (SCB)	0.81	0.52
Wav2vec2 (Ben-10 Finetuned)	0.83	0.53
Tugstugi (Ben-10 Finetuned)	0.70	0.36

Table 3: Fine-tuning on Ben-10 reduces WER by ~14% compared to the best SCB model.

### Regional Variance:

- Best Performance: Sylhet & Tangail (Stronger training data/closer prosody).
- Worst Performance: Chittagong & Sandvip (High OOV and rapid speech).

## Why Models Fail: Error Analysis

Foundation models exhibit a strong SCB-Bias. When they encounter a dialect word, they often hallucinate the phonetically closest Standard Bengali word instead of transcribing the actual dialect.

Region	Ground Truth	Tugstugi	Tugstugi (Ben-10)
Rangpur	ধইৰছ [d̪oir̚bɔn]	দুঃসা [du'ʃɔa]	দিতো [d̪ito]
	গেসনু [gesnu]	<>	গেছেনু [ge'chenu]
Kishoreganj	কইতাম [kɔit̚am]	<>	করোনি [koroni]
	দুইনার [dujn̚er]	<>	দুইনের [dujn̚ner]
Narail	মাইসোতো [ma'is̚oto]	মায়ের সেতু [meer set̚u]	মাইনে ও তো [m̚in̚se o to]
	ম্যালাডিক [meledik]	মেলাডি [meledi]	মেলাডি [meledi]
Chittagong	বেককুন [bekkun̚]	এখনো [ekh'ono]	একখনো [ekh'ono]
	হারাপ [herap]	<>	খালা [k̚hala]
Sandvip	লাডিয়ালা [lepielen]	রিলেন্ডা [rilend̚a]	হেলা [hel̚r]
	এগেরে [eggere]	<>	এগেরে [egere]
Sylhet	ছোখো [c̚ok'ho]	সোকে [ʃo'ke]	ছো খেকে [c̚h oke]
	তাইনের [tein̚er]	<>	তাই [tei̚]
Habiganj	আছিলা [ac'hila]	চাইছিলা [caic'hile]	চাইছিলা [caic'hile]
	খিচুলিন [k̚i'c̚ul̚in]	পুরিখীরি [pr̚i'k̚ib̚ri]	ফিতুলিন [fitudin̚]
Narsingdi	খেলাই [k̚elāl̚ai]	ফেলে [fele]	ফেলায় [fele]
	ডাহাত [dehat̚]	ডাকতেছে [dk̚teq̚e]	ভাত [det̚]
Tangail	দাহে [dehe]	দেহে [dehe]	দেহে [dehe]
	থোয় [t̚hɔi̚]	হ্য [ho̚]	দেখায় [dek̚v̚i̚]
Barishal	কাইলগো [kulgo]	<>	কাজ-মাজ [kej-mej]
	বাইরাইয়া [ba'erai̚]	বাইরে [ba're]	বাইরে [ba're]

Table 4: Examples of dialectal words that are incorrectly transcribed by models. SCB words are underlined.