

# Human-1 by Josh Talks : A Full-Duplex Conversational Modeling Framework in Hindi using Real-World Conversations

Bhaskar Singh  
bhaskarsingh@joshtalks.com  
JoshTalks

Shobhit Banga  
shobhit@joshtalks.com  
JoshTalks

Pranav Sharma  
pranav@joshtalks.com  
JoshTalks

## Abstract

Full-duplex spoken dialogue systems can model natural conversational behaviours such as interruptions, overlaps, and backchannels, yet such systems remain largely unexplored for Indian languages. We present the first open, reproducible full-duplex spoken dialogue system for Hindi by adapting Moshi, a state-of-the-art duplex speech architecture, using a custom Hindi tokeniser and training on 26,000 hours of real spontaneous conversations collected from 14,695 speakers with separate speaker channels, enabling direct learning of turn-taking and overlap patterns from natural interactions. To support Hindi text generation, we replace the original English tokeniser and reinitialise text-vocabulary-dependent parameters while retaining the pre-trained audio components. We propose a two-stage training recipe—large-scale pre-training followed by fine-tuning on 1,000 hours of conversational data. Evaluation through the prompted dialogue continuation paradigm with both automatic metrics and human judgements demonstrates that the resulting model generates natural and meaningful full-duplex conversational behaviour in Hindi. This work serves as a first step toward real-time duplex spoken dialogue systems for Hindi and other Indian languages.

**Index Terms:** spoken dialogue system, full-duplex, speech-to-speech, language adaptation

## 1 Introduction

Human conversation is inherently full-duplex: speakers routinely overlap, produce backchannels, and interrupt one another in ways that half-duplex systems—which wait for one speaker to finish before responding—cannot capture [1–3]. Recent end-to-end speech dialogue models have begun to address this limitation by jointly modeling both speaker streams in parallel [4–7].

However, progress in full-duplex spoken dialogue has been concentrated almost exclusively on English, with limited exploration in other languages. This leaves a significant gap for languages with large speaker populations, distinct phonological systems, and different conversational norms. Hindi, spoken by over 600 million people, is one such language. Hindi conversations are characterised by frequent backchannels, code-mixing with English, and prosodic patterns that differ substantially from those found in English dialogue [8, 9].

Building a full-duplex dialogue system for Hindi presents several technical challenges. First, existing duplex architectures use tokenisers trained on English (or other Latin-script languages), which are inefficient for Devanagari script and produce excessive token fragmentation. Second, replacing the tokeniser necessitates reinitialisation of vocabulary-dependent model parameters, fundamentally changing the training dynamics from standard fine-tuning to partial re-training, as these parameters must be learned from scratch while the remaining pre-trained weights are preserved. Third, and most critically, training full-duplex models requires large-scale stereo recordings with separate speaker channels capturing natural turn-taking, overlaps, and backchannels. The quality and scale of such data is arguably the most important factor in the performance of duplex dialogue systems, yet such large-scale stereo conversational resources remain largely unavailable for Hindi and most Indian languages.

In this paper, we address these challenges. We collect a large-scale corpus of 26,000 hours of real Hindi spontaneous conversations from 14,695 speakers with separate channel recordings—to our knowledge, the largest purpose-built conversational speech corpus for any Indian language—and present a framework for adapting full-duplex speech dialogue architectures to Hindi. Our results demonstrate that *Real world natural spontaneous conversational data* is the critical enabler for building effective full-duplex dialogue systems in new languages.

## 2 Related Work

The dGSLM model [1] was among the first to model two-channel spoken dialogue jointly. Moshi [4] advanced this by combining a 7B-parameter text LLM with a neural audio codec (Mimi) and a hierarchical RQ-Transformer to achieve real-time full-duplex dialogue in English. Other recent efforts include SyncLLM [5] and OmniFlatten [6]. These systems have been developed primarily in English.

Cross-lingual adaptation has been explored for ASR [12] and TTS [13], but adapting full-duplex dialogue models to new languages remains unexplored. SpeechGPT [10] and AudioPaLM [11] demonstrated multilingual speech capabilities but in a turn-based setting. Neural audio codecs such as Mimi are trained predominantly on English; whether they generalise to

typologically different languages is an open question. To our knowledge, no prior work has presented an open full-duplex spoken dialogue system for any Indian language.

## 3 Method

### 3.1 Base Architecture: Moshi

We build upon Moshi [4], a full-duplex spoken dialogue model with three components. **Mimi** is a neural audio codec that encodes 24 kHz speech into discrete tokens at 12.5 Hz using 8 codebook layers—layer 1 captures semantic content while layers 2–8 capture acoustic detail. The **RQ-Transformer** is a hierarchical architecture: the *Temporal Transformer* (7B parameters, based on Helium) models 17 parallel streams per timestep (1 text + 8 Moshi audio + 8 user audio), producing a hidden vector  $z_s$  from which a text token is sampled via the Text Linear layer; the *Depth Transformer* then autoregressively generates 16 audio tokens conditioned on  $z_s$ . PAD tokens are inserted at timesteps without corresponding text, and a one-step acoustic delay stabilises generation.

Preliminary evaluation confirms that Mimi achieves acceptable reconstruction on Hindi without retraining (Section 5.1), so we **freeze Mimi** throughout and train only the RQ-Transformer, substantially reducing computational cost.

### 3.2 Hindi Tokeniser and Parameter Reinitialisation

The original architecture uses a SentencePiece tokeniser [16] trained on English text with a 32,000-word vocabulary. This tokeniser is fundamentally incompatible with Hindi’s Devanagari script, producing excessive fragmentation that degrades both training efficiency and language modelling quality.

We replace it with a SentencePiece model trained on a large Hindi text corpus, also with a 32,000-word vocabulary. This replacement requires reinitialisation of all vocabulary-dependent parameters which includes Text token embedding tables in both the Temporal and Depth Transformers. And The Text Linear projection layer that maps hidden states to token probabilities.

This reinitialisation has a critical consequence: while audio processing components retain their pre-trained weights, all text representations must be learned from scratch. The task is therefore *partial pre-training* rather than fine-tuning, which requires substantially different hyperparameter choices (Section 3.4).

### 3.3 Data

#### 3.3.1 Pre-training Data

A key contribution of this work is the collection of a large-scale Hindi conversational speech corpus specifically designed for full-duplex dialogue modelling. We

Table 1: Summary of collected Hindi conversational corpus.

Characteristic	Value
Total duration	26,000 hours
Unique speakers	14,695
Recording type	Spontaneous conversations
Channels	Stereo (separate per speaker)
Style	Spontaneous, unscripted
Quality control	Trained annotators + manual checks

collected 26,000 hours of real two-person Hindi spontaneous conversations involving 14,695 unique speakers. The data was gathered through a dedicated collection effort with hired participants engaging in spontaneous, unscripted conversations on diverse topics. This corpus forms part of our broader, ongoing infrastructure program to collect 1 million hours of natural, spontaneous conversational speech in Hindi - building the scale of data required to support next-generation conversational and full-duplex speech systems. Each conversation was recorded with separate speaker channels, providing natural stereo dialogue without the need for artificial speaker diarisation—a significant advantage over prior approaches that construct pseudo-stereo data from monophonic recordings [4].

**Domain-matched data.** A critical design decision was collecting data in a *conversational* setting rather than compiling read speech, broadcast news, or monologue recordings. Full-duplex dialogue models must learn turn-taking, backchannels, interruptions, and overlaps—phenomena that only occur naturally in real conversations. Our corpus directly captures these dynamics in its native stereo format, eliminating the noise and artifacts introduced by post-hoc diarisation of monophonic sources.

**Quality assurance.** The collected data underwent extensive quality checks by trained annotators who evaluated recordings against predefined criteria including recording clarity, speaker balance, transcription accuracy, and conversational naturalness. Recordings failing quality thresholds were excluded. This rigorous curation pipeline ensures that model training is driven by high-quality signal rather than noisy data that can limit convergence and output quality.

Transcriptions were obtained through a combination of manual and automatic methods and token-level timestamps were obtained using WhisperX [17] to align text and audio token sequences, with PAD tokens inserted at timesteps without corresponding text tokens. Table 1 summarises the dataset characteristics.

Each sample is represented as 2,048 timesteps ( $\approx 2.7$  minutes) across 17 token streams. The Hindi data exhibits a PAD token ratio of approximately 75%, which is higher than reported for English (65%) [4]. This is likely because Devanagari characters and Hindi words encode more phonemic content per token, making text tokens sparser relative to audio.

### 3.3.2 Fine-tuning Data

From the collected corpus, we manually curate a subset of 1,000 hours selected for clear pronunciation, low background noise, balanced speaker participation, and natural conversational prosody. For validation, we hold out a small set of conversations selected for clear pronunciation, low background noise, and natural conversational prosody, which is used for checkpoint selection during fine-tuning. For evaluation, we use separately recorded conversations that were not part of the training corpus, ensuring that the model is tested on fully unseen data. The remaining curated data (approximately 990 hours) is used for fine-tuning.

## 3.4 Training Procedure

### 3.4.1 Stage 1: Pre-training

Due to the reinitialised embedding layers, we use a learning rate of  $3 \times 10^{-5}$ —matching the original Moshi pre-training rate rather than the lower rates typical of adaptation. We use AdamW [18] with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.95$ ,  $\epsilon = 10^{-5}$ , and weight decay of 0.1. Following Moshi [4], PAD token losses are reduced by 50% and the loss ratio between semantic and acoustic tokens is set to 100:1. Mimi parameters are frozen throughout; only the RQ-Transformer is trained.

Training was conducted on  $8 \times$  NVIDIA H100 80GB GPUs with bf16 mixed precision. The effective batch size was 64 samples (per-device batch size of 4, gradient accumulation of 2), corresponding to approximately 2.9 hours of audio per update. We trained for 1 epoch ( $\approx 10,000$  steps) in approximately 13 hours.

### 3.4.2 Stage 2: Fine-tuning

Starting from the pre-trained checkpoint, we fine-tune on the approximately 990-hour curated set using split learning rates:  $2 \times 10^{-6}$  for the Temporal Transformer and  $4 \times 10^{-6}$  for the Depth Transformer. The higher Depth Transformer rate allows faster adaptation of acoustic token prediction while preserving learned text representations.

The effective batch size is reduced to 16 (per-device batch size of 2, 8 GPUs, no gradient accumulation), and warmup is set to 50 steps. We evaluate on the validation set every 802 steps and save checkpoints at the same interval.

## 4 Training Analysis

### 4.1 Pre-training Convergence

Figure 1 shows training loss curves during pre-training. All components exhibit rapid convergence within the first 2,000 steps, stabilising thereafter. Total loss decreased from  $\approx 20$  at initialisation to  $\approx 3.5$  at convergence. Text accuracy on non-PAD tokens reached  $\approx 70\%$ , and audio accuracy stabilised at 48–53% across codebook levels. Training loss plateaued after approximately 4,000 steps, indicating efficient extraction of

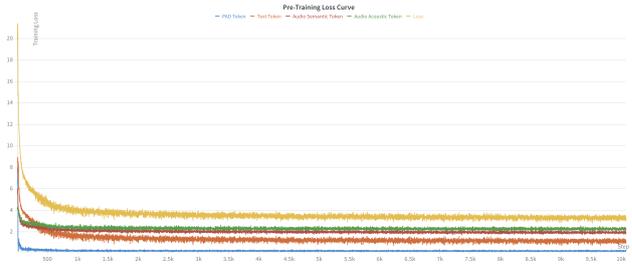


Figure 1: Training loss during pre-training on 26,000 hours of Hindi speech. All metrics converge within 2,000–4,000 steps.

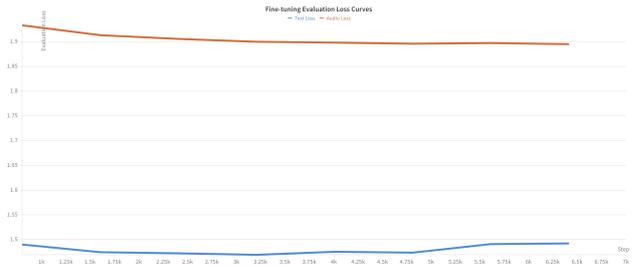


Figure 2: Evaluation loss during fine-tuning. Text loss overfits after step  $\approx 4,800$ , while audio loss continues to improve marginally. Optimal checkpoint is at step 4,812.

patterns from the 26,000-hour dataset within less than half an epoch.

## 4.2 Fine-tuning and Overfitting

Fine-tuning revealed distinct convergence patterns for text and audio components (Figure 2). Text validation loss reached its minimum at approximately step 4,800, after which it increased steadily—a clear indication of overfitting. Audio validation loss continued improving slightly beyond this point but with diminishing returns.

The optimal checkpoint was selected at step 4,812 based on minimum total validation loss (3.370), with text evaluation loss of 1.474 and audio evaluation loss of 1.896. This highlights the importance of validation-based early stopping when fine-tuning on limited curated data, as training loss continued to decrease well past the optimal point for generalisation.

## 5 Evaluation

We evaluate using the prompted dialogue continuation paradigm [1, 4]. Each test dialogue is split into 30-second segments; the first 10 seconds serve as prompt, and the model generates the subsequent 20 seconds at temperatures  $\tau \in \{0.8, 0.9, 1.0\}$ .

### 5.1 Codec Quality on Hindi

To verify that Mimi generalises to Hindi, we compare ground-truth audio with re-synthesised audio (encode

Table 2: Mimi codec quality on Hindi (ground-truth vs. re-synthesis, 654 segments).

Metric	Score
PESQ ( $\uparrow$ )	$2.55 \pm 0.37$
STOI ( $\uparrow$ )	$0.878 \pm 0.027$

Table 3: Perplexity ( $\downarrow$ ) computed using Sarvam-1 (2B) on Whisper-v3 transcriptions of generated speech.

Model	$\tau$	PPL $\downarrow$
Ground-truth	–	237.1
Human-1	0.8	<b>356.9</b>
Human-1	0.9	467.1
Human-1	1.0	640.6

$\rightarrow$  decode through Mimi) using PESQ and STOI on 654 twenty-second segments from the test set.

As shown in Table 2, the high STOI (0.878) confirms that Hindi speech remains largely intelligible after codec processing. The moderate PESQ (2.55) reflects the expected quality–bitrate trade-off at 1.1 kbps, consistent with Mimi’s design for low-latency streaming rather than maximum reconstruction fidelity. These results justify freezing Mimi during training.

## 5.2 Perplexity

Following [1, 4], we measure the perplexity (PPL) of a Hindi language model (Sarvam-1, 2B parameters) on Whisper-v3 transcriptions of generated speech. As shown in Table 3, lower temperatures yield more fluent Hindi, with  $\tau = 0.8$  achieving the best PPL (356.9). Quality degrades progressively with higher temperature as increased sampling randomness reduces linguistic coherence.

## 5.3 Human Evaluation

We conduct human evaluation with 130 native Hindi speakers who completed 2,125 paired comparisons between human and model-generated speech, rating Naturalness and Clarity on 5-point scales and assessing conversational quality through binary rubrics.

As shown in Table 4, the model achieves a naturalness score of 4.10 vs. 4.55 for human speech, with 66.9% of pairwise comparisons resulting in ties—indicating that generated speech frequently approaches human quality. Conversational rubric results show that while the model produces human-like speech in 85% of cases, maintaining contextual relevance (53%) and producing complete responses (42%) remain key areas for improvement.

## 5.4 Turn-Taking Analysis

To examine whether the model acquires Hindi conversational dynamics, we compute turn-taking statistics following [1]: Inter-Pausal Units (IPU), Pause, Gap,

Table 4: Human evaluation: perceptual scores and conversational rubrics.

Perceptual Scores (5-point scale)	
Naturalness (Human / Model)	4.55 / 4.10
Clarity (Human / Model)	4.05 / 3.04
Preference (H. / M. / Tie)	30.0% / 3.1% / 66.9%
Conversational Rubrics (Pass Rate)	
Human-like interaction	$\approx 85\%$
Appropriateness (follows prompt)	$\approx 53\%$
Completion (complete reply)	$\approx 42\%$

Table 5: Turn-taking statistics (per minute) in generated dialogues.

Model	$\tau$	IPU	Pause	Gap	Overlap
Ground-truth	–	35.30	10.49	8.51	3.03
Human-1	0.8	23.12	9.16	6.77	1.67
Human-1	0.9	29.14	9.24	8.54	4.30
Human-1	1.0	38.90	11.67	8.10	9.68

and Overlap duration per minute. We apply energy-based voice activity detection to each speaker channel independently.

Table 5 shows that  $\tau = 0.9$  produces turn-taking dynamics closest to ground-truth, with near-identical gap duration (8.54 vs. 8.51 s/min) and comparable pause patterns (9.24 vs. 10.49 s/min). Lower temperature ( $\tau = 0.8$ ) yields overly conservative turn-taking with reduced IPU count and minimal overlap, while  $\tau = 1.0$  produces excessive overlap (9.68 vs. 3.03 s/min), indicating less controlled conversational dynamics.

## 6 Conclusion

We presented a Hindi full-duplex spoken dialogue model by adapting Moshi through tokeniser replacement and two-stage training on a novel 26,000-hour corpus of real Hindi telephone conversations from 14,695 speakers. Our finding that Mimi generalises to Hindi without retraining enabled an efficient adaptation strategy focused on the RQ-Transformer alone. Crucially, our results demonstrate that high-quality, domain-matched conversational data is the key enabler for full-duplex dialogue in new languages—the model acquired natural turn-taking behaviour directly from real conversations without synthetic augmentation. Analysis of training dynamics provides practical guidance for adapting Moshi to other languages. Given that our model achieved coherent dialogue with 26,000 hours while the original Moshi leveraged 7 million, we believe that scaling quality-controlled conversational data represents the most promising path toward highly fluent full-duplex dialogue systems for Hindi and other Indian languages.

## References

- [1] T. A. Nguyen, E. Kharitonov, J. Copet, Y. Adi, W.-N. Hsu, A. Elkahky, P. Tomasello, R. Algayres, B. Sagot, A. Mohamed, and E. Dupoux, “Generative spoken dialogue language modeling,” *Transactions of the Association for Computational Linguistics*, pp. 250–266, 2023.
- [2] P. Wang, S. Lu, Y. Tang, S. Yan, W. Xia, and Y. Xiong, “A full-duplex speech dialogue scheme based on large language model,” in *Proc. NeurIPS*, 2024.
- [3] Z. Ma, Y. Song, C. Du, J. Cong, Z. Chen, Y. Wang, Y. Wang, and X. Chen, “Language model can listen while speaking,” *arXiv preprint arXiv:2408.02622*, 2024.
- [4] A. Défossez, L. Mazaré, M. Orsini, A. Royer, P. Pérez, H. Jégou, E. Grave, and N. Zeghidour, “Moshi: a speech-text foundation model for real-time dialogue,” *arXiv preprint arXiv:2410.00037*, 2024.
- [5] B. Veluri, B. N. Peloquin, B. Yu, H. Gong, and S. Gollakota, “Beyond turn-based interfaces: Synchronous LLMs as full-duplex dialogue agents,” in *Proc. EMNLP*, pp. 21390–21402, 2024.
- [6] Q. Zhang, L. Cheng, C. Deng, Q. Chen, W. Wang, S. Zheng, J. Liu, H. Yu, and C. Tan, “Omniflaten: An end-to-end GPT model for seamless voice conversation,” *arXiv preprint arXiv:2410.17799*, 2024.
- [7] W. Yu, S. Wang, X. Yang, X. Chen, X. Tian, J. Zhang, G. Sun, L. Lu, Y. Wang, and C. Zhang, “SALMONN-omni: A codec-free LLM for full-duplex speech understanding and generation,” *arXiv preprint arXiv:2411.18138*, 2024.
- [8] R. Hayashi, “Simultaneous talk—from the perspective of floor management of English and Japanese speakers,” *World Englishes*, vol. 7, no. 3, pp. 269–288, 1988.
- [9] M. Stubbe, “Are you listening? Cultural influences on the use of supportive verbal feedback in conversation,” *Journal of Pragmatics*, vol. 29, no. 3, pp. 257–289, 1998.
- [10] D. Zhang, S. Li, X. Zhang, J. Zhan, P. Wang, Y. Zhou, and X. Qiu, “SpeechGPT: Empowering large language models with intrinsic cross-modal conversational abilities,” in *Findings of EMNLP*, pp. 15757–15773, 2023.
- [11] P. K. Rubenstein, C. Asawaroengchai, D. D. Nguyen, et al., “AudioPaLM: A large language model that can speak and listen,” *arXiv preprint arXiv:2306.12925*, 2023.
- [12] A. Conneau, A. Baevski, R. Collobert, A. Mohamed, and M. Auli, “Unsupervised cross-lingual representation learning for speech recognition,” in *Proc. Interspeech*, pp. 2426–2430, 2020.
- [13] X. Zhang, Z. Tan, R. Huang, et al., “VALL-E X: Speak foreign languages with your own voice,” *arXiv preprint arXiv:2303.03926*, 2023.
- [14] A. Défossez, J. Copet, G. Synnaeve, and Y. Adi, “High fidelity neural audio compression,” *Transactions on Machine Learning Research*, 2023.
- [15] N. Zeghidour, A. Luebs, A. Omran, J. Skoglund, and M. Tagliasacchi, “SoundStream: An end-to-end neural audio codec,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 30, pp. 495–507, 2021.
- [16] T. Kudo and J. Richardson, “SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing,” in *Proc. EMNLP: System Demonstrations*, pp. 66–71, 2018.
- [17] M. Bain, J. Huh, T. Han, and A. Zisserman, “WhisperX: Time-accurate speech transcription of long-form audio,” *arXiv preprint arXiv:2303.00747*, 2023.
- [18] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *Proc. ICLR*, 2019.
- [19] H. Touvron, L. Martin, K. Stone, et al., “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.
- [20] S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He, “ZeRO: Memory optimizations toward training trillion parameter models,” in *Proc. SC*, 2020.