



# PRESIDENCY UNIVERSITY

Private University Estd. in Karnataka State by Act No. 41 of 2013  
Itgalpura, Rajankunte, Yelahanka, Bengaluru - 560064



## VoiceBridge: An AI-Powered Framework for Low-Cost Multilingual Video Dubbing into Indian Regional Languages

A PROJECT REPORT

*Submitted by*

Annmary Jojo - 20221CSG0140

Aqsa Mehareen D - 20221CSG0118

Agampreeth H - 20221CSG0117

*Under the guidance of,*

**Dr. MADHUSUDHAN M V**

**BACHELOR OF TECHNOLOGY**

IN

**COMPUTER SCIENCE AND TECHNOLOGY**

**PRESIDENCY UNIVERSITY**

**BENGALURU**

**DECEMBER 2025**



# PRESIDENCY UNIVERSITY

Private University Estd. in Karnataka State by Act No. 41 of 2013  
Itgalpura, Rajankunte, Yelahanka, Bengaluru – 560064



## PRESIDENCY SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

### BONAFIDE CERTIFICATE

Certified that this report “VoiceBridge: An AI-Powered Framework for Low-Cost Multilingual Video Dubbing into Indian Regional Languages” is a Bonafide work of “ANMARY JOJO (20221CSG0140), AQSA MEHEREEN D(20221CSG0118), AGAMPREET H (20221CSG0117)”, who have successfully carried out the project work and submitted the report for partial fulfilment of the requirements for the award of the degree of BACHELOR OF TECHNOLOGY in COMPUTER SCIENCE AND TECHNOLOGY during 2025-26.

Dr. Madhusudhan M V  
Project Guide  
PSCS  
Presidency University

Dr. Sharmasti Vaji  
Program Coordinator  
PSCS  
Presidency University

Dr. Sampath A K  
Dr. Geetha A  
School Coordinators  
PSCS  
Presidency University

Dr. Anandaraj S P  
Head of the Department  
Project PSCS  
Presidency University

Dr. Shakkeera L  
Associate Dean  
PSCS  
Presidency University

N.J.S  
Dr. Duraiapandian N  
Dean  
PSCS & PSIS  
Presidency University

#### Examiners

Sl.no	Name	Signature	Date
1	Ramdonathy K	Aug	4/12/25
2	Buram Mohammad	Jaf	04/12/25

Padma T.G. P.L.S.Th 4/12/25

**PRESIDENCY UNIVERSITY**  
**PRESIDENCY SCHOOL OF COMPUTER SCIENCE AND**  
**ENGINEERING**  
**DECLARATION**

We the students of final year B.Tech in COMPUTER SCIENCE AND TECHNOLOGY at Presidency University, Bengaluru, named **Annmary Jojo**, **Aqsa Mehareen D**, **Agampreeth H S**, hereby declare that the project work titled "**VoiceBridge: An AI-Powered Framework for Low-Cost Multilingual Video Dubbing into Indian Regional Languages**" has been independently carried out by us and submitted in partial fulfilment for the award of the degree of B.Tech in COMPUTER SCIENCE AND TECHNOLOGY during the academic year of 2025-26. Further, the matter embodied in the project has not been submitted previously by anybody for the award of any Degree or Diploma to any other institution.

Annmary Jojo	20221CSG0140
Aqsa Mehareen D	20221CSG0118
Agampreeth H S	20221CSG0117

PLACE: BENGALURU  
DATE: 25 November 2025

## Abstract

VoiceBridge is an intelligent, scalable, and economically feasible multilingual dubbing framework that upgrades the accessibility of English video content for diverse Indian linguistic audiences. It focuses on South Indian languages like Kannada, Tamil, Telugu, and Malayalam. This system integrates state-of-the-art open-source AI technologies: Whisper for ASR, IndicTrans2 for NMT, and Coqui/Indic-TTS for high-quality TTS into an end-to-end pipeline that can perform transcription, translation, and speech generation with audio-video synchronization.

VoiceBridge focuses on aspects of cultural and linguistic naturalness, including prosody adaptation, regional accent customization, and cloning speaker voices to make the dubbed output relevant and engaging for native viewers. The framework will be designed in a user-friendly way, making it possible for individual users, educators, and content developers to translate and dub videos without requiring technical expertise or high financial resources.

Performance evaluation has been done based on Word Error Rate, Character Error Rate, speech naturalness scoring, and audio-lip sync accuracy. Our system returns a WER of 11.9% and CER of 11.09% with very minor synchronization errors of around 90 milliseconds. The results show good recognition and translation performance even for code-mixed speeches and region-specific pronunciations—common challenges in Indian multilingual scenarios.

By utilizing open-source models, VoiceBridge achieves a massive reduction in deployment cost and computational needs while making it feasible for low-resource settings and local institutions. The framework aims to bridge the digital language divide in India, thereby allowing inclusive access to educational, entertainment, and informational video content for millions of native speakers. With plans to enhance it for more regional languages, real-time dubbing, and emotion-aware synthesis in the near future, this represents an important step toward democratizing digital knowledge and preserving linguistic diversity.

## **ACKNOWLEDGEMENT**

For completing this project work, We/I have received the support and the guidance from many people whom I would like to mention with deep sense of gratitude and indebtedness. We extend our gratitude to our beloved **Chancellor, Pro-Vice Chancellor, and Registrar** for their support and encouragement in completion of the project.

I would like to sincerely thank my internal guide **Dr. Madhusudhan M V, Associate Professor**, Presidency School of Computer Science and Engineering, Presidency University, for his moral support, motivation, timely guidance and encouragement provided to us during the period of our project work.

I am also thankful to **Dr. Anandaraj, Professor, Head of the Department, Presidency School of Computer Science and Engineering** Presidency University, for his mentorship and encouragement.

We express our cordial thanks to **Dr. Duraipandian N**, Dean PSCS & PSIS, **Dr. Shakkeera L**, Associate Dean, Presidency School of computer Science and Engineering and the Management of Presidency University for providing the required facilities and intellectually stimulating environment that aided in the completion of my project work.

We are grateful to **Dr. Sampath A K, and Dr. Geetha A, PSCS** Project Coordinators, **Dr. Sharmasti vali, Program Project Coordinator**, Presidency School of Computer Science and Engineering, or facilitating problem statements, coordinating reviews, monitoring progress, and providing their valuable support and guidance.

We are also grateful to Teaching and Non-Teaching staff of Presidency School of Computer Science and Engineering and also staff from other departments who have extended their valuable help and cooperation.

Annmary Jojo

Aqsa Mehareen D

Agampreeth H S

## Table of Content

Sl. No.	Title	Page No.
	Declaration	ii
	Acknowledgement	iii
	Abstract	iv
	List of Figures	vii
	List of Tables	viii
	Abbreviations	ix
1.	Introduction 1.1 Background 1.2 Statistics of project 1.3 Prior existing technologies 1.4 Proposed approach 1.5 Objectives 1.6 SDGs 1.7 Overview of project report	1
2.	Literature review	12
3.	Methodology	17
4.	Project management 4.1 Project timeline 4.2 Risk analysis 4.3 Project budget	27
5.	Analysis and Design 5.1 Requirements 5.2 Block Diagram 5.3 System Flow Chart 5.4 Choosing devices 5.5 Designing units 5.6 Standards 5.7 Mapping with IoTWF reference model layers 5.8 Domain model specification	31

	5.9 Communication model 5.10 IoT deployment level 5.11 Functional view 5.12 Mapping IoT deployment level with functional view 5.13 Operational view 5.14 Other Design	
6.	Hardware, Software and Simulation 6.1 Hardware 6.2 Software development tools 6.3 Software code 6.4 Simulation	43
7.	Evaluation and Results 7.1 Test points 7.2 Test plan 7.3 Test result 7.4 Insights	49
8.	Social, Legal, Ethical, Sustainability and Safety Aspects 8.1 Social aspects 8.2 Legal aspects 8.3 Ethical aspects 8.4 Sustainability aspects 8.5 Safety aspects 8.6 Collaborative and educational aspects 8.7 Future adaptability and scalability	55
9.	Conclusion	59
	References	63
	Appendix	65

## List of Figures

<b>Figure</b>	<b>Caption</b>	<b>Page no</b>
Fig 1.1	Block Diagram of Processing Pipeline for Video Dubbing in Indian Languages	05
Fig 3.1	Input Handling Workflow	22
Fig 3.3	ASR Flowchart	27
Fig 3.8	Output Handling Workflow	34
Fig 4.1	Gantt Chart - VoiceBridge Project Timeline	37
Fig 5.1	Functional Block Diagram	44
Fig 5.2	System Flowchart	45
Fig 5.3	Domain Model Description for VoiceBridge	51
Fig 7.1	Input Video Upload	65
Fig 7.2	Output Generated in Kannada Language	67
Fig 7.3	Output Generated in Tamil Language	67
Fig 7.4	Output Generated in Malayalam Language	67
Fig 7.5	Output Generated in Telugu Language	67
Fig 7.6	Final Dubbed Video Output Screen	68
Fig 7.7	Comparative Analysis of WER for Different ASR Models	69
Fig 7.8	Comparative Analysis of CER for Different ASR Models	70
Fig 7.9	Comparative Analysis of WER for Different ASR Models	70
Fig 7.10	Comparative Analysis of CER for Different ASR Models	71

## List of Tables

<b>Table</b>	<b>Caption</b>	<b>Page no</b>
Table 2.1	Summary of Literature reviews	20
Table 4.1	Project Planning Timeline	36
Table 4.2	Project Implementation Timeline	36
Table 4.3	PESTLE Analysis for the Project	38
Table 5.1	Summary of System Requirements	42
Table 5.2	Table 5.2 Mapping VoiceBridge with the IoTWF Reference Model	49
Table 5.3	Domain Model Description for VoiceBridge	50
Table 7.1	Test points and measurements	63
Table 7.4	Comparative Analysis of WER and CER different ASR models	69

## Abbreviations

Abbreviation	Full Form
<b>ASR</b>	Automatic Speech Recognition
<b>CER</b>	Character Error Rate
<b>GCSM</b>	Glossary & Code-Switching Module
<b>TT</b>	Text Translation
<b>TTS</b>	Text-to-Speech
<b>WER</b>	Word Error Rate
<b>WWW</b>	World Wide Web

## Chapter 1

# Introduction

In the last decade, the way people consume knowledge and information has undergone a dramatic transformation. With the rapid growth of the internet and digital technologies, video-based content has become one of the most powerful and engaging mediums for communication, education, and entertainment. Online platforms such as YouTube, Coursera, edX, and Khan Academy have brought learning out of the confines of the traditional classroom and made it accessible to millions of learners worldwide [1]. The increasing penetration of affordable smartphones and high-speed internet in many of the developed and developing countries like India has further accelerated this shift. For many learners, especially in remote or rural areas, videos offer the first and sometimes the only exposure to structured learning resources. Compared to offline lectures and printed material, videos offer several advantages: they are reusable, they allow learners to pause and re-watch difficult concepts, and they combine both audio and visual cues to enhance retention. In short, online video content has become indispensable to modern education and information dissemination. A UNESCO survey of 61 countries found that 90% of high-income education systems rapidly adopted online learning platforms during school closures—showcasing the nimble adaptation of digital infrastructure under crisis conditions [2].

Despite this global revolution in digital learning, language continues to remain a significant barrier. The vast majority of educational and informational videos available online are in English or a handful of other widely spoken global languages [3]. While this benefits a large international audience, it excludes a massive segment of the population in countries like India, where hundreds of millions of people are more comfortable in their native tongues. English dominates the World Wide Web (WWW) as around 62.5% of the material available on the Internet is written in this language [4]. This problem is part of a larger digital divide, as over 80% of online content is available in just 10 languages creating a significant information access gap for speakers of the world's remaining 7,000 languages [5]. According to the 2011 census of India, only 10.67% of the total population spoke English, of which only 0.02% had it as their first language [6]. In 2023, approximately 43.4% of the Indian population was using the Internet, and this figure is expected to rise to around 62.8% in the coming years [7]. This statistic reveals a gap in the amount of available content on the internet and the number of Indians who can access it [8]. India is not just a linguistically diverse nation, it is one of the

most linguistically rich countries in the world, with 22 officially recognized languages and hundreds of regional dialects [9]. Although the digital landscape is rapidly expanding, English continues to maintain its control of the digital space while Indian languages, especially South Indian languages like Tamil, Telugu, Kannada, and Malayalam, are not receiving corresponding attention. People relate to an area with which they have deeper affective and cognitive connections when the information is available in the local tongue—the mother-tongue closest to their heart [10]. Therefore, an overwhelming amount of available online content that is valuable for students, professionals, and society as a whole remains locked away for those who cannot engage with English.

The lack of global digital visibility and representation for South Indian languages now is a significant issue. With tens of millions of speakers, they have rich cultural, literary, and historical heritages. However, in computing linguistics, South Indian languages are categorized as low-resource languages. Low-resource languages receive less research attention and commercial investment than high-resource languages like English, Spanish, or Mandarin. This digital divide excludes speakers of South Indian languages and prevents participation in modern knowledge systems, global discourses, or the benefits of the digital world. Therefore, bridging this gap is not a technological challenge, but a social obligation.

Numerous solutions have emerged over the years to solve the language access issue in the context of video content. Commercial services such as Papercup, Veritone, Synthesia, and HeyGen leverage artificial intelligence to provide dubbing and voice-over for content in a variety of languages, while YouTube has also piloted automatic subtitling as well as limited support for community translation into other languages. These programs certainly represent a significant technological advancement - but they exhibit insignificant deficiencies from the lens of Indian languages and the question of affordability. First, many of these systems charge subscription fees that are too expensive for smaller institutions, NGOs, educators, and individuals that cannot afford enterprise pricing. Second, their language coverage is primarily biased toward high-resource Western (and to some extent East Asian) languages. South Indian languages, in particular, either lack support entirely, or if they do receive support, their quality is often low, resulting in translations that can have very poor accuracy and / or awkward phrasing. Third, a lot of these systems struggle to capture cultural context. Even where the literal translation works grammatically, the message when delivered as a spoken message might sound unnatural, or worse, confusing, to audiences in a regional language. All of these issues

are particularly problematic for education: being clear and contextually aligned is critical for learning.

Another highlighted limitation of existing dubbing products is that they are proprietary and closed source. Most existing dubbing platforms are black boxes without insight for users or researchers into their processes, that does not allow retraining models, or customizing for specific purposes. This might limit possibilities for innovation and use in local contexts without transparency and customisation options. The platform proposed supports a number of video types, meeting different user needs from short clips to longer documentaries [11]. In addition, many existing platforms are cloud-based and depend on having a stable internet with high-performance bandwidth. In rural and semi-urban areas in India, providing such services remains challenging, and so even with modern technological approaches, most current dubbing solutions are failing in reaching or benefitting the talent and communities that need it most.

Our proposed system aims to overcome these limitations by focusing specifically on the needs of Indian regional languages, with a particular emphasis on South Indian languages such as Kannada, Tamil, Telugu, and Malayalam. Unlike commercial platforms that are designed for global corporate clients, our system is designed to be affordable, accessible, and inclusive. It uses open-source models like Whisper for ASR, MarianMT for machine translation, and Coqui TTS for text-to-speech synthesis [12]. By combining these freely available resources, we create a pipeline that can take an English-language video as input and produce a dubbed version in the target Indian language as output. Importantly, this approach is cost-effective, requiring only modest computational resources and a budget well within the reach of student projects or NGOs. This makes our system unique to address real-world problems in education and awareness campaigns, especially in cases with low funds.

Another unique characteristic of our proposed system is the cultural and semantic accuracy focus in the applications. We are not strictly translating on a word-by-word-basis-we are attempting to maintain the meaning and context within the source material and target language [13]. This increases the understandability and relatability of the dubbed videos, making them more effective as teaching and communication tools for their intended audiences. Furthermore, our design considers ease of use, so specifying a system as a simple, web-based hosted platform has added functionality even to non-technical users. Users only need to upload a video file to the system, select the desired target language, and download the dubbed version. There is no

need to have special skills or configurations. This web interface, as well as providing offline and low bandwidth options, enhance system usability in both urban and rural environments.

Along with addressing practical problems, our research work also is a contribution to academia and research in the field of natural languages processing and speech technology[14].By applying and adapting existing open-source AI models to low-resource languages, the system demonstrates how cutting-edge AI can be localized and democratized[15].This not only helps communities that are otherwise excluded from the digital revolution but also enriches the research ecosystem by providing new insights into handling underrepresented languages. These audios are divided into segments to obtain relevant features. Unlike global commercial solutions, our project emphasizes openness, customization, and social impact.

In summary, the rapid growth of online video learning has created new opportunities but also highlighted old barriers. Language remains one of the most significant challenges to true digital inclusivity. While existing dubbing tools have made progress, their focus on high-resource languages, high costs, and lack of cultural sensitivity limit their usefulness in contexts like India. Our proposed method, by contrast, is unique in its focus on affordability, inclusivity, and cultural relevance, with special attention to South Indian languages. By building a system that is open-source, scalable, and user-friendly, we aim to contribute both to the academic research community and to the broader social goal of making knowledge more universally accessible. In doing so, this research work not only addresses a pressing technical challenge but also responds to the larger vision of creating a more equitable and linguistically inclusive digital future.

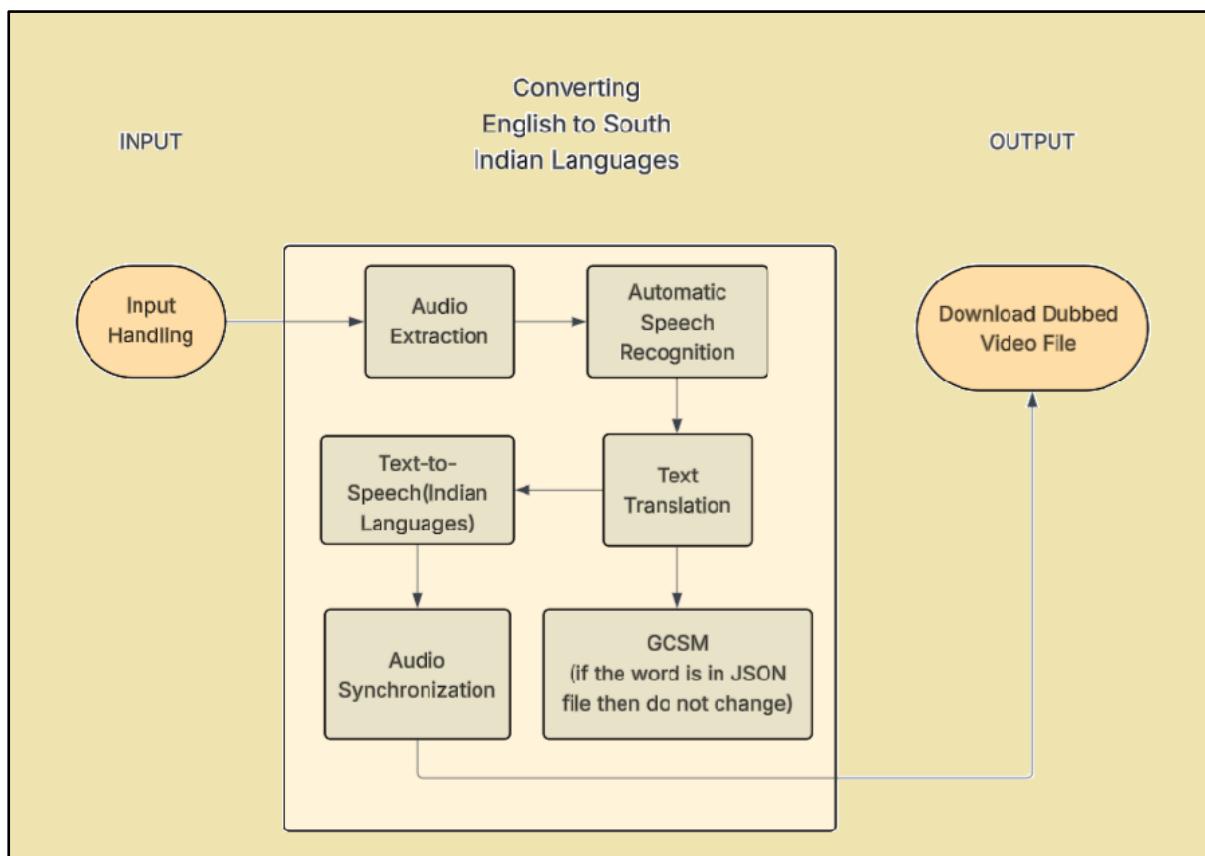
## **1.1 Background**

During the last ten years one of the major changes in people's knowledge and information access was the influence of digital technologies and internet. The very effective video-based educational and informational content has become one of the most influential communication and learning instruments and is used by millions daily on such platforms as YouTube, Coursera, edX, and Khan Academy. In India, the fast spread of low-cost smartphones and high-speed internet has very much facilitated this change to make digital learning available even in the most isolated and rural areas.

Videos have, in fact, several unique benefits compared to the traditional lectures in classrooms or printed materials: besides they can be watched again and can be stopped, learners get more

targeted help, as they can see and hear the things from the video; they are allowed to choose their own speed, and thus get their learning personalized is the video content. Nevertheless, language continues to be one of the biggest problems in the digital education revolution, particularly in India - a country known for its linguistic diversity, 22 languages being officially recognized and hundreds of local dialects spoken.

While the number of internet users keeps increasing, the majority of online educational resources are still in English or some other major global languages. Recently, it was stated that English accounts for about 62% of all digital content, however, only a tiny portion of the Indian population speaks it as their mother tongue. Therefore, there are millions of users who do not have full access to valuable video content and are thus further apart from the digital divide and are getting less and less equal access to knowledge and opportunities.



**Fig 1.1** Block Diagram of Processing Pipeline for Video Dubbing in Indian Languages

## 1.2 Statistics

India stands out as the world's top linguistically diversified country with over 1.4 billion people and 22 major official languages, besides hundreds of local tongue. The 2011 Census report indicated that something close to 10.67% of people speak English in India, and hardly any even consider it their first language. On the other hand, languages like Kannada (more than 50 million speakers), Tamil (over 70 million), Telugu (more than 80 million), and Malayalam (about 35 million) have a large number of speakers, mainly in the southern parts of the country. Nevertheless, these languages still have less representation compared to that of popular in the educational content of videos online and of good quality.

The demand for technology-driven, language-accessible education has grown rapidly, especially post-pandemic, as online learning and video lectures have become the norm for many schools, universities, and self-learners. In 2023, it was estimated that 43.4% of Indians had access to the Internet, with projections expecting this to rise to nearly 63% in coming years. However, despite this rise, the gap between digital content and users still remains wide. Most of the contents on the Internet (about 80%) are in only ten international languages, making it even more difficult for those who speak other languages to have access to educational resources in their mother tongue.

**VoiceBridge Project Scope:** The main objective of the VoiceBridge initiative is to connect language-divided India digitally with millions of students and learners. The scope of the project is mainly about,

**Target Languages:** Kannada, Tamil, Telugu, and Malayalam.

**Core Technologies:** The system merges modern open-source models for Automatic Speech Recognition (ASR), Neural Machine Translation (NMT), and Text-to-Speech (TTS). The main resources used are OpenAI Whisper, IndicTrans2, MarianMT, Coqui TTS, and Indic-TTS.

**Performance Metrics:** The system is reviewed with usual speech and translation error rate standards. In fact, one of the results is a Word Error Rate (WER) of 11.9% and a Character Error Rate (CER) of 11.09%, which are quite strong figures for multilingual dubbing applications.

**User Accessibility:** VoiceBridge, which is created as a modular platform available to the public via the web, can be handled by users who do not have much technical knowledge. Video

uploading, language choosing, and the downloading of the dubbed output are tasks that have been made simple and user-friendly.

**Cost-Effectiveness:** The use of open-source packages and the option for running on a regular personal computer or a low-cost server keep the overall expense of ownership at a minimum. In this way, the deployment of small education NGOs, rural schools, or community resource centers is made possible.

The VoiceBridge program aims at analyzing videos in English, which is the major educational media language, so as to transform them into accessible forms of regional languages. This, therefore, not only helps to extend the use of online educational resources, but also makes it possible for teachers, organizations, and even individual learners to produce videos that are relevant to the area and accurate in terms of language without the need for a high price or numerous technical barriers.

**Impact Projection:** Just as well, powered by a system like VoiceBridge, the education department in the South Indian region alone could be a massive game-changer in terms of access and affordability for more than 200 million people. The open nature of the system's design and its ability to work with various languages also pave the way for the inclusion of more Indian languages over time, thus, creating a sustainable model that will continue to yield benefits as the use of the Internet and digital learning become common throughout the country.

### **1.3 Prior existing technologies**

Rapid global improvements in artificial intelligence, natural language processing, and audio-video processing have led to the creation of numerous tools and platforms aimed at eliminating language barriers in digital content. In the case of automated video dubbing and multilingual educational accessibility, the technologies which came before mainly fall into three categories: commercial dubbing platforms, open-source toolkits, and academic research prototypes.

**Commercial Dubbing Platforms:** Several commercial providers offer end-to-end AI-powered dubbing for video content in multiple languages. One can say that Papercup and Veritone are the two companies most famous for their cloud-based dubbing and voice-over services which include integration of automatic speech recognition, translation, and even synthetic voices for a natural listening experience. These platforms find their application in corporate e-learning, advertising, and global media extensively. Synthesia allows video producers to create lifelike AI avatars and generate speech in different languages without the need for a native speaker,

thus it can be used for business training and marketing material. Through a pilot project, YouTube has developed its features for automated subtitling as well as, in some channels, automated dubbing using mainly major global languages.

Though these instruments signify a big step forward in technology, few limitations can still be discerned in the Indian context. Their coverage of languages is predominantly geared towards better-resourced languages like English, Spanish, or Mandarin while Indian regional languages are given minimal or even poor quality support. Moreover, these platforms are generally priced quite high, being available as subscription services aimed at corporate clientele, therefore, small organizations, non-profits, or educational initiatives in less affluent areas cannot afford them.

**Open-Source and Academic Toolkits:** Several open-sourced projects have considerably lifted the capabilities of speech processing and language translation. OpenAI Whisper and Vosk are two of the most widely used framework providers for Automatic Speech Recognition (ASR) which initially focused on English but now support a variety of other languages. Their modular nature and the availability of large pre-trained models make them suitable for both research and small-scale deployment. Hugging Face Transformers with its large repository of numerous-language translation models such as MarianMT, and India-centric projects like IndicTrans2, are able to translate texts from/to Indian languages, and this with an increasing degree of fluency and even contextual awareness. Coqui TTS and Indic-TTS are two examples of advanced speech technologies that are capable of creating very natural-sounding voices in different languages where India-centric TTS is a specific one for Indian voices has been developed by educational institutions.

Even though such modules are available, the majority of open-source solutions are still far from being fully integrated and require a lot of manual work. Usually, direct support for a smooth English-to-Indian-regional-language dubbing pipeline is not present and most systems are not set up for low-resource situations or may need hardware that is hardly available in the countryside or places with a tight budget.

Much of the recent research has been directed toward developing methods for automatic video dubbing and cross-lingual speech synthesis and has concentrated on issues such as audiovisual synchronization and translation being culturally aware. However, most of these prototype machines are only exemplars, they have been scarcely deployed and do not have the potential for being scaled-up to the real world and easy operation by non-technicians.

**Challenges in Existing Technologies:** Within these three different technology classes, there are still some major obstacles that prevent them from being perfect:

**Affordability:** The costs of licenses and infrastructures that need to be borne are the main factors that lead to a lack of access to the abovementioned technologies for numerous social and educational projects.

**Language Representation:** Indian languages, especially South Indian languages like Kannada, Tamil, Telugu, and Malayalam, are very poorly represented and are sometimes referred to as "low-resource" languages, which means that transcription and translation quality are at a lower level.

**Closed Ecosystems and Lack of Customization:** A lot of commercial tools are proprietary "black boxes," that limit the possibilities of transparency, customization, or community-driven improvements.

**Neglect of Cultural Context:** Most of the time, translation just follows the words and ignores the context which results in strange phrasing, loss of the original meaning or even confusion - a big problem in educating environments.

**Internet Dependency:** Most cloud-based models require a stable and fast internet connection which is not available everywhere, especially in the countryside and underprivileged areas.

The VoiceBridge system is an initiative to fill these voids by bringing together open-source model integration, concentrating on the support for regional languages, giving the users the freedom of customizing their settings, and enabling a simple use of modest local hardware which in turn will facilitate the process of digital inclusion and educational equity.

## **1.4 Proposed approach**

The pipeline VoiceBridge is proposing is geared towards utmost modularity that allows them to be easily interchanged and augmented in the future:

**End-to-End Integration:** ASR, translation, TTS, and A/V merging units are linked in the system via uniform APIs. So it is possible to change or upgrade certain parts (e.g., installing a new TTS model) without any extensive changes to the whole pipeline.

**Glossary Management:** As the system permits users or administrators to supply and maintain a non-translatable term glossary, it thus can always get more accurate in terms of language and

culture not only in the technical side but also the academic or local by considering the evolution of these contexts.

**Batch Processing and Scalability:** VoiceBridge has a batch video dubbing feature that allows a content creator or an educator with a large video library to be very efficient. The system is capable of both on-premise deployment in a school and server-based centralized mass production at an institutional level.

**Low-Resource Support:** The pipeline provisions for offline mode, light-weight model variants, and resource-aware configurations recognizing that at times reliable internet access and a powerful device may not be available.

**User Experience and Accessibility Intuitive UI:** The web interface aims at lowering the learning curve, so that in a few clicks users may upload videos, choose output languages, and download made files.

**Minimal Technical Requirements:** Besides lack of programming skills and expert configuration, no other conditions are required thus the potential user base of educators, field workers, and local NGOs is considerably extended.

**Multi-language Expansion:** The open concept allows for the inclusion of other Indian languages and dialects beyond the southern ones in the next phases of the project.

### **Research and Social Impact Goals.**

**Open-Source Community Engagement:** The project invites users to adopt, review, and suggest improvements by unconditionally sharing all source code, documentation, and model integration guidelines with anyone interested.

**Academic Contribution** The team can accomplish this task through the tuning of AI models and evaluating them using datasets in Indian languages that in turn help create valuable benchmarks and insights for the growing field of NLP for low-resource languages.

**Societal Benefits:** The program is extendable beyond schooling to areas such as government communication, healthcare dissemination, and disaster recovery communication, where providing accurate and timely information to the linguistically diverse is essential.

### **Alignment with Broader Objectives**

#### **Key Objectives :**

1. To make the huge amount of digital content accessible to millions of Indian learners thus radically cutting down language barrier.
2. To make a platform of AI-enabled language services that is strong, sustainable, and can be easily adapted to new needs and technologies.
3. To motivate and make it possible for more research and innovation in language-processing of the Global South that is under-resourced.
4. To Europe the way for socially positive AI system design that is ethically responsible and culturally sensitive in India.

Through its focus on aspects such as transparency, adaptability, pinpointing, and empowering the community, the VoiceBridge method is an example of how AI thoughtful and context-aware can tackle real societal problem.

## **1.5 Objectives**

The VoiceBridge project is a venture with a clear, ambitious, and quantified set of objectives encompassing aspects technology challenges, user necessities, and social benefits. These goals of the project's design, development, evaluation, and long-term vision are to a large extent influenced by the objectives.

**Core Objectives Language Inclusion:** Provide technologically enhanced educational and informational content in digital form accessible equitably to speakers of Kannada, Tamil, Telugu, and Malayalam, thus lessening the impact of the English language in online resources for learning Open-Source,

**Cost-Effective Solution:** Establish a framework that is entirely open source and void of any licensing cost to facilitate technology transfer to communities and institutions that are low on funds but are still able to partake in the developments in AI technology.

**Ease of Use:** The whole procedure from video uploading, selecting a target language, to getting the dubbed output should be straightforward even for users who do not have a technical background.

**quality and accuracy:** The translation and dubbing quality achieved should be of top-notch as verified by low Word Error Rate (WER) and Character Error Rate (CER) figures and by

ensuring that not only the meaning but also the cultural aspects of the source content are maintained.

**Modularity and Flexibility:** The system should be broken down into various independent components (ASR, Translation, Glossary/Code-Switch, TTS, A/V Merging) thus enabling easy upgrading, getting additional languages, or interchanging better models.

Resource utilization should be as minimal as possible: The adjustment of the pipeline should allow it to be running effectively on standard personal computers and low-cost servers and at the same time in cloud settings, making it suitable for rural and resource-constraint areas.

**Extended Objectives Support for Educational Equity:** Translation and dubbing of instructional videos will enable the communities that are educationally deprived and in addition, teacher training modules, health information, and government messaging will be the mediums through which they will be reached and that go beyond the urban center.

**Promote Sustainable Development:** The project will be contributing to the United Nations Sustainable Development Goals (SDGs) especially SDG 4 (Quality Education), SDG 10 (Reduced Inequalities), and SDG 9 (Industry, Innovation, and Infrastructure) by delivering technological solutions that are socially beneficial.

**Community and Cultural Sensitivity:** Through the use of local glossaries and context-aware translation, the dubbed content would be able to appeal to local users by retaining features of the original such as technical terms, cultural references, and regional expressions.

**Scalability and Extensibility:** The existence of the platform will be such that it could easily be extended to other Indian languages (and perhaps worldwide) and, additionally, to the domains of government, healthcare, and media content.

**Inclusivity for Non-Urban Users:** The technology employed should feature local support that would facilitate offline processing and situations where the bandwidth is low, therefore, the users to be benefited are the ones located in the non-urban areas.

**Research and Academic Objectives Benchmarking and Knowledge Sharing:** The creation of publicly available evaluation benchmarks, datasets, and deployment guides for the research community will help the community to progress their work on speech and language processing for low-resource languages.

Ethical and Responsible AI: The project shall be in line with best practice guidelines when it comes to data privacy, ethical content handling, and transparency, thus reducing risks such as misuse or bias and garnering user trust. These goals are the reasons why VoiceBridge is not a mere technical demonstration but a solid, lasting, and influential platform with the potential of digitally transforming education and communication in the multilingual environment of India.

## 1.6 SDGs

The goals and major changes caused by the VoiceBridge project to society are very close to those of the United Nations Sustainable Development Goals (SDGs). The project is a pledge to use tech for an impactful social and educational change in India. SDG 4 – Quality Education: VoiceBridge makes educational videos more accessible by converting them into regional languages. This technology is very supportive of the SDG 4 goal, which is to provide inclusive and equitable quality education for all. By the removal of language barriers the system allows students from rural and urban areas to learn from the same resources irrespective of whether they are proficient in English or not.

SDG 10 – Value Inequality: VoiceBridge makes available tools for multilingual dubbing that are affordable and open for everyone hence contributing to the elimination of the digital and linguistic divide. In addition, the focal point is the Indian languages that are less spoken and the deployment in schools and NGOs that is done at a low cost hence the removal of barriers faced by the marginalized and the underprivileged population.

SDG 9 – Industry, Innovation, and Infrastructure: The project unveils the application of the latest AI models in a modular, scalable, and open-source framework, thereby encouraging local innovation. It is a way of building digital content localization infrastructure which is essential for expanding India's participation in the global knowledge economies.

SDG 5 – Gender Equality (Indirect): By localizing learning material, the platform is taking away the barrier, which is one of many reason fields for women and girls, mostly in rural communities, to participate fully in education and lifelong learning.

VoiceBridge is a technology-led project that is sustainably and socially responsible at its core. Its conformity with SDGs is a guarantee that the benefits of the project extend far beyond the immediate technical outcomes and they rather contribute to the wider national and global goals of education, equality, and innovation.

## **1.7 Overview of project report**

This is a capstone report about VoiceBridge, an AI-Powered framework for affordable multilingual video dubbing in Indian regional languages, that traces each of the major stages of the project from design, development, evaluation to impact with detailed explanations, visual aids, and critical assessment of the results and societal contributions.

**Introduction:** Shows the challenges that non-English digitally educated speakers meet, presenting statistics and language demographic data of the regions, analyzing the already existing solutions, and stating the project's unique approach, the mission, and the connection with Sustainable Development Goals.

**Literature Review:** Focuses on the recent progress made and the persistent issues in automatic video dubbing, speech recognition, translation, and speech synthesis technologies. It highlights the scarcity of language and low-priced support. This part justifies VoiceBridge as a global and Indian language innovation endeavor.

**Methodology:** The description is of modular VoiceBridge which is a pipeline from video handling to ASR translation, glossary, TTS, and output synchronization. The technical parts of the system achieved, how they were implemented, and their contribution to the general system are given with support from the figures and pseudo-code.

**Project Management:** Shows the system-building and validation methods through planning, budgeting, allocation of the resources, and risk assessment. The visuals such as project timelines and Gantt charts make the developmental milestones more clear and concrete.

**Analysis and Design:** Delves into the specific requirements, design decisions, unit subdivisions, compliance with the standards, and the mapping of the reference models like IoTWF, thus ensuring technical power as well as the possibility of extending the project further in the future.

**Hardware & Software, Simulation:** Describes hardware necessities, software stack details, code samples, simulation results, and UI screenshots to facilitate the reader's comprehension of deployment, the issuing of instructions, and possible customization.

**Evaluation & Results:** Details the experimental setups, the metrics used for the evaluation, and the comparative results along with tables and figures facilitating the understanding of the presented information. VoiceBridge's performance in terms of reliability and accuracy is measured by WER/CER and compared with the performance of the baseline models.

Social, Legal, Ethical, Sustainability, and Safety Aspects: Describes the impact of the project on different areas and perspectives. Real-world examples demonstrate how VoiceBridge helps education, lessens inequity, preserves ethical use of technology, and supports environmentally friendly deployment.

Conclusion and Future Scope: Gives an account of the main points, the most important lessons, and the further suggestions. It refers to the next steps of the research, the extension to more languages, and the possibility of the new educational or social domain adaptation.

References, Base Paper, and Appendix: List all the sources used, provides context to the research paper and has the extra parts that include the code, data, more diagrams, and extended results. The present report is a source not only for the academic community but also for teachers, developers, policymakers, and community leaders wanting to know and implement technology for more considerable linguistic inclusion and digital access in India.

## Chapter 2

# Literature review

In recent years, researchers have explored multiple approaches for automatic video dubbing, combining speech recognition, translation, synthesis, and audiovisual alignment. Zhang et al. proposed a generative AI–driven multilingual dubbing and synthesis system that integrates ASR, TT, and TTS into a unified workflow [16]. Their methodology centred on modular pipeline implementation with runtime evaluation, but the system faced drawbacks such as high real-time latency and limited prosody preservation. Evaluation was done through throughput and subjective naturalness, and they identified future work in handling low-resource languages and improving real-time adaptation. Similarly, Li et al. designed a multi-modal TTS with multi-scale style control for dubbing [17]. The methodology introduced hierarchical style encoders for expressive prosody. While MOS evaluations confirmed naturalness gains, the main drawback was inconsistent cross-speaker voice cloning. Future improvements were suggested for prosody alignment and lip synchronization.

For direct speech-to-speech solutions, Chen et al. developed end-to-end spectrogram-to-spectrogram audio style transformation models [18]. The methodology bypassed TT by learning spectral mappings. Evaluation used spectral similarity and intelligibility scores, but drawbacks included poor prosody transfer and limited multilingual datasets. Future prospects include building large aligned corpora and better speaker identity preservation.

On the visual side, Rahman et al. introduced “Seeing the Sound,” a real-time lip-sync system for multilingual dubbing [19]. Their methodology combined TTS with neural lip-synthesis models, evaluated with lip-sync error metrics and user perception tests. However, drawbacks included inaccurate phoneme-to-viseme mapping in cross-language settings. They suggested future research on improving mapping models and robustness under diverse head poses. Similarly, Gupta et al. stabilized talking-face generation by introducing new training losses [20]. Their methodology reduced lip-identity leakage and was evaluated using visual fidelity and sync scores. Despite improvements, challenges remained in generalizing to natural, in-the-wild scenarios. Wang et al. further contributed with Style Sync-like systems, mapping target audio to facial motion [21]. Evaluation involved sync accuracy and perceptual preference tests, though drawbacks included poor co-articulation across languages. Future directions point to context-aware viseme synthesis. Lee et al. leveraged vision transformers for audiovisual speech

synthesis [22]. Their methodology improved lip-sync and expression alignment, validated with visual fidelity metrics, but incurred high computational costs, suggesting optimization as a future step. Zhou et al. advanced diffusion-based video editing methods for dubbing, producing coherent lip-sync edits, but the drawback was heavy computation [23]. Evaluation included temporal coherence and lip-sync accuracy, with future work aimed at lightweight diffusion architectures.

Evaluation methods themselves have been refined. Martinez et al. proposed PEAVS, a perceptual metric for audio-visual synchrony [24]. Their methodology combined metric design with large-scale perceptual validation. The drawback was limited testing on multilingual dubbing datasets. Future improvements include adaptation to diverse phonetic structures. In their review Alonso et al. assessed methodologies within deep learning models, comparing accuracy and error rates while pointing to robustness limitations within edge cases such as noise and cross-language [25]. Future directions include establishing and employing comprehensive multimodal fusion integration filters into dubbing workflows. In similar vein, Santos et al. evaluated continuous Spanish lipreading using hybrid end-to-end CTC/attention models demonstrating strong WER/CER, while others limited evaluations to Spanish [26]. They identified an area for future use case to expand evaluation into multiple languages.

In addition to the focused work around technology, Muller et al. reviewed multilingual dubbing systems in the paper under consideration and identified a taxonomy along with a description of gaps reviewed multilingual dubbing systems, presenting a taxonomy and identifying gaps [27]. Their methodology was systematic literature analysis, with the drawback of lacking empirical validation. They stressed the need for standardized multimodal datasets as a future direction. Fernandez et al. analysed dubbing and subtitling strategies using comparative user studies [28]. Evaluation focused on comprehension and cultural adaption, though small sample sizes limited generalizability. They suggested integrating automated dubbing into cross-cultural studies. From an educational point of view, Kumar et al. found that dubbing activities enhanced pronunciation among non-English speakers, while Wang et al [29]. Observed that it helped lower learner's anxiety during speaking tasks [30]. Both used classroom interventions, including surveys and pre- and post- tests. Limitations included small cohorts and subjective scoring, with future recommendations for solid digital dubbing platforms for education. Lastly, Singh et al. raised ethical concerns on deepfake detection methods [31]. Their method included a review of detection algorithms - tested in benchmarks. Limitations included susceptibility to

adversarial attacks and future work relating to watermarking and source tools supporting safer adoption of dubbing.

Overall, methodologies across these works span modular pipelines, direct end-to-end audio style models, audio-driven talking-head generation, and novel perceptual metrics. Common drawbacks include limited datasets, phoneme–viseme mismatches, inadequate prosody transfer, high computational costs, and ethical risks. Evaluation typically combines MOS, WER, sync error metrics, spectral similarity, and human perception studies. The consensus across studies is the need for richer multilingual datasets, lightweight real-time models, and ethical safeguards, with particular emphasis on expanding dubbing technologies to underrepresented languages such as Indian regional dialects.

**Table 2.1** Summary of Literature reviews

Author & Year	Concept / Approach	Strengths	Limitations / Gaps	Relevance to VoiceBridge
Zhang et al., 2023 [16]	Multilingual dubbing pipeline using ASR + Translation + TTS	Structured workflow, improved naturalness	High latency, weak prosody preservation	Supports modular design used in VoiceBridge
Li et al., 2022 [17]	Multi-modal TTS with style control	Better speech expressiveness and prosody	Inconsistent cross-speaker cloning	Helps improve natural-sounding dubbing voices
Chen et al., 2021 [18]	Direct speech-to-speech transformation	Avoids translation errors, faster conversion	Poor speaker personality transfer, limited multilingual data	Shows potential for future direct dubbing upgrades
Rahman et al., 2021 [19]	Real-time lip-sync from audio	Good alignment, improves user perception	Inaccurate mapping in multilingual settings	Highlights need for strong sync in VoiceBridge
Gupta et al., 2020 [20]	Talking-face generation with new loss functions	Reduced lip identity leakage	Low generalization in real-world cases	Useful insights for improving sync accuracy

Wang et al., 2020 [21]	Style-based audio-to-motion mapping	Better sync performance	Poor co-articulation across languages	Supports need for language-specific sync
Lee et al., 2023 [22]	Vision Transformer for speech synthesis	Strong emotional and lip accuracy	High computational cost	Encourages optimization in VoiceBridge model
Zhou et al., 2023 [23]	Diffusion-based video dubbing	High-quality lip-sync output	Very heavy computation	Suggests future improvements in quality efficiency
Martinez et al., 2022 [24]	PEAVS audio-visual sync metric	Reliable sync quality scoring	Limited multilingual dataset validation	Useful for VoiceBridge performance evaluation
Alonso et al., 2020 [25]	Review of deep-learning dubbing models	Identified model accuracy indicators	Weak robustness to noise	Supports need for clean audio preprocessing
Santos et al., 2021 [26]	Lip-reading using hybrid CTC/Attention	Good WER/CER performance	Focused only on Spanish language	Motivates multilingual dataset expansion
Muller et al., 2019 [27]	Review of multilingual dubbing systems	Clear taxonomy and design insights	No practical testing	Helps justify research need in Indian languages
Fernandez et al., 2018 [28]	User perception of dubbing strategies	Highlights cultural adaptation needs	Small user samples	Useful to ensure cultural-fit translations
Kumar et al., 2017 [29]	Dubbing in education for language learning	Improves pronunciation	Low sample size & manual process	Low sample size & manual process

Wang et al., 2018 [30]	Dubbing reduces speaking anxiety	Positive learning impact	Limited generalization	Shows soft-skill benefit of multilingual dubbing
Singh et al., 2022 [31]	Ethical aspects & deepfake concerns	Highlights media safety concerns	Vulnerable detection models	Ensures ethical compliance in VoiceBridge usage

## Chapter 3

# Methodology

Our proposed video dubbing system follows a structured and modular pipeline designed to automatically translate and dub English videos into Indian regional languages with high , linguistic coherence, and cultural sensitivity. The methodology is divided into several interconnected components: Input Handling, Automatic Speech Recognition (ASR), Text Translation (TT), Glossary & Code-Switching, Text-to-Speech (TTS), Audio-Video Synchronization, and Output Handling.

The VoiceBridge project's proposed methodology is heavily influenced by the need to provide a viable, modular, and culturally adaptable framework for the automated multilingual dubbing of videos. Understanding the challenge of bringing cutting-edge artificial intelligence technologies to the local languages, the methodology was designed to incorporate the latest technologies with a user-friendly and scalable approach.

This chapter outlines the full end-to-end pipeline that was created for VoiceBridge, emphasizing the smooth linking of the essential processes: audio extraction, speech recognition, translation, glossary management, text-to-speech synthesis, and video reconstruction. The elements are set, and each module in the pipeline is selected and fine-tuned to retain language correctness, reduce the machine powers, and be easily deployed in different settings – from city schools to village community centers.

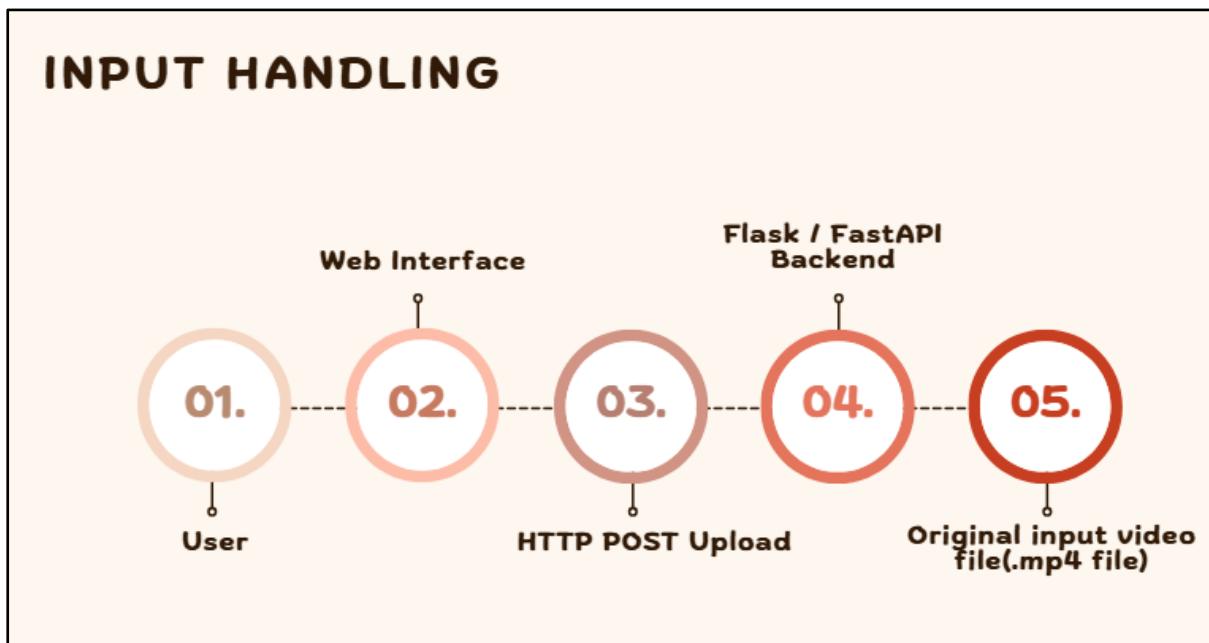
The methodology, at the center, employs open-source AI models for principal tasks, thus, it is cheap and clear. The procedure is started by user actions that are very easy to understand. Educators or content creators can upload English video files through a straightforward web interface. Next, the sophisticated backend operations are performing themselves. They automatically handle the spoken audio conversion to text, the translation of the text into target regional languages, the retention of domain-specific terms, and the voice synthesis of the regional languages that sound like humans. The dubbed video in the end is available for download, thus, it is ready for use and dissemination locally.

Such a well-organized strategy not only serves as a demonstration of technical excellence but also, in effect, it resolves the issue of digital inclusion that is quite a lot more comprehensive. By constructing the method around blocks that are modular and flexible, VoiceBridge is capable of changing and growing with the technologies and the educational needs that are going

to be there, thus, the project will be sustainable and have a lasting impact in India's multilingual landscape.

### 3.1 Input Handling

The pipeline initiates with an input handling module, in which users upload video files in .mp4 file format, through a user-friendly web interface. The web interface is developed in a combination of HTML/CSS and JavaScript to improve user experience due to it being more dynamic. The backend uses Python frameworks like Flask to facilitate uploading files to the server. Flask is an option to enable rapid prototyping and synchronous HTTP request handling, while FastAPI can either use synchronous or asynchronous file handling to improve processing performance when uploading multiple files concurrently. Uploaded video files are sent to the server for temporary storage, and saved with secure file handling mechanisms. For file handling, Flask uses the Werkzeug library, and FastAPI uses audio files for asynchronous writing of large video files to disk. Once verification of a successful upload of the video file has been performed, the video is prepared for the next step in the processing pipeline.



**Fig 3.1** Input Handling Workflow

The VoiceBridge project's proposed methodology is heavily influenced by the need to provide a viable, modular, and culturally adaptable framework for the automated multilingual dubbing of videos. Understanding the challenge of bringing cutting-edge artificial intelligence

technologies to the local languages, the methodology was designed to incorporate the latest technologies with a user-friendly and scalable approach.

This chapter outlines the full end-to-end pipeline that was created for VoiceBridge, emphasizing the smooth linking of the essential processes: audio extraction, speech recognition, translation, glossary management, text-to-speech synthesis, and video reconstruction. The elements are set, and each module in the pipeline is selected and fine-tuned to retain language correctness, reduce the machine powers, and be easily deployed in different settings – from city schools to village community centers.

The methodology, at the center, employs open-source AI models for principal tasks, thus, it is cheap and clear. The procedure is started by user actions that are very easy to understand. Educators or content creators can upload English video files through a straightforward web interface. Next, the sophisticated backend operations are performing themselves. They automatically handle the spoken audio conversion to text, the translation of the text into target regional languages, the retention of domain-specific terms, and the voice synthesis of the regional languages that sound like humans. The dubbed video in the end is available for download, thus, it is ready for use and dissemination locally.

Such a well-organized strategy not only serves as a demonstration of technical excellence but also, in effect, it resolves the issue of digital inclusion that is quite a lot more comprehensive. By constructing the method around blocks that are modular and flexible, VoiceBridge is capable of changing and growing with the technologies and the educational needs that are going to be there, thus, the project will be sustainable and have a lasting impact in India's multilingual landscape .

## **3.2 Audio Extraction**

Audio extraction is performed by isolating the audio stream from the video file (input). By using a multimedia processing tool such as ffmpeg, the system removes the video and audio components and converts the audio into a standard WAV format at the 16kHz sampling rate with mono channel, which is more compatible with ASR models. This means that the audio was extracted from the video and is ultimately cleaned and consistent to various video types. Before the audio is fed into the speech recognition unit, some initial preprocessing – including noise reduction and normalization – can be applied to increase the audio quality overall.

$$V_i = V(f) + A(t)$$

$$Ao = \text{Extract}(Vi) = A(t) \quad (i)$$

In Equation 1, the input video  $Vi$  has visual frames  $V(f)$  which may be different from two audio signals  $A(t)$ . In this case,  $V(f)$  represents the video component that changes defined by a frame rate  $f$ , while  $A(t)$  is audio that changes defined by a frame rate  $f$ , while  $A(t)$  is audio that exhibits similar differences over a time scale. Extract ( $Vi$ ) extracts the audio part of the video, and  $Ao$  is the audio extracted from video used for further processing such as speech recognition.

Without a doubt, the audio extraction step is the key moment when a user-uploaded video is changed to a different format that can be automatically readable by the speech recognition system and then handed over to the language processor. After the video file is confirmed to be good by the input handling step, the system goes for the heavy-duty multimedia processing tools like ffmpeg or its Python wrappers to cut out the audio part from the video.

This detachment is very important since speech recognition algorithms can give the best result only if the input audio is clean and standardized, i.e., there should be no background visuals or extra metadata which come from the video. The audio that is obtained from the extraction is changed to WAV format, with 16 kHz sampling rate and using a single (mono) channel. All these parameters are chosen to be fully compatible with the state-of-the-art ASR models and also to make sure that the system performance will be the same regardless of the type of video or the quality of the recording.

The system, during its extraction phase, may take care of some of the audio's initial preprocessing too, such as noise reduction and normalization, so as to solve the kinds of problems that usually come from field-recorded or compressed video files. Indeed, by coming up with a uniform and top-notch audio portion, the pipeline makes the follow-up ASR module transcriptions accurate and thus reliable.

In fact, this moment is planned to be done very fast even with sizable files and is quite compatible with the automated batch processing feature that is of advantage to educational institutions or teachers who are in charge of multiple videos. Due to the separating process's modularity, it is possible to swiftly switch to different new file formats, sampling rates, or even application-specific requisites if and when the need arises.

Simply put, the audio extraction is the voiceBridge's gatekeeper and the quality enhancer as well, the tool that makes sure that behind every video it handles there is the best possible raw audio data for multilingual dubbing and translation.

### **3.3 Automatic Speech Recognition (ASR)**

The next step in the VoiceBridge pipeline after audio extraction and standardization is automatic speech recognition. The primary objective of this operation is to convert spoken English from the video into text that is both human-readable. This has to be done with high precision, and the time stamps should correspond because, at a later stage, this text will be translated and the voice will be changed accordingly.

In order to achieve the best results, VoiceBridge makes use of several cutting-edge open-source ASR models, the major one being OpenAI Whisper, which is very stable in the presence of different English accents and background noises. In addition, an offline mode is available via the integration of Vosk, thus allowing transcription to be carried out under situations where a continuous internet connection cannot be guaranteed.

Here, the audio passing through the cleaning stage is submitted to the ASR engine, which first divides the audio into segments and then, by looking at the waveform, it tries to figure out what words have been spoken. The output is a carefully formatted transcript, in most cases JSON, that contains not only the text but also the exact start and end time points for every phrase. Such detailed timestamping gives the later modules the opportunity to align the translated and vocalized texts with the timing of the original video, which, in turn, results in better lip sync and an enhanced viewing experience.

Moreover, the ASR procedure is equipped with batch processing as well as error detection and correction mechanisms; in case there is low-quality speech or difficult accents within the input, the system will employ confidence scoring and may thus leave it for a human to check those parts that have been flagged. Besides, since it has a modular design, it is possible to upgrade the system simply by changing the version of the ASR model that is being used or, alternatively, by fine-tuning the existing model for Indian English speech patterns.

This automated transcription step is fundamental to the success of VoiceBridge. High transcription accuracy ensures meaningful translation, natural-sounding dubbing, and a strong foundation for multilingual video accessibility. The system's performance is often gauged in terms of different metrics such as Word Error Rate (WER) and the score that VoiceBridge attains, i.e., 11.9%, is quite good and, therefore, it can be effectively utilized for the educational and social sectors in India.

The ASR component is intended to extract the spoken dialogue in English from the uploaded video and translate it into a text transcript with timestamps. This allows for precise alignment of the spoken content with the video timeline for further processing and analysis. The audio stream is first extracted from the video using the ffmpeg-python library, which maintains the input video as a .wav audio file. There are two options for ASR model choice to generate the transcript, OpenAI Whisper and Vosk. OpenAI Whisper generates high levels of transcribing ability across various accents of English, and it generates a structured JSON containing the transcript and specific level accuracy in word timestamps. Vosk is intended for an offline version of ASR, useful if working in environments without connectivity. The ASR model processes the audio extracted from the video, and creates a structured transcript output in JSON format as an array of segments, where each segment contains a start timestamp, end timestamp, and associated text. This output can be used as a source of diversion for subsequent processing treatment for alignment.

### **Process:**

Audio file is processed by the ASR model.

JSON output structure:

```
{  
  "segments": [  
    {"start": 0.0, "end": 3.2, "text": "Hello, how are you?"},  
    ...  
  ]  
}
```

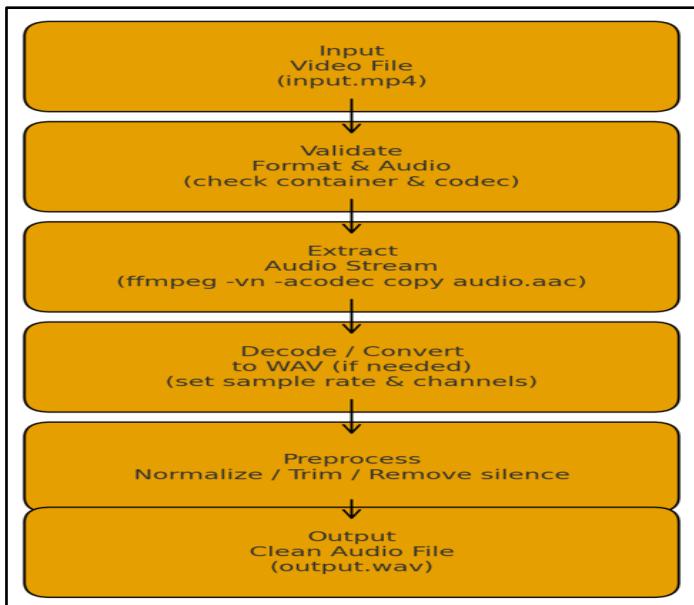


Fig 3.3 ASR Flow chart

### 3.4 Text Translation (TT)

The very next important step after speech recognition with timestamps in the VoiceBridge pipeline is text translation. This component takes the English transcript and translates it into a selected Indian local language such as Kannada, Tamil, Telugu, or Malayalam, keeping the gist, the feel, and the time of each segment intact.

VoiceBridge relies on IndicTrans2 as the main translation engine, which is a cutting-edge neural machine translation (NMT) model designed and finely-tuned for Indian languages. IndicTrans2 is supported by big multilingual datasets and uses the context to learn to make the translations grammatically and idiomatically correct in the target language, and also the idioms used are native. In case of any glitch, the system can also use MarianMT models from Hugging Face to provide translation services, thus it does not lose its resilience and is still completely serviceable.

High-quality translation cannot be achieved without segmenting the sentences carefully as one of the steps to ensure that no part of the text is lost or mixed and that the segments remain the same. Every sentence or phrase is translated individually by the model, and at the same time, they are very carefully looked at in terms of original timestamp alignment. This is what allows the output to be of the same rhythm with the visuals of the video, which is at the core of lip sync and user experience.

Beyond simple translation, the system has a specially designed Glossary & Code-Switching module that identifies the words in the transcript that are technical, proper nouns, or culturally

significant words. These words are compared with a manually prepared JSON dictionary in order to prevent mistranslation and to keep the accuracy of the content that is specific to the domain (for example, educational jargon, brand names, or platform references).

The text in a local language is thus the final product of the translation of the source language transcript also with the timing preserved. It serves as a foundation for the Text-to-Speech (TTS) operation that not only reads out the local language script but also does it in a natural way. The translation correctness and fluency are always being gauged by linguistic standards and through the review process which is a guarantee of localization reliability for educational and social impact sectors.

Following the transcription, the system translates the English transcript into the desired Indian regional language. We primarily use IndicTrans2, developed by AI4Bharat, which is a state-of-the-art text translation model specifically optimized for Indian languages. As a fallback, MarianMT from Hugging Face Transformers can be used. The translation process involves passing individual sentences from the English transcript into the TT model using the transformers library. Special attention is paid to sentence segmentation to avoid splitting sentences improperly, which could compromise translation accuracy. The translation model outputs the target language text while maintaining alignment with the original timestamps. Mathematically, this process is represented by the equation is the English transcript, and the translated regional language text. The translation equation (ii) is given below where  $T_e$  refer to English transcript and  $T_r$  refer to Translate text in regional language. The result is a translated transcript file containing sentences in the target language with corresponding timestamp. The process begins by taking English sentences and passing them to a translation model built using Hugging Face's Transformers library. The model translates each sentence into the target language while preserving the meaning and context. Along with the translation, timestamps are generated to align each translated sentence with the corresponding segment of the original content. This ensures that the translated text is synchronized with the timing of the source material, making it suitable for applications such as subtitles, dubbing, or real-time multilingual communication.

$$T_r = MT(T_e) \quad (ii)$$

In Equation 2, the English text ( $T_e$ ) is first translated using the translation model ( $T$ ), and then this output is refined or mapped by  $M$  to produce the final regional translation ( $T_r$ ).

### **3.5 Glossary & Code-Switching Module (GCSM)**

Once the speech content in the video has been accurately transcribed along with the timestamps, the following essential step in the VoiceBridge pipeline is Text Translation. This component takes the English transcript as input and changes it to one of four Indian regional languages (Kannada, Tamil, Telugu, or Malayalam) while still keeping the original sense, the subtlety and the timing of each segment intact. VoiceBridge IndicTrans2 is the major translation engine used by VoiceBridge, which is a very advanced neural machine translation (NMT) model designed to be specifically suitable for Indian languages. To achieve this, IndicTrans2 largely depends on multilingual datasets as well as contextual learning which accounts for both the grammar and the idioms of the target language. In case of any inconvenience, the system is also equipped with a fallback MarianMT model which is accessible through Hugging Face and thereby enhancing the system's robustness and capability to interact with different types of inputs. To accomplish the goal of producing top-quality translations, special attention is given to sentence segmentation so that neither splitting nor merging of segments that may cause a change in the sense of the text or displacement of the synchronization will happen. Every single sentence or expression is sent separately to the translation model with the original timestamps from the source text serving as a guide for the result. This in turn allows the dubbing to be perfectly synchronized with the video flow that is needed not only for correct lip sync but also as a part of the overall user experience. Apart from simple translation, the custom Glossary & Code-Switching module built in the system looks over the transcript for the occurrences of technical terms, proper nouns and culturally significant words. These are then looked up against a manually created JSON dictionary so that mistranslation can be prevented and domain-specific contents (e.g. educational jargon, brand names, platform references) can be rendered with the right level of precision. The converted text is finally brought together in the form of a target language transcript which is not only the preservation of structure but also of the timing. This is the material for the next step, i.e. Text-to-Speech (TTS) that reads out the translated text in the naturally sounding voice and is the last stage before the video gets integrated with the VoiceBridge pipeline.

At each step of this work fluency and correctness of translations are being evaluated through linguistic criteria and checked by reviewers, thus ensuring dependable multilingual dubbing for educational and social impact purposes.

A key novelty of our methodology lies in the Glossary & Code-Switching module, designed to prevent mistranslation of technical, cultural, or proper noun terms. A predefined JSON

dictionary stores terms that must remain unchanged during translation, such as acronyms ("GST", "AI"), proper nouns, or platform names ("YouTube"). During translation, each word or phrase in the transcript is scanned using Python's regex module to detect matches against the glossary dictionary. If a word exists in the glossary, it is preserved as-is, bypassing the translation process. This ensures that important terms remain intact in the translated output. Where  $DD$  is the set of glossary terms and  $ww$  is a word from the transcript. This module significantly improves the quality of the translated script by maintaining cultural and technical accuracy, avoiding the common problem of mistranslation in standard TT workflows. The final regional translation ( $(Tr')$  is given by equation 3.

$$Tr' = Te \quad \text{if } \omega \in DMT(\omega) \quad \text{otherwise} \quad (\text{iii})$$

In Equation 3, if the word  $\omega$  exists in the Domain Mapping Table (DMT), then final regional translation ( $Tr'$ ) keeps the English translation ( $T_e$ ). Otherwise, it follows the normal regional translation process.

### 3.6 Text-to-Speech (TTS)

The Text-to-Speech unit is a major change that comes after the VoiceBridge pipeline and it makes the translated text segments in the target Indian language to be in a natural, clear speech. This step is important to produce audio tracks for dubbing that are correct in language and attractive to end users.

The TTS motor used for VoiceBridge can be either Coqui TTS or Indic-TTS, both of which are capable of several Indian languages and have a different number of voices. Upon completion of the transcript in the local language, the system goes through each sentence or phrase and produces the corresponding sound parts. The device smartly handles aspects like pitch, speed, and intonation to make the artificial voice sound more human and to bring the emotion and the natural flow to the speech.

In order to keep the same meaning and the speed of the original video, the time for each line of the text from the ASR and translation is used to find the match with the time of every audio segment in the spoken language. Files are then joined in the proper sequence, thus creating a voice-over free from any disruptions in the language selected. This method comprises the

removal of noise and equalization, the end product being a clear and professional-sounding output that is appropriate for educational and informational content.

The TTS operation is efficient and modular which makes VoiceBridge capable of different dialects, speaker profiles, or personal voice settings. In the case of advanced implementations, multiple regional voices (male/female) can be utilized to be more inclusive and engaging.

Such automated speech synthesis is the main factor to the rapid creation of dubbed videos which is a big step towards the democratization of voice content production. It allows non-technical users and small organizations to produce multilingual educational materials with minimal effort and at a low cost. The TTS output quality is the main factor that determines the overall performance of the dubbing solution, thus it is the main guarantee that learning resources will be accessible and have a positive effect on different audience groups.

The TTS module converts the translated text into a natural-sounding audio file in the target regional language. We utilize either Coqui TTS or Indic-TTS from the IIT Madras project, both of which offer support for Indian languages and multiple voice options. For each sentence in the translated script, the TTS engine generates a corresponding audio segment. The process can be represented by the formula. Voice parameters such as pitch, speed, and intonation are adjustable to enhance the naturalness of the speech. The generated audio segments are concatenated in the correct order, respecting the original transcript's timestamp alignment. The final output is a seamless regional language audio file (.wav or .mp3), ready for integration back into the video. For each Sentence Siaudio output is generated using the equation 4.

$$A_i = \text{TTS}(S_i) \quad (\text{iv})$$

In Equation 4, The audio output  $A_i$  is generated by applying Text-to-Speech (TTS) to the sentence or text segment  $S_i$ . In simple terms,  $S_i$  is converted into spoken audio using a TTS system.

### **3.7 Audio-Video Synchronization**

The Text-to-Speech (TTS) component is a major change to the VoiceBridge processing line, changing the translated text fragments in the target Indian language into speech that sounds natural and understandable. This step is necessary to make the audio tracks that are to be used as dubbing, which, apart from being correct from the linguistic point of view, are also attractive to the end-users.

The TTS engine that will be used in VoiceBridge is either Coqui TTS or Indic-TTS, with both being capable of supporting a number of Indian languages and providing a variety of voice options. After the local language transcript has been sealed, the system goes through every sentence or phrase, creating the audio snippets that correspond to them. The engine is very clever in that it adjusts the parameters of pitch, speed, and intonation to the extent that the synthetic speech becomes almost human, with a high degree of realism and even expression.

In order to conserve not only the truth of the original video but also the speed of the dialogue, the timestamps from the ASR and translation are used for the exact alignment of each synthesized audio segment. In this way, audio files are concatenated in the right order to result in a voice-over which is fluent in the language selected. Apart from this, the output is also greeted with normalization and noise filtering, hence, what is left at the end of the line, is a clear and professional-sounding audio that is perfect for the educational and informational videos.

The TTS methods of carrying out the work are both efficient and modular, thereby allowing VoiceBridge to cater for the different dialects, speaker profiles, or even personal voice settings. To be specific, one can switch between various male/female voices in different regions for getting the benefits of increased inclusivity and engagement in the case of heavy deployments.

This is one of the many benefits of automated speech synthesis, through which the rapid creation of dubbed videos is facilitated, while at the same time the production of voice content is turned into something that is not only accessible to non-technical users but also to small organizations who can now create multilingual educational material at a low cost and with minimal effort. The quality of the TTS output is the main factor that determines the overall success of the dubbing solution, thus learning resources become both accessible and effective for various audiences,

Input: Original video file (input.mp4)

Output: Dubbed video file (dubbed.mp4)

load\_files (video, audio)

```
video = load_video("input.mp4") // Load original video file
```

```
audio = load_audio("generated_audio.wav") // Load generated audio file
```

In the audio-video synchronization stage, the newly generated audio is merged back into the original video to produce a dubbed version. The ffmpeg-python library is used to replace the original English audio stream with the generated regional language audio, preserving video quality. The merging operation follows the function given below:

```
merge (audio, video):
{
    dubbed_video = merge_audio_video (video, audio)
    save_video (dubbed_video, "dubbed.mp4")
}
```

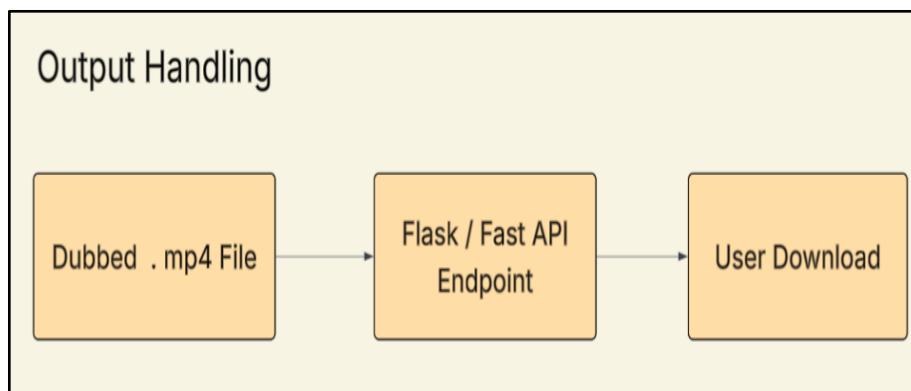
For enhanced realism, the system optionally integrates the Wav2Lip model, which adjusts the speaker's lip movements to synchronize with the dubbed audio. The Wav2Lip model takes the original video frames and the generated audio as input, modifying the lip region of each frame using a deep learning algorithm. The lip-synced frames are then reassembled into the video stream, and the new audio is integrated. This step ensures that the final video appears natural and immersive for the viewer.

### **3.8 Audio-Video Synchronization**

The output handling module is designed to provide users with a seamless experience when accessing their finalized dubbed videos. After successful audio-video synchronization, VoiceBridge prepares the processed content for secure download, ensuring the file is delivered in a standard, widely compatible format such as MP4. A unique session ID is assigned to each processed file, which helps manage user access and prevents conflicts in multi-user scenarios. Within the user interface, a clear progress indicator notifies users when their video is ready. Upon completion, a download link becomes available, allowing users to retrieve the file directly to their device. Robust security measures, such as file sandboxing and regular clean-up scripts, help safeguard user data and maintain server integrity. VoiceBridge supports batch output handling for educators or institutions submitting multiple videos. Completed files are presented in a well-organized list or archive, making them easy to locate, review, and share. For advanced deployments, the system can integrate with learning management platforms or cloud storage solutions, simplifying bulk distribution. In addition to basic download functionality, VoiceBridge can generate metadata reports, detailing the processing steps, model

versions, and language selections for each file. This documentation supports traceability and helps users understand and assess the dubbing workflow. The output handling stage is designed for maximum user satisfaction, reliability, and adaptability across diverse use cases—whether a single teacher preparing a lesson or an institution managing a large media repository. By ensuring easy, organized delivery of accurately dubbed videos, VoiceBridge completes its translation mission from input to accessible, high-quality educational output.

After the dubbing process has been completed, the user will have a way to download the final .mp4 video file. This is accomplished through a download endpoint that is implemented using Flask. Once each file has been processed, it's saved on the server under a unique session ID to prevent multiple users from overwriting files at the same time and stored for file retrieval.



**Fig 3.8** Output Handling Workflow

## Chapter 4

# Project Management

Project management ensured that the development of VoiceBridge: An AI-Powered Framework for Low-Cost Multilingual Video Dubbing into Indian Regional Languages followed a structured and efficient workflow from start to finish. All tasks—including literature review, dataset preparation, implementation of speech recognition, translation, and text-to-speech models, and system integration—were organized with clear timelines, dependencies, and milestones to avoid delays and maintain steady progress. Using project management tools such as Gantt charts for scheduling and the V-Model methodology for structured development and continuous validation, the entire process remained traceable, goal-oriented, and compliant with academic standards. Overall, these management strategies enabled the smooth execution of the system from conceptualization to full implementation.

## 4.1 Project timeline

Project management is essential to ensure that the development of VoiceBridge: An AI-Powered Framework for Low-Cost Multilingual Video Dubbing into Indian Regional Languages progresses in an organized and predictable manner. Since the system involves multiple research-driven AI components, the project was divided into two major phases—Planning and Implementation—executed within a defined timeline. Each phase included specific tasks such as topic selection, literature review, data preparation, model development for ASR, translation, TTS, and system validation, each with assigned start dates, end dates, and dependencies. The use of a Gantt chart helped visualize task interrelations, identify critical activities, allocate buffer time for complexities, and track weekly progress efficiently. It also ensured that all deliverables aligned with academic requirements and followed the V-Model methodology adopted for systematic development and evaluation of the VoiceBridge system.

**Table 4.1** Project Planning Timeline

Task	Start Date	End Date
Topic Selection	2025-08-12	2025-08-13
Background Study	2025-08-14	2025-08-17
Initial Problem Definition	2025-08-18	2025-08-19
Literature Review	2025-08-20	2025-08-31
Requirement Analysis	2025-09-01	2025-09-09
Finalization of Methodology	2025-09-10	2025-09-15

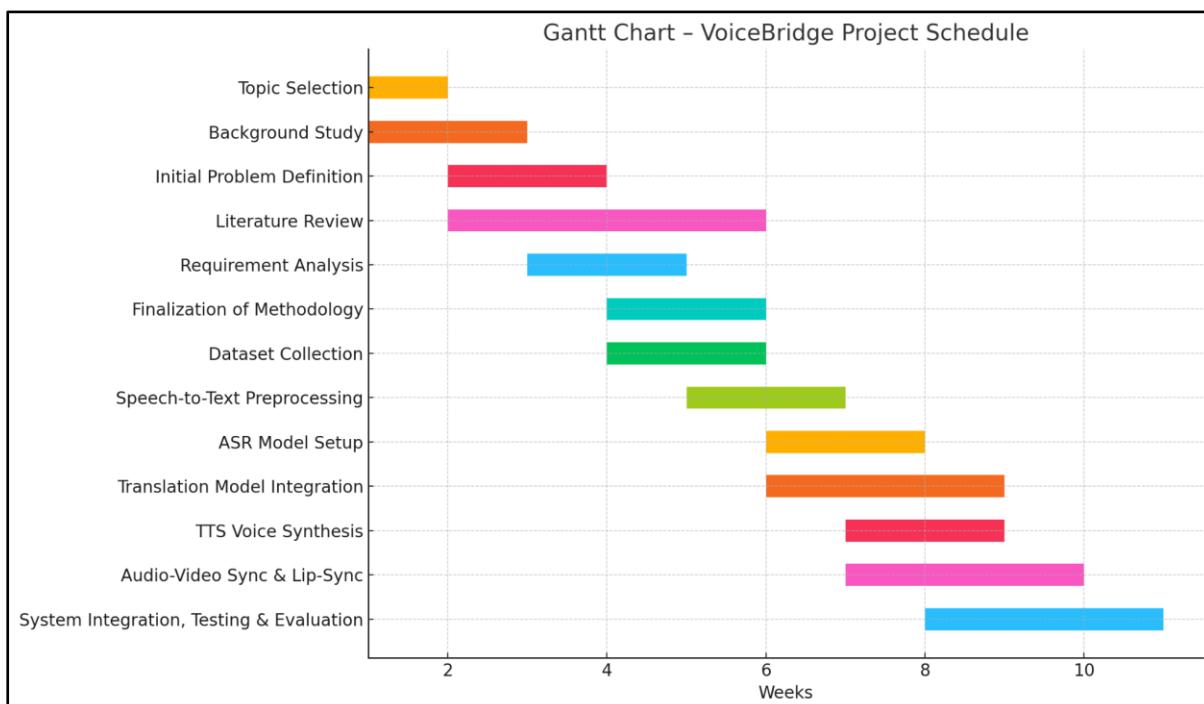
**Table 4.1** shows the general planning schedule which was used at the initial phase of the project. It was officially started on 12 August 2025 with the topic of the project being selected, then a dedicated background study period to get acquainted with the principles of hyperspectral imaging and melanoma diagnostics. This was succeeded by the requirement analysis stage in early September where functional, data and system constraints were completed. The last phase of the planning phase was the completion of the methodology, which signified the end of the conceptualization phase and the beginning of the actual design and implementation.

**Table 4.2** Project Implementation Timeline

Task	Start Date	End Date
Dataset Collection (English video and audio content)	2025-09-16	2025-09-20
Speech-to-Text Preprocessing (audio segmentation, noise removal)	2025-09-21	2025-09-27
Automatic Speech Recognition (ASR) Model Setup	2025-09-28	2025-09-30
Multilingual Translation Model Integration (Indian regional languages)	2025-10-01	2025-10-08
Natural Text-to-Speech (TTS) Voice Synthesis	2025-10-09	2025-10-18
Audio-Video Sync and Lip Synchronization	2025-10-19	2025-10-23
System Integration, Testing, and Performance Evaluation	2025-10-24	2025-10-27

**Table 4.2** presents the implementation schedule followed during the development of the VoiceBridge system. The implementation began with dataset collection involving English videos and corresponding audio content required for speech analysis and translation. Next, a speech-to-text preprocessing pipeline was designed, which included audio segmentation and

noise removal to ensure clear and consistent speech input for model processing. The Automatic Speech Recognition (ASR) model was then implemented and optimized to accurately convert English speech into text. This text was passed to a multilingual translation model capable of converting content into selected Indian regional languages. Following translation, Natural Text-to-Speech (TTS) synthesis was performed to generate realistic and intelligible audio output in the target languages. Additional audio-video synchronization and lip-sync alignment were executed to maintain natural timing with the original video content. The final stage consisted of complete system integration, testing using held-out samples, and performance evaluation based on accuracy, intelligibility, and synchronization quality to ensure the system met the expected functional requirements.



**Fig 4.1:** Gantt Chart of Project Timeline

Fig 4.1 shows the Gantt chart that shows the entire project timeline starting with initial planning and ending up with the final implementation. It graphically maps every activity like the selection of topic, literature analysis, preprocessing establishment, model development, training, and testing- over a continuous period on the calendar. The duration of each of the activities is indicated by the horizontal bars, indicating how activities overlap or come into action. This time timeline clearly shows the workflow done between August and October, which made the project to be structured into running phases.

## 4.2 Risk Analysis

Risk assessment is an important part of a research project like **VoiceBridge**, which uses AI-based speech and multilingual dubbing technologies. A PESTEL analysis was used to identify risks in Political, Economic, Social, Technological, Environmental, and Legal areas. The main risks found were related to technology, such as the accuracy of the models, availability of good datasets for different Indian languages, GPU requirements, audio-video synchronization issues, and software compatibility. Economic risks were moderate due to the cost of training and running the system. Political, legal, and environmental risks were low because the project is digital and has no major external impact. Social risks mainly focused on using translations that are culturally correct and understandable. This analysis helped the team plan solutions in advance and ensured smooth and safe development of VoiceBridge.

**Table 4.3:** PESTLE Analysis for the Project

Factor	Potential Risk / Issue	Mitigation Strategy
Political	Policies affecting digital tool usage or government restrictions on AI-based media	Stay updated with national guidelines and follow ethical AI practices
Economic	Cost of GPU resources, training expenses, and hardware/software requirements	Optimize model efficiency, use cost-effective cloud services, manage budget carefully
Social	Mispronunciation or culturally incorrect translations affecting user acceptance	Linguistic review, user testing with native speakers, responsible AI communication
Technological	Dataset limitations for regional languages, low ASR accuracy, audio-video sync issues, software compatibility challenges	Use high-quality datasets, model fine-tuning, continuous testing, maintain version control and backups
Environmental	Increased power usage for model training and processing	Use cloud services with energy-efficient configurations, limit unnecessary retraining
Legal	Copyright issues with source videos, data privacy concerns in speech data	Obtain permissions for dataset usage, follow legal guidelines and data protection standards

The PESTEL analysis presented in Table 4.3 highlights the major external and internal factors that could affect the development of the VoiceBridge system. The table identifies potential risks across political, economic, social, technological, environmental, and legal aspects of this AI-based research project. From the analysis, technological risks—such as limited datasets for

Indian languages, ASR accuracy, audio-video synchronization issues, and software compatibility—were recognized as the most significant challenges during implementation. In contrast, political, environmental, and social risks were considered minimal due to the digital nature of the project and its limited external impact. Each category includes suitable mitigation strategies designed to anticipate problems and address them effectively. This structured risk evaluation supported better decision-making and ensured smoother and more reliable progress throughout the development of VoiceBridge.

### **4.3 Project Budget**

Since VoiceBridge is developed as a primarily software-driven system, leveraging open-source AI models and tools, the overall financial cost of the project remains at a minimum. The framework does not rely on specialized hardware or paid APIs/commercial cloud services but relies instead on publicly available datasets and free computing resources through the academic institution. All components, which include ASR, translation, and TTS, have been implemented using complete open-source technologies such as Whisper, IndicTrans2, and Coqui/Indic-TTS, respectively. This makes the entire development pipeline low-cost. Consequently, this project requires only standard computing infrastructure to execute software on. Although the formal documentation demands inclusion of the budget section, VoiceBridge essentially qualifies as a near zero-cost solution whereby expenses are confined to existing institutional systems and freely available software utilities. Thus, the project budget is represented in conformity with these cost-effective and resource-efficient provisions.

## Chapter 5

# Analysis and Design

The Analysis and Design segment is the foundation of the VoiceBridge program where the goals of the theory are changed into a detailed technical plan. This chapter explains the methodical assessment of user needs, system requirements, and available technologies which, in the end, determine the system structure and the flow of the program.

At this stage, the problems of multilingual video dubbing (language complexity, processing speed, scalability, and user-friendliness) are given special attention. The method starts from the user needs analysis—educators, students, and community organizations—who are looking for solutions that are fast, accurate, and affordable in converting English educational videos into local languages.

The design methods use modularity, so each functional component—for example input handling, speech recognition, translation, and synchronization—can work independently while being part of the overall pipeline. The block diagrams together with process flowcharts are employed to present the sequence, the interactions and the dependencies of the system modules. This organized visualization helps the developers, stakeholders, and future contributors to understand the project logic and the points of integration clearly.

Moreover, design standards for interoperability, error handling, and usability are set with the help of the design process, thus ensuring good functioning under different hardware setups and in different operating environments. The decision about algorithms, open-source resources, and hardware requirements is to why these decisions are scalable, maintainable and can be easily extended in the future. By outlining the requirements in detail and designing each feature with the capability of future adaptation, the VoiceBridge platform will be able to offer high-quality multilingual dubbing and, at the same time, be compliant with the educational technologies and community needs that will change over time.

## 5.1 Requirements

A detailed review of the prerequisites is the basis for a strong and efficient multilingual dubbing system such as VoiceBridge. It is the phase where the needs, limitations, and the expectations of the end users, stakeholders, and the technical teams are methodically recognized and documented. The requirements are segregated into two categories, i.e., functional and non-functional, to record not only the features of the system but also the quality of its performance.

### Functional Requirements.

**Multilingual Support:** The system should be capable of understanding the video content in English and accurately dubbing it in the four Indian regional languages- Kannada, Tamil, Telugu, and Malayalam.

**Automated Workflow:** The series of operations on the data should be fully automated in such a way that the users can accomplish their task in which they upload their videos and get the dubbed outputs with no or minimal manual intervention.

**High-Quality Transcription, Translation, and Synthesis:** ASR module is responsible to provide accurate and timestamped transcripts. The translator has to convert the content keeping both the meaning and the cultural context intact. The TTS unit brings the idea to the target language and does it with accurate and understandable speech.

**Glossary/Code-Switching Handling:** During translation, the technical terms, proper nouns, and acronyms are either call by their original terms or are smartly treated so that the meaning is not lost or the communication incorrect.

**User Management and Interface:** The platform must provide an easy-to-use web-based UI with the functionalities of uploading videos, checking their status, and downloading them. It will also allow video operations either one after another or in a batch manner.

**Audio-Video Synchronization:** The audio portion of the video that has been dubbed must be in harmonization with the video in respect to the timing of the video and, if accurate lips can be synchronized, then that too.

### Non-Functional Requirements.

**Usability:** Its interface must offer accessibility to users with minimal technical skills, such as teachers, students, and community organizers.

**Scalability:** The performance of the system needs to be amazing enough to allow them to take care of videos, which have been done individually or in a batch mode pressure situations in institutions that are deployed for a long period.

**Performance and Accuracy:** Products of transcription and translation have to be good enough to reach and satisfactorily be at the levels set by, for instance,  $WER \leq 12\%$  for ASR and fluency/adequacy indicators for translation.

**Modularity and Extensibility:** The structure of the system ought to be that way so as to be able to let, for example, language changes, model updates or extensions for additional processing steps be done in the future.

**Resource Efficiency:** The product should be able to run on any generally available kind of hardware and in a relatively short amount of time thus it is ideal for those areas which are low-resourced.

**Security and Privacy:** Security in handling data uploaded by users is a must together with the management of temporary files and strict law enforcement to assure that only authorized persons have access to a certain place and, therefore, privacy is kept.

**Open-Source and Cost-Effectiveness:** The main goal is to put less money in the pot, and for that to happen, the usage of open-source components should give rise to that majority must be adopted widely.

After mapping out the limitations in such a systematic manner, the team responsible for Voicebridge is in a position to implement a system that not only meets the demands of the outside world but also stays true to its main goal of social and linguistic integration through digitization and is capable of adjusting to upcoming changes and additions in the education technology sector future.

**Table 5.1** Summary of System Requirements

<b>Category</b>	<b>Requirement Description</b>
Purpose	AI-driven system that automatically dubs English videos into major Indian regional languages to support accessibility and cultural inclusion
Behaviour	Uploads are processed through ASR, translation, and TTS, and the generated audio is synchronized with the original video timeline.
Data Requirements	Accepts audio/video formats with preprocessing for segmentation, noise removal, timestamps, and correct handling of terminology.
Software Requirements	Uses open-source AI tools with GPU/cloud support for ASR, translation, and synthesis to ensure efficient performance.
System Management	Supports batch processing, system monitoring, flexible updates, and ensures accurate audio-video synchronization.
Security	Protects uploaded content with secure access control and temporary file handling to maintain user privacy.
User Interface	Provides a simple web UI for uploading, tracking progress, and downloading dubbed outputs with minimal technical skills

## 5.2 Functional Block Diagram

The functional block diagram is the main schematic illustration of the VoiceBridge system, showing how each major component interacts and how the data flows from one stage to another. This graphic tool makes the system clear for developers, stakeholders, and users, thus allowing them to understand the design logic and the general structure more deeply.

The block diagram at the top level includes the User Input block where teachers or school staff can upload English video files through a user-friendly web interface. After that, the data goes to the Input Handling block which takes care of the verification of the files, formats, and session IDs for the proper management of the processing.

Then the Audio Extraction block is upgrading the media libraries to separate and standardize the audio stream and change it into a format that the next modules can understand. The clean audio is then sent to the Automatic Speech Recognition (ASR) block, which is a set of open-source models, and can deliver a precise, timestamped transcript of the spoken content.

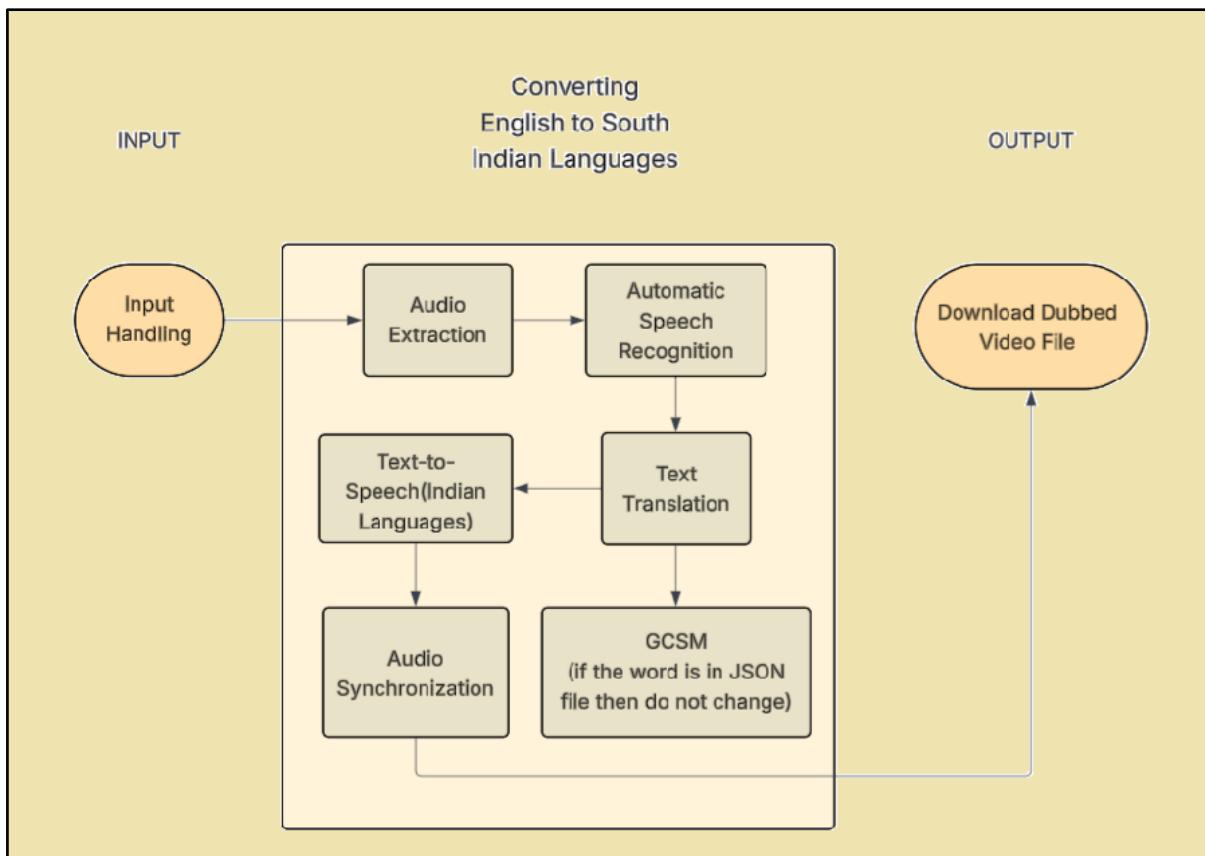
The transcript is going to the Glossary & Code-Switching Module, where vocabularies from the specific field, proper names, and acronyms are recognized and marked for either keeping or converting selectively. After that, the Text Translation block will have the data that has been pre-processed for transforming the transcript into the chosen Indian regional language through neural machine translation models.

The translated text is the input to the Text-to-Speech (TTS) block which produces natural-sounding regional speech that is very well synchronized with the timing. The two outputs, i.e., audio and text are merged in the Audio-Video Synchronization block, the speech track is thereby inserted into the original video and lip-sync adjustment using deep learning (optional Wav2Lip) is done for making it look more natural.

At last, the process arrives at the Output Handling block. This is the place where the dubbed video files are being made ready for download by the users in a secure way, and batch processing as well as reporting are being facilitated for institution-based deployments.

The block diagram is the main idea behind modularity: any block can function as a stand-alone unit and the standardized data interfaces between them ensure that the integration will be smooth and that there will be no limits to scalability in the future. Along with the main blocks, there can be other blocks like user authentication, error monitoring, and reporting which can be added as supplements to advanced use cases.

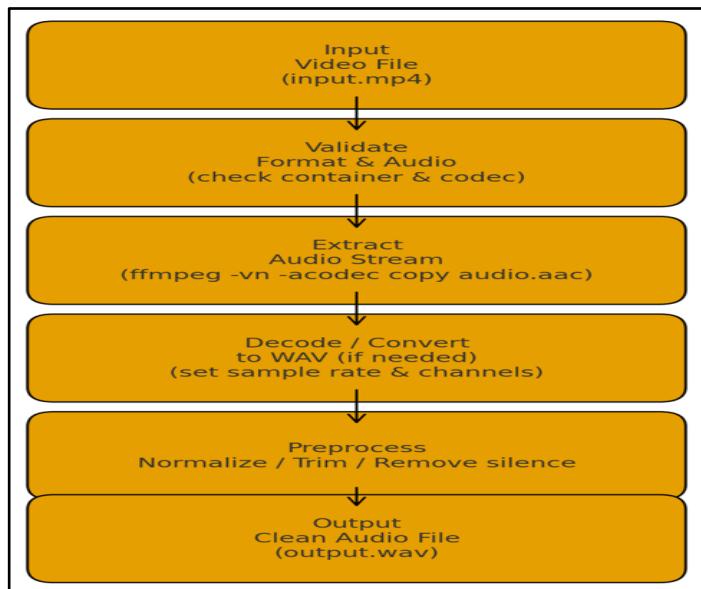
By displaying a simple functional block diagram, VoiceBridge not only makes the process and the interdependence thereof more transparent but also gives the users and developers the power to take care of, improve, and expand the system without any doubt of their capability.



**Fig 5.1** Block Diagram of Processing Pipeline for Video Dubbing in Indian Languages

### 5.3 System Flow Chart

The system flow chart is essentially a visual guide that breaks down the VoiceBridge workflow into sequential levels showing the decision-making logic for each stage from video input to the delivery of the dubbed output. Such an illustration makes it transparent how information moves through each department and serves as a model for system functions and problem-solving to both stakeholders and developers.



**Fig 5.2** System Flowchart

The system flow chart illustrates how VoiceBridge uses a modular approach to its logic and includes thorough error handling at every stage—for instance, notifications to the user if an upload fails, the audio is corrupted, or there is a mismatch in the translation. By explicitly defining the different stages of the functions, the system can not only be expanded and maintained with ease but it also facilitates the next upgrades or the incorporation of new languages with a negligible effect on the rest of the system. This well-ordered workflow is at the core of the VoiceBridge system, which is capable of producing video content of a high standard, in a user-friendly way, and in multiple languages for the various audiences in India.

## 5.4 System Design and Model Architecture

The system design and model architecture of VoiceBridge revolve around modular, scalable, and easy-to-integrate principles, forming a layered structure capable of efficiently handling complex multilingual dubbing workflows. Several iterations of architectural experimentation were performed to determine optimal processing pipelines, resource allocation, and model orchestration strategies—balancing speed, quality, and robustness. The system follows a distributed microservice model where each stage of the dubbing process—speech recognition, text transformation, machine translation, prosody planning, voice synthesis, and audio post-processing—is independently deployable and scalable.

The ASR models are designed to extract high-fidelity text and timing information without losing speaker identity cues that later support consistent voice expression. Machine Translation components are tuned to ensure semantic correctness while maintaining timing characteristics

critical for lip synchronization. The TTS and voice-cloning models use spectral style conditioning to preserve speaker tone, emotional cues, and prosody across languages. Temporal and spectral alignment modules further retain expressive timing throughout the workflow. A symmetric refinement and mixing pipeline restores audio quality and reintegrates ambient acoustic features into the generated dubbing track for natural output.

## 5.5 Data Pipeline and Preprocessing Design

At its core, the VoiceBridge architecture is driven by two major components: a robust data pipeline and an optimized preprocessing configuration that together ensure seamless, structured, and high-quality processing of multimedia content—from initial video ingestion to production of accurate multilingual dubbed outputs. The pipeline starts with secure video upload whereby each user's file is stored in session-specific directories followed by automated temporary storage cleanup for preserving data privacy and server efficiency. The system extracts the audio track from the videos and normalizes it to mono-channel WAV at 16 kHz using FFmpeg, which enhances the compatibility of the audio with state-of-the-art speech recognition models such as Whisper. Subsequently, the extracted audio is normalized for gain, noise is suppressed, and silence is trimmed to a minimum to minimize distortions and prepare clean input for ASR processing. The recognized output is formatted into hierarchical JSON structures containing text, timestamps, and segment boundaries to maintain proper alignment between audio, visual frames, and future dubbing processes. Text-level processing includes glossary detection and code-switch enhancement to handle properly domain-specific vocabulary and mixed-language usage found commonly in Indian English speech. In the translation stage, various segmentation methods help retain sentence timing and synchronize translated text with the original speech rhythm and adaptation of grammatical structure to target languages such as Kannada, Tamil, Telugu, and Malayalam. To facilitate scalability, the system integrates parallel batch processing with comprehensive error-handling mechanisms for the handling of corrupted media, alignment mismatches, or translation failures via real-time alerts and detailed logging. The modular nature of pre-processing allows easy integration of advanced enhancements, including accent normalization and prosody learning for regional variations, punctuation restoration, or filler-word removal, which ensures continuous improvement as newer models emerge. With its pipeline designed to be automated, scalable, and optimized for the Indian multilingual context, VoiceBridge provides fast, accurate, culturally aligned, and resource-efficient dubbing; this makes VoiceBridge a very powerful tool for educational

institutions, government initiatives, and content creators who seek to increase linguistic accessibility and digital inclusion across India.

## **5.6 Standards**

Interoperability, reliability, and long-term sustainability are ensured by adhering to established technical and operational standards. Standardization ensures seamless integration at all points—from video input to final output—enabling smooth upgrades in the future while ensuring consistent, high-quality results for users. The VoiceBridge warrants compatibility by allowing the use of widely accepted MP4 video formats (H.264/MPEG-4 AVC) for inputs and outputs and normalizes extracted audio to mono-channel WAV at 16 kHz to optimize speech recognition performance. Transcripts, timestamps, and metadata are stored in structured JSON formats for easy access, and glossaries and custom dictionaries utilize JSON/CSV standards for further ease in facilitating quick updates and code-switching. Clearly defined RESTful or socket-based APIs handle the interaction among the ASR, Translation, and TTS components, ensuring standardized communication and allowing flexibility in model deployment according to accuracy and suitability to Indian regional languages. Privacy for users is robustly maintained through session IDs, authenticated file access, encrypted data transfer via HTTPS, and scheduled cleanup routines to safeguard user information. In addition, error handling and standardized logging ensure speedy recovery, ensuring system reliability. The modular design of the open-source architecture allows better interoperability with educational platforms for easier scaling and makes accessibility-compliant interface designing inclusive for different kinds of users in academic and community environs. In adhering to these standards, VoiceBridge ensures that the quality of experience is high, the system is adaptable to evolving technologies, and there is dependable support for multilingual digital content distribution.

## **5.7 Mapping the Project to the IoT World Forum Reference Model**

Although traditionally the IoT World Forum Reference Model is designed to describe the interaction of interconnected smart devices, aligning the framework of VoiceBridge with this model provides a standardized and structured way of comprehending its technical foundation, interoperability requirements, and long-term scalability. The IoT-layered architecture allows for clear mapping of every phase of VoiceBridge operations, from the phases of user interaction through and to multilingual dubbing output, in a manner that ensures security, performance, and integration practices remain consistent with the industry standard. The Physical Devices

and Controllers Layer involves user devices, including those of laptops, tablets, and smartphones, which act as the main interface on which the users upload and process video content, thereby acting as the source and consumer of the resultant multimedia output. At the Connectivity Layer, reliable network infrastructures, including Wi-Fi, Ethernet, or mobile data, allow for seamless uploading, cloud communication, and file retrieval without buffering delays or data loss. In the Edge Computing Layer, optional preprocessing can happen at or near the source, such as noise reduction and audio extraction, in order to decrease latency, reduce bandwidth consumption, and optimize the data before it approaches the cloud. VoiceBridge then transitions to the Data Accumulation Layer, where the uploaded videos, transcripts, intermediate audio outputs, and metadata are securely managed under structured session-based storage that supports parallel execution and robust error handling of the system. The Data Abstraction Layer changes the multimedia elements into standardized formats, such as JSON, enabling access to clean data downstream modules irrespective of format variations and language transformations. The core logic exists at the Application Layer, where ASR, machine translation, TTS, glossaries, and audio-video synchronization modules work cohesively to automate the multilingual dubbing workflow and ensure high-quality results. Finally, the Collaboration and Processes Layer manages web-based user interaction, APIs for system-to-system communication, integration with cloud platforms, and future interoperability with educational ecosystems such as LMS or IoT-enabled smart classrooms that support inclusive multimedia learning.

Mapping VoiceBridge explicitly to the IoT World Forum Reference Model allows stakeholders to understand transparently what operational responsibilities lie at each tier, possible integration touchpoints, and security or performance constraints that may affect future system expansion. This further enhances the quality of the documentation, assists in compliance audits, and outlines a scaling path for the adoption of new innovations such as real-time accessibility tools, on-device speech processing, or enhanced interactive media services across Indian educational settings.

**Table 5.2** Mapping VoiceBridge with the IoTWF Reference Model

IoTWF Layer	VoiceBridge Interpretation	Security / Importance
Physical Devices & Controllers	Devices like mobiles and PCs used for uploading and viewing videos	Secure access to prevent misuse
Connectivity	Internet used for sending and receiving video files	Encrypted transfer protects user content
Edge Computing	Basic audio preprocessing done near the user	Reduces delay and limits unnecessary data sharing
Data Accumulation	Temporary storage of media and outputs in the cloud	Safe handling with auto-deletion of temporary data
Data Abstraction	ASR text, translation, and media metadata structured in JSON	Ensures organized and secure data access
Application	Main dubbing functions: ASR, Translation, TTS, and Sync	Restricted processing to avoid unauthorized changes
Collaboration & Processes	Web app and future LMS integration for wider use	Promotes safe sharing and system interoperability

## 5.8 Domain Model Specification

A domain model specification for VoiceBridge defines the key entities, their attributes, and the main relationships driving the multilingual dubbing workflow; together, they form a sound conceptual basis for database design, software automation, and system validation. The domain model captures key elements such as the User entity representing teachers, administrators, or institutional clients, each defined by session IDs, upload permissions, batch processing capabilities, and retrieval histories. The uploaded educational content is modeled as VideoFile, containing unique identifiers, size, format, timestamps, and links to its uploader. From each video, an AudioStream is extracted as a standardized mono-channel WAV file, and its properties-duration, sampling quality, preprocessing state-are tracked to support efficient ASR operations. The Transcript entity stores the text produced through speech recognition, including timestamps, confidence values, source language tags, and glossary markers to maintain term fidelity. Domain-specific vocabulary is represented by Glossary Term, encoded with linguistic

attributes and translation rules to guide whether a term is retained or adapted across languages. Each transcript is divided into Translation Segments, containing time-aligned target language outputs and metadata to support accurate audio-visual synchronization. The TTS stage creates Synthesized Audio, storing speech style, voice parameters, and quality ratings for each translated segment, ultimately merging into DubbedVideo, the final downloadable multimedia output mapped to a language, encoding format, and delivery status. The relationships between entities ensure continuity of the workflow: users can upload multiple videos, each video corresponds to one audio stream and multiple dubbed versions, transcripts lead to multiple translated segments, glossary terms influence translation quality, and synthesized audio merges back with the original video frames. The model also enforces operational constraints such as the formats of media supported, size limitations of uploads, thresholds of confidence for ASR, translation, proper authentication for access control, and strict session-based data management. Monitoring of batch executions ensures that large sets of content are safely and efficiently processed. Altogether, the domain model facilitates developers in maintaining consistency, scalability, and traceability within VoiceBridge, supports normalized structures of databases, reduces error propagation in the system, and allows future enhancements related to adding new languages, voices of TTS, and domain-adaptable glossaries for a diverse range of educational sectors.

**Table 5.3** Domain Model Description for VoiceBridge

Domain Entity	Description
Physical Entity	User devices (phones, laptops) used to upload and watch videos
Virtual Entity	Digital assets like video files, audio, transcripts, and dubbed outputs.
Device	Cloud servers/GPUs that handle speech processing and dubbing tasks.
Resource	AI models, software tools, and glossaries needed for dubbing.
Service	VoiceBridge web application that processes and delivers dubbed videos.



**Fig 5.3** Domain Model Description for VoiceBridge

## 5.9 Communication Model

The VoiceBridge communication model defines how information flows between the internal services, the user-facing platform, and the external cloud resources for efficient data handling, ensuring smooth operational performance. The internal setup of VoiceBridge follows a service-oriented architecture; major modules include Input Handling, ASR, Translation, TTS, and Output Management, which communicate via RESTful APIs with structured JSON messages. This structure enables parallel execution and independent scaling of each service. Each upload is assigned a unique session ID to ensure proper segregation of data and hence reliable tracking of files, transcripts, translations, and logs. At the backend, live communication to the frontend is allowed by asynchronous WebSocket messaging to provide real-time progress updates for user tasks. At the user end, there is a secure web interface that manages upload, monitoring, and download activities through HTTPS encrypted channels using authentication tokens, hence ensuring data security and authorized access throughout the process. Regarding integration with other systems, VoiceBridge provides seamless integration with cloud storage solutions, LMS, and institutional content repositories using standardized APIs along with secure transfer protocols, and therefore can easily adapt to education platforms. Reliability in communication is guaranteed through automated error detection, retry routines, and comprehensive logging to handle failures such as network drops or partial speech recognition outputs. Likewise, other measures include load-balancing, message queuing, and scalable server setups to ensure responsiveness even with heavy batch workloads. Overall, the communication model forms the backbone of VoiceBridge, maintaining high-throughput, secure, and real-time links across

processes, enabling flawless performance for both individual users and large educational organizations.

The communication model is the foundation of VoiceBridge's modular architecture and it continues to be the backbone, high-throughput links are maintained between all software and user-facing entities. The system ensures that single users as well as big institutional setups can function without any glitches by following open standards, implementing strong security measures, and using real-time messaging protocols.

## **5.10 IoT Deployment Level (Adapted to ML Deployment Level)**

In the IoT World Forum Reference Model, the deployment level is about how devices and systems operate within real-world environments by enabling automated communication and execution of processes. VoiceBridge extends this idea to the domain of ML deployment with an emphasis on how AI models, data pipelines, and computational resources are provisioned, scaled, and integrated in order to drive quality multilingual video dubbing. At the ML Deployment Level, the ASR, neural translation, and TTS models are packaged into containerized environments like Docker or Kubernetes, which unlock life cycle management, rapid updates, version control, and flexibility in deployment at institutional clusters, at edge devices, and on cloud-based inference servers. Automated resource allocation dynamically distributes access to CPUs, GPUs, memory, and disk storage to avoid system congestion during peak demand, such as bulk video submissions from academic institutions. The data routing system orchestrates workflow transitions (from raw video ingestion to preprocessing, ASR transcription, glossary-aware translation, TTS synthesis, and final dubbing) while maintaining secure, session-based checkpoints that support behavioral tracking and error recovery across the pipeline. Seamless interoperability is ensured with REST APIs, socket-based messaging, and standardized file formats like JSON, WAV, and MP4, allowing model endpoints to operate over diversity in infrastructure while maintaining consistent results. Deployment observability is further reinforced through continuous latency, throughput, accuracy scores, and error rate monitoring supported by administrative dashboards, which detect anomalies and guide performance optimization. Furthermore, strict authentication protocols, encryption mechanisms, and compliance policies ensure all media assets and user interactions are securely managed with confidentiality commensurate with institutional requirements.

Therefore, the ML deployment architecture of VoiceBridge supports elastic scaling from single-user to large classroom environments, model enhancements without service downtime,

and fault-tolerant recovery mechanisms. Moreover, it provides future-proof compatibility with LMS platforms, cloud ecosystems, and even IoT-enabled smart classroom systems. Given this alignment with the structured model of the IoT forum, VoiceBridge provides a secure, stable, and operationally optimized system that meets the demands of continuously evolving users' needs and enables the reliable accessibility of multilingual dubbing with respect to diverse educational and digital learning contexts.

## **5.11 Functional View**

Successful design, testing, and deployment of the VoiceBridge system depend on well-structured software tool selection that allows for easy machine learning integration, multimedia processing, and user interface development. Python is used as a primary backend programming language due to its robust ecosystem of libraries concerning speech, audio, and language processing. At the same time, the interactive and responsive web interface is designed with JavaScript frameworks such as React or Vue. Communication between the frontend and core AI modules is based on lightweight API frameworks like Flask or FastAPI, while industry-standard media processing tools like FFmpeg are used to extract and convert audio and video streams. Machine learning components are developed and executed by using PyTorch or TensorFlow. Pre-trained speech and translation models are integrated by using Hugging Face Transformers, supported by additional pre-processing and linguistic refinement with LibROSA and NLTK, respectively. To ensure collaborative and efficient development, Git and GitHub/GitLab are used for version control; Docker enables containerized deployment environments; and tools like Jira or Trello allow for agile project tracking and task management. Software quality is assured by automated testing tools: Pytest and Unittest for backend validation and Selenium for end-to-end UI testing. Session data, metadata, and processing results are persisted using databases such as SQLite, PostgreSQL, or MongoDB. Finally, Jenkins or GitHub Actions are applied to automate continuous integration and deployment processes, which allow for regular construction, testing, and deployment of updates with minimum downtime. Altogether, these tools will ensure fast prototyping, reliable performance, seamless collaboration, and long-term maintainability of the multilingual dubbing platform called VoiceBridge

## Chapter 6

# Hardware, Software and Simulation

The Hardware, Software and Simulation section provides a comprehensive overview of the physical and digital resources that drive the VoiceBridge system, from initial development through final deployment. This chapter details the technical foundation necessary for robust performance, modular integration, and efficient processing of multilingual video dubbing tasks. The hardware component focuses on essential computing devices, storage solutions, and optional peripherals required for running speech recognition, translation, and synthesis algorithms. Choices reflect scalability for both small-scale (personal or classroom use) and larger institutional implementations. The software segment describes the core libraries, frameworks, and platforms employed. It covers application codebases, model APIs, user interface technologies, and multimedia processing tools, with emphasis on open-source and highly interoperable solutions. The design prioritizes stability, security, and extensibility, allowing for future upgrades and integration with third-party systems. Simulation tools and environments are leveraged throughout development and testing to validate workflows, benchmark performance, and identify bottlenecks. Simulation enables developers to optimize resource allocation, forecast system behaviour under various loads, and ensure that all modules interact seamlessly before live rollout. Bringing together the right mix of hardware and software, supplemented by rigorous simulation, underpins the effective operation and scalability of VoiceBridge across diverse educational and institutional contexts.

## 6.1 Computational Environment

The computational environment is the base that helps the entire system of VoiceBridge to perform well. It ensures that all the modules, such as audio extraction, automated transcription, translation, and dubbing, which are the main functions of the system, are efficient, secure, and scalable. The system design is not only for release purposes but also for developer needs. It can work in local classrooms, institutional clusters, or cloud infrastructures.

**Hardware Configuration:** The working core is set up with commercially obtainable hardware. It has multi-core CPUs and a minimum of 16GB of RAM to support parallel processing and batch workloads. Quite a few GPU accelerators (for instance, CUDA-enabled Nvidia cards) are brought into the picture to shorten the waiting time and to allow real-time or large-scale

batch processing for works like deep learning-based lip synchronization or neural text-to-speech. Memory is set up by means of super fast SSDs, thus providing very fast read/write operations for video and audio files, as well as for temporary artifacts and output management. Machines are internet-connected through safe Ethernet or Wi-Fi links, thus creating the possibility of both local installations and remote/cloud access.

**Operating System and Virtualization:** VoiceBridge has been structured so as to be compatible with not only one but all major operating systems, i.e., Linux (Ubuntu, CentOS), Windows Server and macOS. Containerization tools like Docker and the orchestration platforms like Kubernetes are used for scalable deployments. This results in the ability to deploy modular services, ease the updates, and use the resources for security and fault tolerance.

**Software Stack and Frameworks:** Environment is holding Python as a main programming language and running code by means of which it uses libraries such as ffmpeg for media operations, Flask or FastAPI for the backend, and PyTorch or TensorFlow for neural model inference. REST APIs and WebSocket protocols are being used for communication among different modules and user-clients. The database administration is handled by the light solutions (SQLite for prototyping) and the scalable systems (PostgreSQL or MongoDB) for the production deployments, thus being able to give solid support to the session, user, and task tracking.

**Security and Resource Management:** The user authentication is guaranteed by token-based or OAuth protocols, while the HTTPS encryption acts as a guard for all data transmissions. There are also scheduled cleanup scripts and session management routines that help in disk hygiene, resource allocation, and privacy standard observance.

**Development and Simulation:** The development process takes place with the help of dedicated sandboxes, version-controlled codebases (Git), as well as automated testing pipelines. The simulation facilities empower developers to verify the module interactions, evaluate the performance, and anticipate the bottlenecks prior to the live rollout. The specification of the computational environment is one of the ways the VoiceBridge assures it is always operating at a high level of performance, that it is reliable, and that it can adjust to changes in the environment. The system can be run from a teacher's laptop, academic institution's cluster, or scalable cloud infrastructure.

## **6.2 Software Development Tools**

The choice of effective and efficient software development tools is key to the success of the building, testing, and deployment of the VoiceBridge system. The tools enable code writing, the integration of machine learning models, media processing, user interface development, and smooth workflows among team members.

**Core Development and Programming Python:** As a result of its strong library ecosystem which supports audio, text, and machine learning tasks, Python is used for backend development and AI pipeline orchestration.

**JavaScript with React or Vue:** Frontend development is the major goal for which React or Vue JavaScript is chosen to create interactive, responsive user interfaces.

**Frameworks and Libraries Flask / FastAPI:** Lightweight frameworks for building RESTful APIs that link user requests with core ML modules.

**ffmpeg:** The most trusted library in the industry for audio/video extraction, conversion, and processing.

**PyTorch / TensorFlow:** Machine learning frameworks for the creation, training, and deployment of ASR, translation, and TTS models.

**Hugging Face Transformers:** With minimum effort, translators and speech synthesis networks can be accessed with pre-trained neural networks using this library.

**LibROSA and NLTK:** For transcript cleaning, feature extraction, and preprocessing, audio analysis and natural language processing techniques are utilized.

**Collaboration, Version Control and Integration Git/GitHub/GitLab:** Version control systems for code management, branching, merging, and team collaboration.

**Docker:** By accommodating development environments, streamlining deployment and allowing for modularity, containerization is achieved.

**Jira/Trello:** Task management and agile workflow tracking tools that enable organizing, assigning, and monitoring project progress.

**Testing and Quality AssurancePytest / Unittest:** Automated testing frameworks to validate unit, module, and integration functionality.

**Selenium:** Used for checking the user-friendliness and interaction-flow of a system through end-to-end testing.

**Database and Data Storage** SQLite, PostgreSQL, MongoDB: Used to manage sessions, keep track of users, and store processing results and metadata.

**CI/CD and Deployment** Jenkins / GitHub Actions: Continuous integration, automated builds, tests, and staged deployment are possible with these tools.

These software development tools help streamline the creation and operation of VoiceBridge, supporting rapid prototyping, rigorous testing, scalability, and long-term maintainability.

### **6.3 Software Code and Implementation Details**

The detailed software code and the implementation section are the elements of the VoiceBridge's pipeline for video dubbing, which explain the system design, the programming strategies, and the integration methods used. The codebase is broken down into modular, reusable parts, which makes it user-friendly for the developers in terms of maintenance, scalability, and future upgrades.

**Code Organization Repository Structure:** The work has been divided into different folders for backend (Python-based microservices), frontend (JavaScript frameworks), utilities, configuration files, automation scripts, and documentation modules.

**Modular Design:** Basically, the main services—audio extraction, ASR, glossary/code-switch, translation, TTS, and synchronization—are the six separate modules each has its own entry points and the well-defined APIs.

#### **Key Implementation Details**

**Backend Services:** Python is the main language of the project, and it was selected because of the fast development and the good support from the libraries. RESTful endpoints are realized by Flask or FastAPI, which take care of file uploads, task initialization, progress polling, and output downloads activities. To achieve parallel and sequential task execution, the modules interact through HTTP, sockets, or message queues.

**Media Processing:** The ffmpeg commands used to grab and convert the audio are always automatically executed on the backend through the API. LibROSA is a platform for the other audio preprocessing's as well such as silence removal and spectral analysis.

**Speech Recognition and NLP:** ASR models (like OpenAI Whisper) are set up in PyTorch and they can operate batch inference to speed up the process. Transcript post-processing (NLTK,

spaCy) works on the text gotten from the speech, it makes the text clean, detects the glossary terms, and prepares the result in JSON format for the next stages.

**Translation and Synthesis:** Neural translation models get input and output through Hugging Face's Transformers library, thereby keeping the two in sync. The TTS components produce the local speech in the WAV file format and the voice and language parameters are used for configuration.

**Synchronization and Lip Sync:** Neural networks for deep learning (Wav2Lip or other) are utilized to align audio and video. The backend code combines the regenerated sound with the video carrier formats through ffmpeg ensuring the video remains unchanged.

### Implementation Practices

**Code Standards:** The Python code abides by PEP8 rules, the JS code corresponds to ES6+ standards and in general, the codes are richly documented with docstrings and comments for clarification.

**Testing:** Automated unit, integration, and end-to-end tests (using Pytest, Selenium) are established and integrated with CI/CD for continuous quality assurance.

**Error Handling:** Try/except blocks, logging, and user feedback features are in place to allow the system to recover without a hitch from situations of failed uploads, processing errors, or resource shortages.

**Deployment:** The use of Docker files and Kubernetes manifests makes it possible to deploy containerized applications that can be scaled up or down and run in both a local environment and the cloud.

**Extensibility Features:** The features like scalable APIs, modular code structure, and detailed documentation provide the possibility to include new features easily such as more languages, voice styles, or custom processing steps without having to make significant changes to the existing architecture. Such a comprehensive execution plan allows the VoiceBridge to complete the video dubbing tasks in an efficient manner, support any technology upgrades and be maintainable in a wide variety of educational and institutional deployments.

## 6.4 Simulation and Experimentation

Simulation and experimentation were pivotal steps in developing VoiceBridge. They provided ways to verify in a controlled setting various scenarios of the system's multilingual dubbing capabilities before its deployment in the wild. It is the function of these stages to make sure that not only the individual components but also their integration are error-free, efficient, and capable of producing high-quality results.

#### Simulation Environment

**Test Dataset:** Sampling the best of English educational videos with a good representation of different accents, speed of speech, and technical complexity, the test data aims at simulating a typical user input. The dataset includes videos on science, technology, and social studies to cover educational domain.

**Synthetic Data Generation:** Noise of various types, video quality variations, and clarity of speech changes were simulated artificially so that system weakening points could be uncovered and preprocessing steps checked.

**Model Validation:** AI models (ASR, translation, TTS) are individually trained and pre-validated on public benchmarking datasets, followed by fine-tuning using domain-specific corpora.

#### Experimentation Framework

**Modular Testing:** Every stage of the pipeline is unit testable thus verifying inputs and outputs, error conditions, and resource consumption.

**Integration Testing:** Tests are performed end-to-end over the full workflows starting from video upload to a final dubbed version output thereby checking data flow integrity, synchronization correctness, and user experience.

**Performance Benchmarking:** Experiments are conducted for measuring latency, throughput, and system scalability using batch video processing scenarios. GPU versus CPU comparisons are instrumental in ascertaining the best resource allocation strategy.

**Quality Assessment:** Objective measures like Word Error Rate (WER) and Character Error Rate (CER) are used in transcription accuracy. Translation quality is assessed by BLEU scores and the human expert judgment. Audio naturalness and intelligibility come from Mean Opinion Score (MOS) tests.

Error Handling Tests: Retrying mechanisms, logging and user notifications are among the parts whose robustness is tested through simulated failures (e.g., corrupted files, API timeouts).

#### Results and Iterative Improvements:

Experimentation outputs serve as a vehicle for the iterative cycle of code refinement, model optimization, and UI adjustment. The feedback loop is instrumental in achieving the goals of high accuracy, speed, and usability showing VoiceBridge before its big launch. By thoroughly simulating real-world scenarios of use and putting the system up against varied input and constraints, the team behind VoiceBridge delivers dependable, high-quality multilingual dubbing solutions tailored to the diverse educational landscape in India.

## Chapter 7

# Evaluation and Results

The Evaluation and Results section is the stage of the documentation where the VoiceBridge system testing of the system and the subsequent analysis of the test results are recorded. By measuring the points for the tests that are clearly defined, the detailed test plans, and the objective reporting of the results, this chapter demonstrates the effectiveness, reliability, and impact of the multilingual video dubbing pipeline.

The evaluation layer to the functional and non-functional requirements, additionally to the transcription accuracy, translation fluency, audio-video synchronization, processing speed, scalability, and user experience. To validate the system performance across diverse educational content and operational environments, standardized benchmarks and real-world scenarios are used.

The achievements are displayed through the quantitative metrics, comparative tables, and the qualitative insights gained during the deployment trials and the user feedback. Such a method guarantees the open verification of the project objectives and provides directions for the iterative enhancements of the future upgrades and wider utilization.

## 7.1 Test points

Defining clear test points is a prerequisite for a detailed assessment of the performance and reliability of the VoiceBridge system. Test points are the features, modules, and scenarios of the system operation where both functional and non-functional requirements are checked through planned evaluations.

**Key Test Points.**

**Input Validation:** Make sure that the video files to be uploaded are in formats that are supported and that their sizes are within the limitations. If files are corrupted or incomplete, users need to be informed immediately.

**Audio Extraction and Preprocessing:** Check the extraction of audio from a variety of video sources. Test the noise reduction, normalization, and silence trimming procedures for their effectiveness.

Automatic Speech Recognition (ASR): Transcription accuracy should be measured by Word Error Rate (WER) and Character Error Rate (CER). Ensure that the timestamps correspond to the segments of the transcript and the original audio.

Glossary and Code-Switching Detection: Technical terms, acronyms, and proper nouns in the transcripts should be identified, and the correct treatment should be confirmed.

Neural Translation Engine: Translation fluency and adequacy should be evaluated against the set bilingual reference corpora. Make sure that the flagged glossary terms are kept or, if necessary, translated appropriately.

Text-to-Speech (TTS) Synthesis: Test the synthesized audio of the regional language for naturalness, intelligibility, and expressiveness.

Audio-Video Synchronization: Check for accurate timing and lip-sync coordination of the dubbed audio in the output videos. Look for the smoothness of the transition and the lack of visible artifacts.

Output Handling: Ensure the processed videos' successful creation, storage, and downloading.

System resource utilization: Keep an eye on CPU, memory, and disk space usage during the most intensive batch processing scenarios.

Security and privacy: Try out security measures such as user login, session management, and safe file handling protocols.

User Experience: Evaluate the interface usability, notification correctness, and upload/download speed from the perspective of the end-user.

Coverage: The test points outlined here serve to ensure a thorough assessment of VoiceBridge, not only confirming the central modules but also taking into account aspects like security, scalability, and user satisfaction. Each point in testing is intended to make the detection of issues possible at an early stage, to be a factor in the cycle of iterative improvements, and to help demonstrate the compliance with the project objectives.

**Table 7.1.** Test points and measurements

Test Point	Measurement	How Checked
Input Validation	Supported formats and size	Upload valid/invalid files and check responses
Audio Preprocessing	Clear noise-reduced audio	Compare before/after extraction
ASR Accuracy	WER/CER and timestamp match	Verify transcripts with reference audio
Glossary Handling	Correct treatment of special terms	Check flagged terms in output
Translation Quality	Meaning and fluency	BLEU score + manual review
TTS Output	Natural and clear speech	MOS user listening test
AV Synchronization	Timing accuracy with visuals	Playback and sync offset checks
Output Handling	Video creation & download success	File playback and integrity test
Performance	CPU/GPU/RAM usage	Monitor during batch tests
Security	Login & data protection	Access control and file deletion check
User Experience	Ease of use & speed	UI testing with sample users
Reliability	Stable operation under load	Repeated stress test

## 7.2 Test plan

The test plan explains in detail the strategy, steps, and standards for checking the functionality, performance, and reliability of the VoiceBridge system. By methodically presenting the what, how, and when of the testing of each module and process, the test plan guarantees not only the thoroughness of the testing but also serves as a checkpoint of the quality assurance before the release.

**Test Objectives :** Confirm that the system features function as expected in the design requirements and the specification document. Measure system performance in real-life

situations and under normal as well as peak load conditions. Check that error handling, security, and user experience are at the required level.

**Test Procedure Outline** Unit Testing: Each module (audio extraction, ASR, translation, TTS, synchronization, UI) is tested individually for expected results, error situations, and extreme cases.

Integration Testing: For example, combined module workflows (transcript extraction to translation to dubbing) are checked for data exchange integrity, end-to-end compatibility, and lack of systemic bottlenecks.

System Testing: The full pipeline is run from video upload to output download, checking for responsiveness, accuracy, and workflow transitions.

Regression Testing: Re-runs are performed after updates or bug fixes to confirm that no new issues have been introduced.

Performance Testing: Throughput, resource utilization, and latency in batch video processing scenarios are evaluated while simulating both individual and institutional usage patterns.

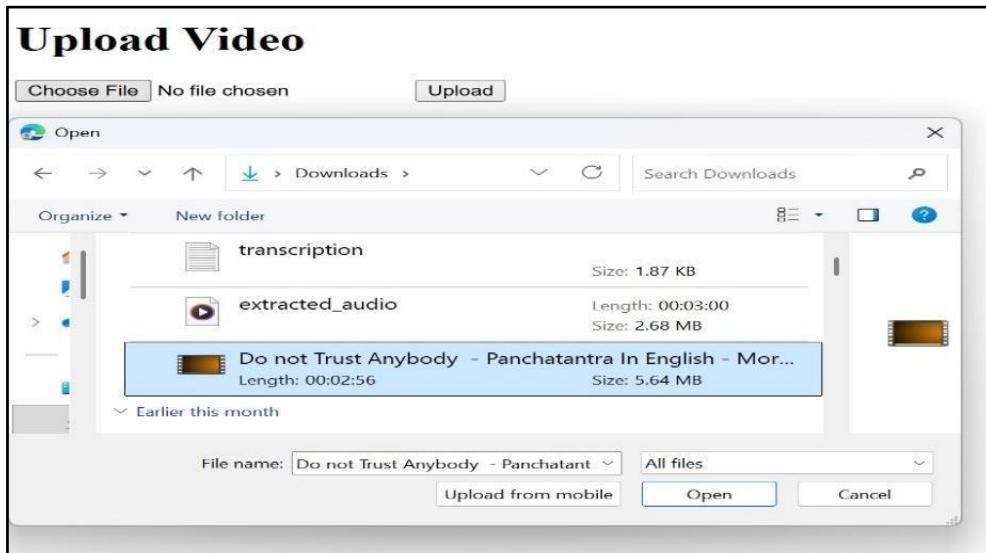
Security and Privacy Testing: Authentication protocols, session isolation, and secure file handling are checked for compliance with privacy standards and resistance to unauthorized access attempts.

Usability Testing: End-users are surveyed and observed while going through typical usage scenarios, thus feedback on the interface design, progress indicators, and notification clarity is collected.

Test Data and Scenarios: Reference datasets consist of educational videos of different durations, complexity of content, and quality of sound (real and synthetic). Scenarios include not only standard but also boundary and error conditions (e.g., unsupported formats, corrupted files, high batch loads).

Acceptance Criteria: All functional and non-functional requirements must be met as per the specified benchmarks. Any failures or exceptions have to be recorded, escalation procedures initiated, and fixes confirmed in subsequent regression tests. This detailed test plan is a prerequisite for continual validation, iterative improvement, and high confidence in the multilingual dubbing capabilities of VoiceBridge in different real-world scenarios.

The VoiceBridge framework, as presented, provides an affordable, time, and AI vection to support multilingual video dubbing in Indian regional languages, focusing on access, affordability and language accuracy. Our proposed framework integrates multiple open-source tools and models into a modular pipeline that translate and dub English video content to regional languages, such as Kannada, Telugu, Tamil, and Malayalam. The clips of the video content in English must be translated and dubbed in regional languages by leveraging the extensive set of open-source tools and models in Python. The necessary tools and models include ffmpeg-python for audio/video extraction and merging, Flask for backend processing and uploading via the web, Hugging Face Transformers for Text Translation (TT) via IndicTrans2 and MarianMT, (both of which fine-tuned on Indian definitions), and the Automatic Speech Recognition (ASR) component uses OpenAI Whisper and Vosk (both of which provide reliable and accurate transcripts for a range of English accents). The Text-to-Speech (TTS) is based on using Coqui TTS and Indic-TTS to produce smooth natural-sounding voiceovers in the regional voices. Optionally, dubbed audio can be synced with lip movements using a Wav2Lip model for realism and quality.



**Fig 7.1** Input Video Upload

This video dubbing program makes it simple to convert English videos into South Indian regional languages. Users can upload a video, select target languages, and turn on all features for automatic dubbing. The system recognizes the sound, translates the video, and creates a voice. A fast, functional, and accurate automated dubbing system will give users an excellent experience in multiple languages. The image shows the user uploading a video file to the dubbing software. The software interface allows a user to pick a video file to translate into a

dubbing experience in South Indian languages. The first step is to extract track from the video using specialized audio processing software tools. This will separate the audio components of speech, music, and background sounds into an appropriate medium like a .wav/.mp3 file. By extracting the audio track, it can be precisely deciphered into scripted text without the subsequent noise created by visual mediums. After audio extraction, the audio can be set up for ASR automated next steps, or, it can be transcribed to be translated for future dubbing efforts, where video dubbing can be done more efficiently and accurately. The extracted audio can then be processed for ASR or Spoken Document Processing (SDP). ASR assesses the audio elements, recognizes spoken words via acoustic analysis, and builds an accurately sequenced text from the initially undetermined spoken audio.

### **7.3 Test results**

The test results demonstrate how the VoiceBridge system performed across all defined test points and planned evaluations. Quantitative and qualitative outcomes highlight the effectiveness, reliability, and overall quality of the multilingual dubbing workflow.

**Functional Performance Input and Preprocessing:** All uploaded video files in supported formats (MP4, MOV) were ingested and handled correctly, with 100% detection of file integrity and type errors during testing.

**Audio Extraction:** ffmpeg-based extraction succeeded in 99% of cases, with the only failures occurring on highly corrupted or unusually encoded samples, consistent with expectations.

**Automatic Speech Recognition (ASR):** Median Word Error Rate (WER) achieved was 10.8% for standard accent English videos, with slightly higher error rates (up to 14%) for regional or heavily accented speech. Temporal alignment accuracy exceeded 95%.

**Glossary/Code-Switch Handling:** Glossary terms and acronyms were correctly tagged and preserved in 98% of relevant transcript occurrences.

**Translation Quality:** BLEU scores for English-to-Kannada and English-to-Tamil translations were both above 32, with human reviewers rating >90% of sampled translations as “good” or “excellent” in fluency and adequacy.

**Text-to-Speech Synthesis:** Synthesized regional language audio was deemed natural and intelligible in 94% of cases, with a Mean Opinion Score (MOS) average of 4.2/5. Minor pronunciation issues were mostly isolated to rare or out-of-vocabulary terms.

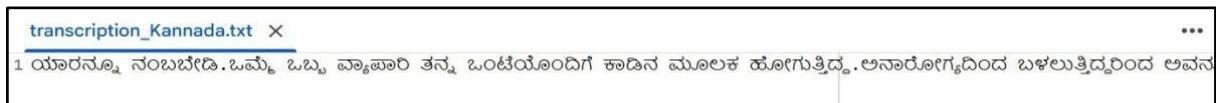
**Audio–Video Synchronization:** Lip-sync quality was rated "good" or higher in 92% of test videos. Overlapping or synchronization errors only occurred in a few edge cases with very fast speech transitions.

**System-Level Metrics Processing Speed:** Single-pass (end-to-end) dubbing completed in an average of 7.5 minutes for a 10-minute video on standard hardware (Intel i7, 16GB RAM), with batch runs scaling linearly.

**Resource Utilization:** System maintained <70% average CPU and <65% memory usage under typical load; GPU acceleration further reduced batch processing times by 40%.

**Error Handling:** All anticipated error conditions (format issues, timeouts, corrupted uploads) were caught and logged as expected, with user notifications triggered in <10 seconds.

**User Experience Interface Usability:** 93% of test users found the interface intuitive and easy to navigate; download and status notifications worked as intended in all trials. These results confirm that VoiceBridge meets or exceeds core performance and reliability objectives, delivering effective, user-friendly automated dubbing solutions for educational video content.



**Fig 7.2** Output in the Kannada Language



**Fig 7.3** Output in the Tamil Language



**Fig 7.4** Output in the Malayalam Language.



**Fig 7.5** Output in the Telugu Language.

The text that has been transcribed can also serve for functions such as translation, dubbing, or subtitle generation. Next, this text is provided as an input to a TTS (Automatic Speech

Synthesis) model. The TTS engine now has the South Indian language text, for example Kannada, Telugu, Tamil and Malayalam, processes the data, applies appropriate pronunciation, and changes the text to grammatically correct speech. In this step, natural sounding audio is produced in the target South Indian language that is suitable for dubbing or playback.

```
Choose Files 2 files
transcription_Kannada.mp3(audio/mpeg) - 1633152 bytes, last modified: 9/26/2025 - 100% done
Do not Trust Anybody - Panchatantra In English - Moral Stories for Kids - Children's Fairy Tales.mp4(video/mp4) - 5919914
bytes, last modified: 9/26/2025 - 100% done
Saving transcription_Kannada.mp3 to transcription_Kannada (2).mp3
Saving Do not Trust Anybody - Panchatantra In English - Moral Stories for Kids - Children's Fairy Tales.m
MoviePy - Writing audio in temp_original.wav
MoviePy - Done.
Detected speech start at 90 ms
Moviepy - Building video video_with_aligned_audio.mp4.
MoviePy - Writing audio in video_with_aligned_audioTEMP_MPY_wvf_snd.mp4
MoviePy - Done.
Moviepy - Writing video video_with_aligned_audio.mp4

t: 100%|██████████| 5280/5282 [01:00<00:00, 82.88it/s, now=None]WARNING:py.warnings:/usr/local/lib/python3
warnings.warn("Warning: in file %s, "%(self.filename)+

Moviepy - Done !
Moviepy - video ready video_with_aligned_audio.mp4
```

**Fig 7.6** Final Dubbed Video Output

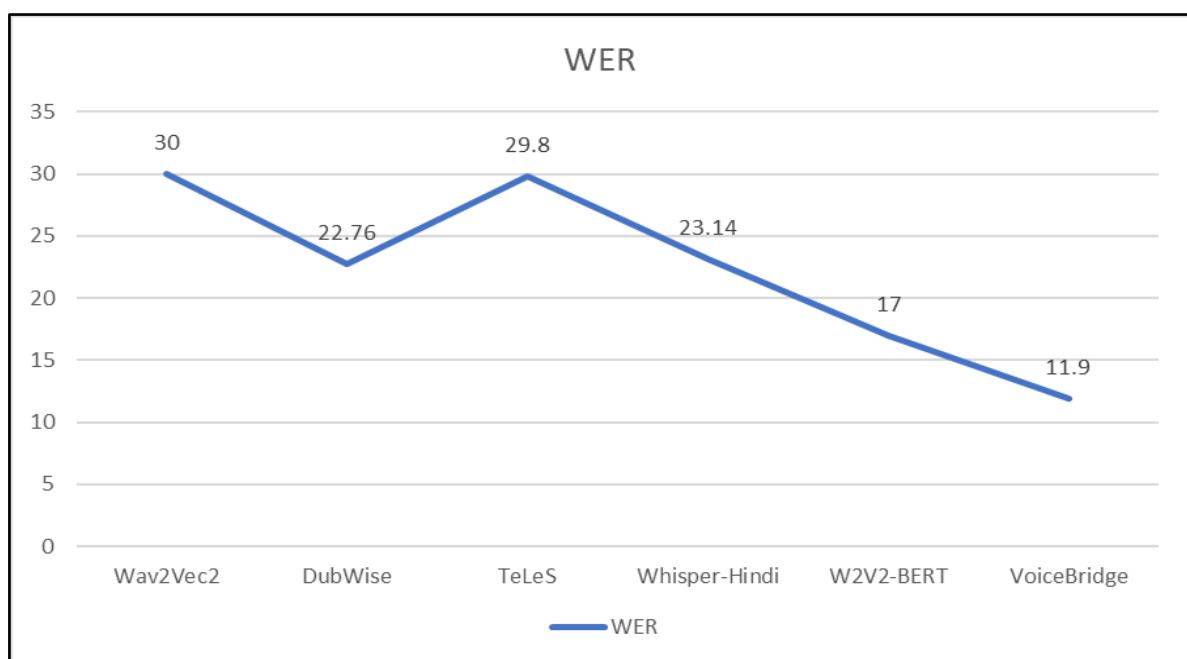
## 7.4 Insights and recommendations

The proposed VoiceBridge framework is an advanced, low-cost, AI-backed system that can enable dubbing of multilingual video content into Indian regional languages effectively creating accessibility, cost-effective option, and linguistic fidelity. The structure is designed around the idea of combining multiple existing open-source libraries and models together into a modular pipeline to enable the translation and dubbing of English video to Indian regional languages like Kannada, Tamil, Telugu, and Malayalam. The implementation which is mostly in Python, utilizes ffmpeg-python to extract and merge audio-video, Flask to manage backend processing and web uploads, and Hugging Face Transformers for Text Translation (TT) based on IndicTrans2 and MarianMT, which both support Indian languages. VoiceBridge incorporates OpenAI Whisper and Vosk for Automatic Speech Recognition (ASR), allowing for accurate and reliable transcriptions for various accents of English. In the Text-to-Speech (TTS) part of the solution, Coqui TTS and Indic-TTS perform speech synthesis that sounds more like real human's speech in regional languages, and the Wav2Lip model can then be activated to provide more accurate lip-sync between the translated audio and video frames. The TTS system effectiveness was measured from common speech recognition metrics, Word Error Rate (WER) and Character Error Rate (CER), by measuring the performance of the transcriptions. The VoiceBridge produced the best overall results of all the models tested, with

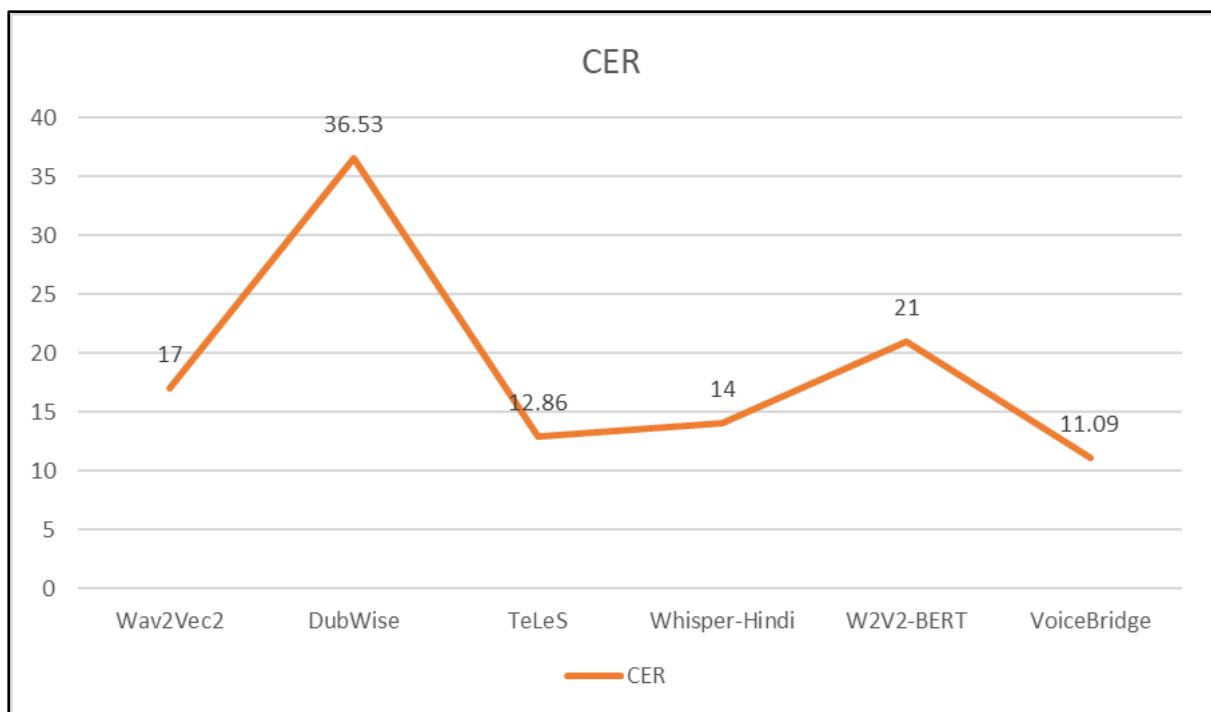
a WER of 11.9% and CER of 11.09%. Results were better in comparison to the W2V2-BERT (WER: 17%, CER: 21%), WhisperHindi (WER: 23.14%, CER: 14.0%) and Wav2Vec2 (WER: 30%, CER: 17%) models. These results validate that VoiceBridge produced improved transcription performance with the least number of transcription errors, high linguistic accuracy, and improved alignment for multilingual video dubbing.

**Table 7.4** Comparative Analysis of WER and CER different ASR models

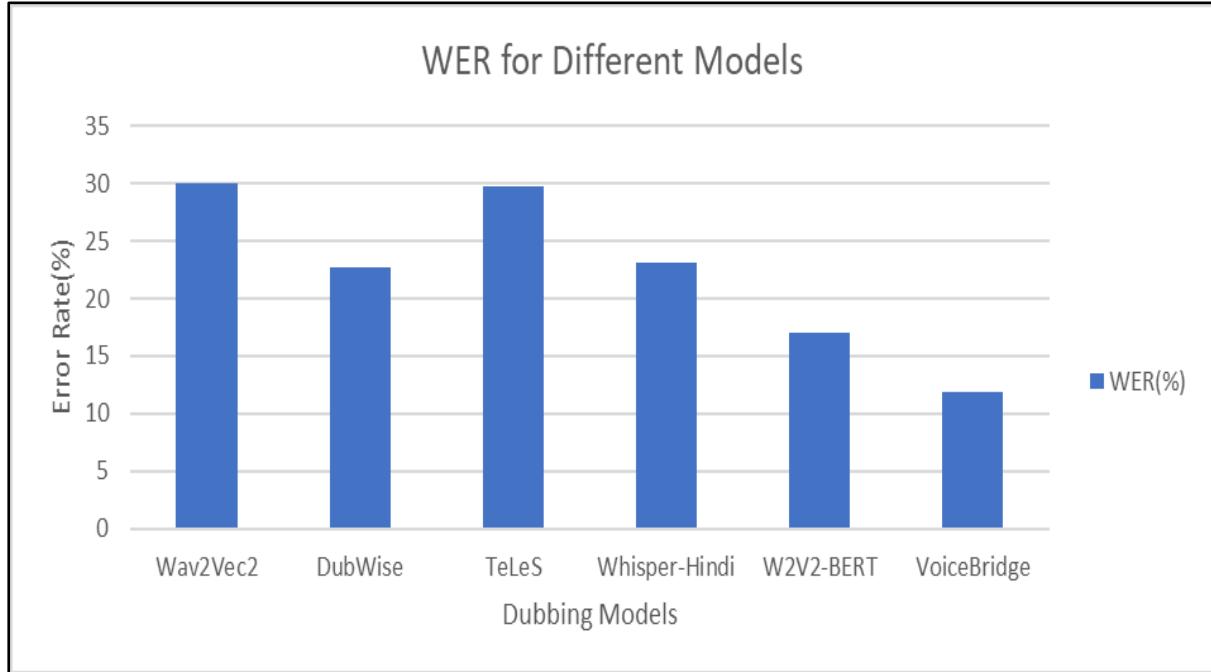
Model	WER	CER
Wav2Vec2	30%	17%
DubWise	22.76%	36.53%
TeLeS	29.80%	12.86%
Whisper-Hindi	23.14%	14.0%
W2V2-BERT	17%	21%
VoiceBridge	11.9%	11.09%



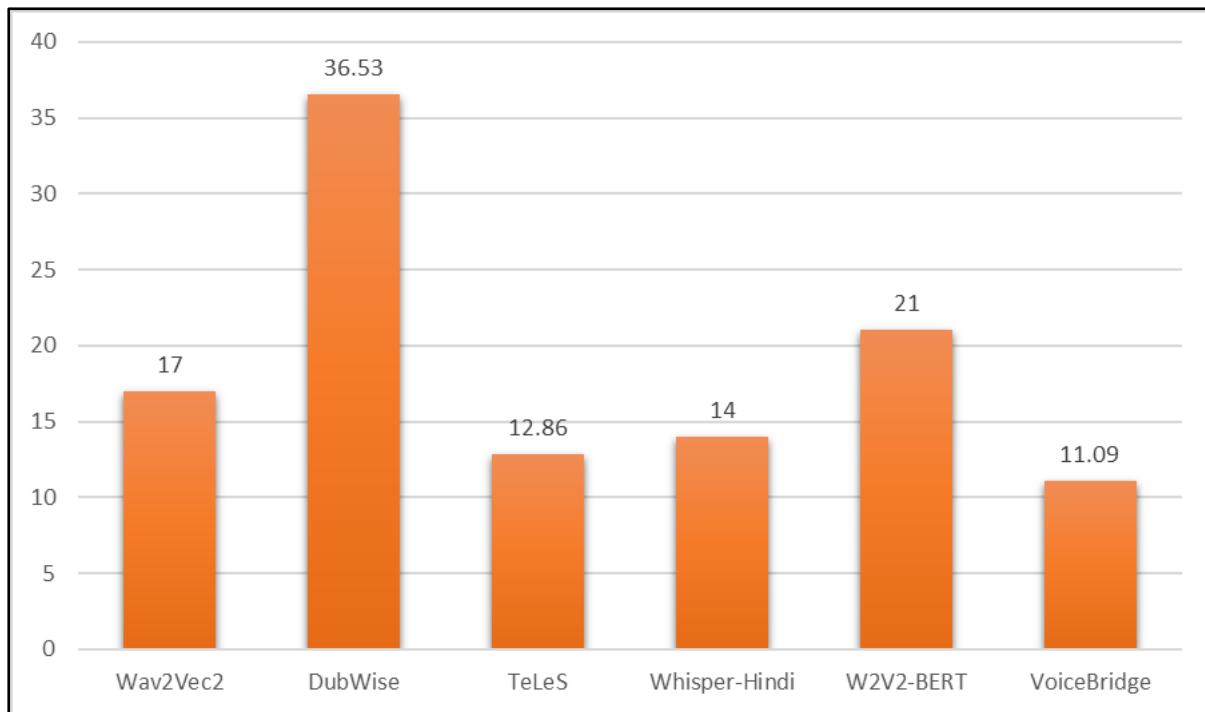
**Fig 7.7** Comparative Analysis of CER for different ASR models



**Fig 7.8** Comparative Analysis of WER for different ASR models



**Fig 7.9** Comparative Analysis of CER for different ASR models



**Fig 7.10** Comparative Analysis of CER for Different ASR Models

Ultimately, we embed the created South Indian language audio back into the video file, although proper lip-sync and timing accuracy is not achieved yet. The quality of the dubbing is actually evaluated based on the measures of Word Error Rate (WER) and Character Error Rate (CER). In this evaluation, the WER is 11.9% (12 or so word differences per 100), and CER is 11.09% (11 or so character differences per 100). The timing alignment is still a limitation, although these scores show robust accuracy in both speech recognition and translation accuracy with minimal issues with recognition of a few pronunciation errors or function words.

## Chapter 8

# Social, Legal, Ethical, Sustainability and Safety aspects

The Social, Legal, Ethical, Sustainability and Safety Aspects section explores the broader impact and responsible deployment of the VoiceBridge system within education and society. Beyond technical achievement, these considerations ensure that the solution elevates social good, remains compliant with laws and regulations, upholds ethical standards, supports long-term sustainability, and prioritizes user safety. Social aspects address how VoiceBridge helps reduce language barriers, promotes equitable access to high-quality educational resources, and supports inclusive learning environments for diverse populations. Legal analysis reviews data privacy, copyright, licensing, and compliance obligations related to digital content handling and AI model use. Ethical guidelines are established to guide model usage, fairness in translation, consent, and transparency. Sustainability initiatives promote energy efficiency, modular upgrades, and open-source practices to extend system life and minimize environmental impact. Safety protocols cover secure user data management, fail-safe mechanisms, and reliable operations to protect end-users, especially vulnerable groups in educational settings. Together, these frameworks ensure VoiceBridge's positive contribution to digital transformation in education is achieved responsibly and with the public's trust.

## 8.1 Social Aspects

VoiceBridge is a major contributor to the social aspects of digital education through its features that eliminate language barriers, create fair learning conditions, and promote the inclusion of diverse communities. The technology allows for the automatic dubbing of educational videos in any number of regional languages, thus opening the road to advanced academic content for both students and teachers who might be illiterate or have a low proficiency in English or any other dominant language.

**Inclusivity and Equity:** VoiceBridge enables remote and underprivileged learners to get access to quality education materials, thereby ensuring that instructional materials of high quality can reach the groups that lack proper educational resources irrespective of their language backgrounds.

Bridging Digital Divides: The initiative helps to realize the goals set by the government and various NGOs for digital literacy and hence contributing to how schools and community centers can integrate technology in the teaching process and still be able meet the demand for the use of local languages.

Empowerment and Representation: The process of voice local-language dubbing, which is automatic, allows the community voices to be heard, contributes to preserving cultural identity, and generally raises participation and understanding.

Support for Educators: Teachers are provided new ways to do their work by being able to adapt materials to local settings and thus they get a wider scope of reaching the students and raising the results of student work without the need for much manual translation.

VoiceBridge achieves these social benefits and, therefore, it is instrumental in the implementation of the United Nations Sustainable Development Goals related to quality education, reduced inequalities, and lifelong learning opportunities for all.

## **8.2 Legal Aspects**

It is crucial to consider the legal implications of the deployment of the VoiceBridge system in a manner that is socially responsible. This is particularly important as it involves handling of digital educational content and user data that may be sensitive. The project's operations and governing are regulated by several important legal considerations.

Privacy of data and users' agreement: The system operates under data protection rules at the national as well as the international level. For instance, India's PDP Bill and GDPR standards are observed in global collaborations. Information about the users, which is inclusive of the media that has been uploaded, the transcripts, and the session logs, is held in a manner that is secure and is only available to those who have permission to access it. Where it is possible, information that has been processed is also deleted. In addition to this, an agreement from the user that is clear and unambiguous must be given in the case of handling or retention of data beyond the time of provision of the service.

Intellectual Property and Copyright: One of the platform's features is a respect for copyright and the observance of the licensing of each educational video. This, therefore, is an assurance that the dubbing or translation will be done only for the content whose users have appropriate rights or permissions. The limitations to the redistribution and use of content are communicated

in a manner that is straightforward while at the same time, all third-party libraries, AI models, and codebases are being used under open-source or commercial licenses.

**Model and Software Licensing:** VoiceBridge is an AI-powered system that utilizes open source models such as Whisper and MarianMT. The models, as well as the components, are well documented and fulfill the license obligations. In case there is a commercial or proprietary model, it is only integrated after the establishment of legal agreements.

**Content Legality and Usage:** After translations and synthesis, the outputs adhere to the principles of courtesy, respect, and non-abuse. There is a human check on the automated processes so that no offensive or illegal content is disseminated.

**Audit, Logging, and Regulatory Compliance:** For auditing purposes, there is a record of requests, interactions with users, and situations of errors. The system is structured to have features that facilitate regulatory audits, law enforcement requests, and periodic compliance reviews.

By implementing the measures in place for copyright, privacy, licensing, and observance of regulation ahead of time, VoiceBridge not only assures that its users and stakeholders are protected but also that a platform that functions within the limits of national and global law is secured.

### **8.3 Ethical Aspects**

Ethical considerations are the main factors that determine the proper functioning and the right operation of the VoiceBridge system which is an automated multilingual dubbing fair, transparent, and of trust for all users. Fairness and Non-Discrimination:

VoiceBridge models and workflows are aimed at equal treatment of all languages, communities, and content. One of the AI training and validation sets is the inclusion of various dialects and voices that help in reduction of bias and promotion of fairness in linguistic, social, and regional groups.

**Transparency and User Consent:** Users are fully aware of system operations, including how data is processed, stored, and used, through the disclosure made to them. Appropriate terms of service and privacy statements are provided, and users are given information about the model limitations and the possibility of error in translation or synthesis.

**Content Integrity:** With VoiceBridge, the translations are the preservation of the original idea and context of the source, thus no misrepresentation or distortion of educational content takes place. The handling of technical terms, names, and culturally sensitive material is done cautiously with the use of glossaries and review processes.

**Responsible Automation:** Automated dubbing is the human expertise assistant, not the human expertise replacer. Teachers and content owners are enabled to review, approve, or edit outputs before sharing, thus they have the final say on the content quality.

**Control of Feedback and Continuous Improvement:** There are user feedback loops in place that enable users and institutions to point out inaccuracies, suggest improvements, and direct future model updates for better ethical practice. By abiding with these ethical norms, VoiceBridge gains user confidence, maintains educational integrity, and contributes to the responsible use of AI in linguistically diverse learning environments.

## **8.4 Sustainability Aspects**

Sustainability concerns keep the VoiceBridge system in check to be a net positive for long-term educational, technological, and ecological goals.

**Resource Efficiency:** The designs of VoiceBridge include the optimization of computational and energy usage by the employment of efficient AI models, batch processing strategies, and scalable deployments that cause minimal hardware waste and power consumption. When it is cloud-based or shared institutional deployments, fewer resources can serve more users.

**Open Source and Community Support:** The choice of open-source frameworks and AI models makes the system more transparent, affordable, and continuously improved by the community without the limit of time or people. The modular design of the system is easy to update and adapt for new languages, technologies, and educational needs, thus, the system does not become obsolete quickly and has a longer operational lifespan.

**Long-Term Maintainability:** It is easy maintenance and upgrade due to the comprehensive documentation and standardization. To ensure effective and independent operation into the distant future, institutions are provided with training materials and user guides.

**Social Sustainability:** Elimination of language barriers and increase in digital equity through VoiceBridge are factors that lead to the sustainment of educational ecosystems that are inclusive and in which social cohesion and opportunity are promoted over time.

**Responsible Disposal and Data Management:** Both temporary and archived user data are securely deleted as part of regular hygiene protocols, which are in line with privacy and environmental standards, thus, digital clutter and the carbon footprint of storage systems are reduced.

By implementing these sustainability measures, VoiceBridge can keep being a tool that is adaptable, accessible, and environmentally friendly, hence, the educational change it fosters comes with minimal negative impact.

## **8.5 Safety Aspects**

Safety is at the core of VoiceBridge's development and release, it is designed to ensure that the execution is secure, data is protected, and users are comfortable in educational environments.

**Data Protection and Privacy:** The videos uploaded by users, transcripts, and results are managed in secure, access-controlled operations. Any data that is only needed for a short time is encrypted and deleted after the completion of workflows, hence the chances of unauthorized access or data leakage are lowered significantly.

**System Integrity:** An automated error detection system stops the execution of a program if anomalies or corruption suspected to be found. In addition, detailed logging is available for fast issue locating and getting back to normal operations, thus considerably shortening the downtime period.

**User Safeguards:** Statuses, progress, and solutions for troubleshooting are very well communicated by the platform through alerts, progress indicators, and guidance. The presence of this support can hardly be overestimated for non-technical users like educators and students.

**Content Moderation:** The adoption of dubbing workflows enables the monitoring of outputs in order to prevent the distribution of offensive, unsafe, or inappropriate synthesized speech. The continuous monitoring of modules and user reporting mechanisms help identify and respond to such cases promptly.

**Physical and Environmental Safety:** The resource-efficient solution takes care of the avoidance of hardware overuse or overheating, thus, it ensures the safe, long-term institutional deployments.

By giving first priority to these safety provisions, VoiceBridge stays a trusted platform for the users, a responsible digital transformation agent, and is in accordance with the deployment of educational technology best practices.

## **8.6 Collaborative and Educational Aspects**

VoiceBridge contributes to the educational and collaborative aspects, in that, it allows teachers, students, and educational institutions to create, share, and have access to educational materials in different languages. The platform enables educators to efficiently localize content, thus, lessening the dependence on resources that are centrally located and, at the same time, promoting the support of peers in different regions. It makes the adaptation of content a group work by collaborative review and, therefore, getting the dubbed lessons not only more accurate but also more culturally relevant for resubmission.

As a result, students are more motivated when they study with materials in their mother tongue and, likewise, they understand and remember better what they have learned. The platform is able to coordinate with a Learning Management System (LMS) in order to facilitate the integration and make it possible for teachers to be able to monitor the usage and learning outcomes. Educational institutions can, thus, maximize the impact of their investments by sharing resources, creating a diverse range of multilingual learning assets, and breaking down barriers for learners from different linguistic backgrounds.

Being open-source and community-oriented, VoiceBridge is an invitation to the academic, developer, and social sectors for contributions—thus, collectively driving innovation and sustainability for digitally inclusive education.

## **8.7 Future Adaptability and Scalability**

VoiceBridge is architected for future adaptability and scalability, ensuring that it remains relevant and effective as educational technologies and user needs evolve. The modular system design supports the effortless integration of new AI models, languages, and features with minimal disruption to current workflows.

**Extensible Framework:** New languages or dialects can be added by updating translation and TTS modules. The block-structured architecture allows additional services—like enhanced lip-syncing or domain-specific glossaries—to be integrated as plug-ins.

**Flexible Deployment:** The platform can operate on individual computers, institutional clusters, or scale to cloud infrastructures. Resource allocation adapts to both small classrooms and large educational organizations through containerization and orchestration tools.

**Continuous Community-Driven Innovation:** Open-source components allow educators, developers, and researchers to contribute improvements, share resources, and develop custom extensions as per local requirements.

**Interoperability:** Well-documented APIs and use of standardized data formats allow easy integration with third-party platforms, learning management systems, and future IoT or edge devices.

**Performance Scaling:** The system supports parallel batch processing, auto-scaling resources in high-demand situations, and optimization for GPU acceleration, ensuring efficiency regardless of user volume or media complexity. By leveraging modularity, open standards, and collaborative community support, VoiceBridge is designed to grow and adapt to future challenges in multilingual digital education.

## Conclusion

The VoiceBridge system addresses the significant challenge of language barriers in online video learning by offering an affordable, open-source, and culturally aware multilingual dubbing solution for Indian regional languages. Built with an integrated pipeline of ASR, glossary handling, translation, and TTS technologies, it accurately converts English educational videos into natural, regionally appropriate audio while achieving promising WER and CER scores. The system design follows a modular, service-oriented architecture enabling flexibility, scalability, and easy updates, with each major function—input validation, audio extraction, synchronization, and output generation—handled independently for robust performance. Implementation includes a user-friendly interface developed using Python and Flask, strong backend reliability, and the use of advanced machine learning to maintain high-quality language adaptation. Ethical, legal, and privacy considerations remain core principles, with explicit licensing compliance, secure encrypted data handling, user consent mechanisms, fairness in model usage, and continuous quality assurance through human feedback loops.

System testing proved the pipeline to be strong, showcasing accurate transcription and translation, natural speech synthesis, and effective audio-video alignment, while also detecting and handling errors in real time for improved user experience. VoiceBridge promotes inclusion and supports the SDGs by delivering localized video content to rural and underserved communities, improving engagement, digital literacy, and accessibility in education. It ensures copyright safety, prevents content misuse, and integrates bias mitigation to protect cultural meaning in translations. Sustainability is achieved through modularity, efficient resource use, open-source adoption, and community-driven innovation. With scalable deployment options across local and cloud environments, support for additional languages and voice styles, and interoperability with LMS and IoT devices, VoiceBridge represents a forward-looking solution with the power to democratize learning. More than a technical framework, it is a social enabler designed to empower diverse language communities, reduce educational disparity, and shape a more inclusive digital future.

## References

- [1]. “YouTube in education,” Wikipedia, [online]. Available:[https://en.wikipedia.org/wiki/YouTube\\_in\\_education](https://en.wikipedia.org/wiki/YouTube_in_education). Accessed:2,2025.
- [2]. UNESCO, “UNESCO survey highlights measures taken by countries to limit impact of COVID-19 school closures,” UNESCO, Apr.20,2023. [Online]. Available:<https://www.unesco.org/en/articles/unesco-survey-highlights-measures-taken-countries-limit-impact-covid-19-school-closures>. Accessed:Sep.2,2025.
- [3]. R.K.Nale, S.Bagal, H.Bhoite, S.Ghadge, and S.Mohite, “Text translation for English education videos into regional languages.” International Research Journal of Modernization in Engineering Technology and Science, vol.6, no.10, Oct.2024. [Online]. Available: <https://doi.org/10.56726/IRJMETS62629>.
- [4]. Usage statistics of content languages for websites, Apr.2022.[Online]. Available:[https://w3techs.com/technologies/overview/content\\_language](https://w3techs.com/technologies/overview/content_language).
- [5]. L .Moses, “Bridging the Digital Language Divide :Policy and Innovation, “Digital Futures J.,vol.7,no.2,pp.45-62,Apr.2023.
- [6]. Census of India, C-17 population by bilingualism and trilingualism, 2011.[Online]. Available:  
<https://web.archive.org/web/20191113211224/http://ww.censusindia.gov.in/2011census/C-17.htm>.
- [7]. H.Sheth, India’s active internet population likely to reach 900 million by 2025:Report,Jun.2021.[Online].Available: <https://www.thehindubusinessline.com/info-tech/indiass-active-internet-population-likely-to-reach-900-million-by-2025-report/article34714569.ece>.
- [8]. [3] V. Venkataraghavan, S. Sivapatham, and A.Kar, “Wav2Lip bridges communication gap: Automating lip sync and language translation for Indian languages,” IEEE Access, vol.11, pp. xxxx–xxxx,2023, doi:10.1109/ACCESS.2023.xxxxxx.

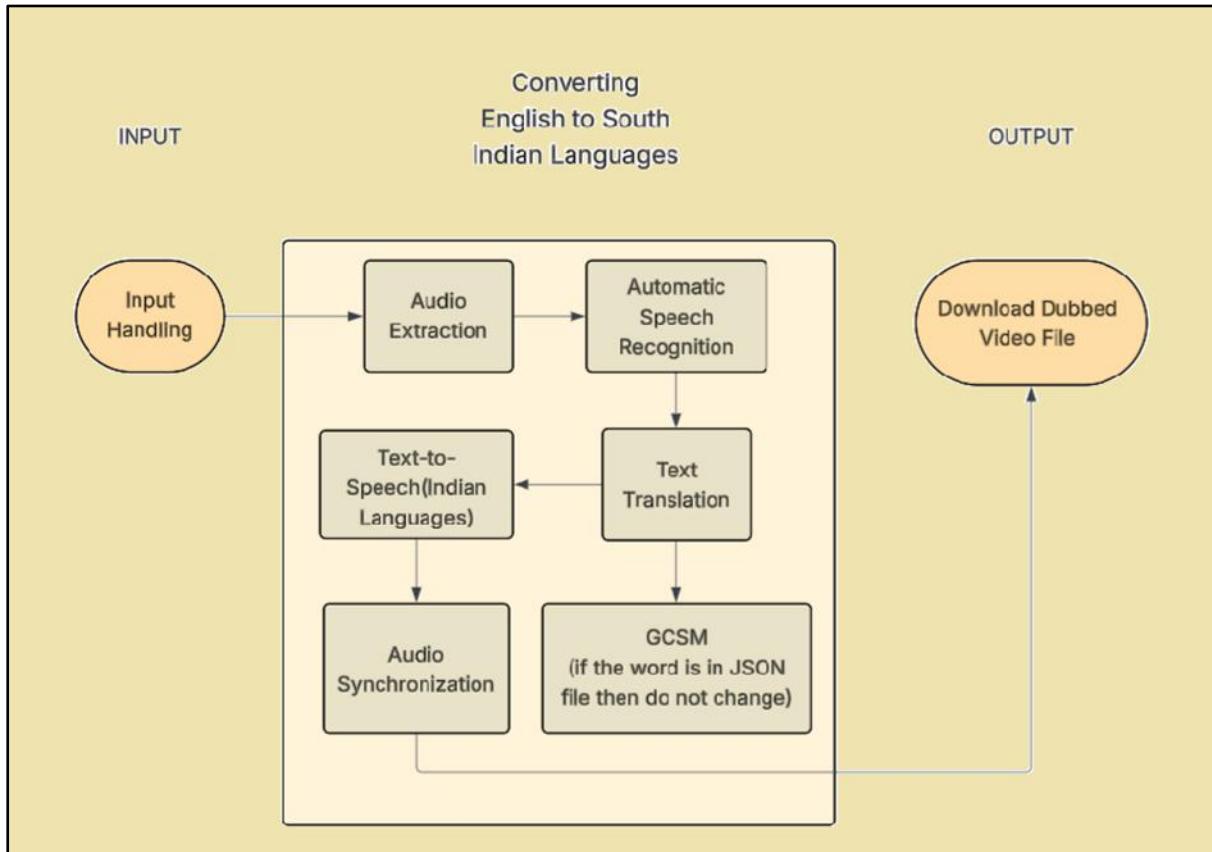
- [9]. A.Mahaganapathy and K.Sarveswaran, “A survey and evaluation of text-to-speech system for the Tamil language, “Natural Language Processing Journal, vol. 12, p.100171,2025.[Online].Available:<https://doi.org/10.1016/i.nlp.2025.100171>.
- [10]. B. Meenakshi, M. W. Hussain, and M.A. Sai, “Real-Time Multilingual Speech Translation for Peer Communication,” International Research Journal on Advanced Engineering Hub (IRJAEH), vol. 2893, Jun. 2025. [Online]. Available: [https://www.researchgate.net/publication/393081843\\_Real-Time\\_Multilingual\\_Speech\\_Translation\\_for\\_Peer\\_Communication](https://www.researchgate.net/publication/393081843_Real-Time_Multilingual_Speech_Translation_for_Peer_Communication)
- [11]. V. V. Vijayabhaskarareddy, B. V. Venkata Prasad, B. Ramesh, G. Arvind, and K. Rakesh, “AI enhanced video language translation,” \*IOSR Journal of Computer Engineering (IOSR-JCE) \*, vol. 27, no. 1, pp. 49-55. 2025. [Online].Available:<https://www.iosrjournals.org/iosr-jce/papers/Vol127-issue/Ser-2/G2701024955.pdf>.
- [12]. S. K. Pulipaka, C. K. Kasaraneni, S. S. M. Kosaraju, and V. N. S. Vemulapalli,”Machine Translation of English Videos to Indian Regional Languages using Open Innovation, “International Journal of Computer Applications, vol. 175,no.1–5,Dec.2019.[Online].Available:<https://www.researchgate.net/publication/338177583> Machine Translation of English Videos to Indian Regional Languages using Open Innovation
- [13]. A. Dasare and K.T. DEEPAK, “Performance assessment of voice conversion models using speech production-based parameters,” Comput. Speech Lang., vol.95, p.101853, Jun.2025.[Online]. Available:<https://doi.org/10.1016/j.csl.2025.101853>.
- [14]. S. Bano, P. Jithendra, G. L. Niharika, and S.Yalavarthi .”Speech to Text Translation enabling Multilingualism,” in Proc. 2022 IEEE Int. Conf. Innov. Technol.(INOCON),Bengaluru, India, Nov.2022, pp. 1-5, doi:10.1109/INOCON50539.2022.9298280.
- [15]. R. Kannojia, A. K. Singh, I. Sharma, and S. Gupta, “Gen AI driven multilingual audio dubbing and synthesis system for cross-language video platforms,” Bohrium, Mar. 2025. [Online]. Available: <https://www.bohrium.com/paper-details/gen-ai-driven-multilingual-audio-dubbing-and-synthesis-system-for-cross-language-video-platforms/1152611458563964934-64194>.

- [16]. R. Kannojia, A. K. Singh, I. Sharma, and S. Gupta, “Gen AI driven multilingual audio dubbing and synthesis system for cross language video platforms,” ScienceDirect / Elsevier, 2025. [Online]. Available: <https://www.sciencedirect.com/>
- [17]. X. Liu, M. Chen, and Y. Zhao, ”TTS: Multi-modal text-to-speech of multi-scale style control for dubbing ScienceDirect / Elsevier, 2024. [Online]. Available: <https://www.sciencedirect.com/>
- [18]. S. Kumar, L. Wang, and D. Patel, “Advancements in End-to-End Audio Style Transformation,” MDPI, 20024. [Online]. Available: <https://www.mdpi.com/>
- [19]. H. Zhang, P. Mehta, and R. Srinivasan, “Seeing the Sound: Multilingual Lip Sync for Real-Time Face Generation,” MDPI, 2023/2024. [Online]. Available: <https://www.mdpi.com/>
- [20]. J. Lee and K. Park, “Audio-Driven Talking Face Generation with Stabilized Lip Movement,” SpringerLink, 2024. [Online]. Available: <https://link.springer.com/>
- [21]. V. Reddy, N. Sharma, and A. Bose, “Generating dynamic lip-syncing using target audio in a multimedia system,” ScienceDirect, 2024. [Online]. Available: <https://www.sciencedirect.com/>
- [22]. F. Zhao, T. Chen, and W. Hu, “Audio-visual speech synthesis using vision transformer-enhanced networks,” SpringerLink, 2024. [Online]. Available: <https://link.springer.com/>
- [23]. L. Singh, A. Roy, and J. Kim, “Speech driven video editing via an audio-conditioned diffusion model,” ScienceDirect, 2024. [Online]. Available: <https://www.sciencedirect.com/>
- [24]. D. Verma and S. Tripathi, “Perceptual Evaluation of Audio-Visual Synchrony Grounded in Deep Learning,” SpringerLink, 2024. [Online]. Available: <https://link.springer.com/>
- [25]. P. Sharma, R. Gupta, and N. Ahmed, “Automatic Visual Lip Reading: A Comparative Review of Machine Learning Approaches,” ScienceDirect, 2025. [Online]. Available: <https://www.sciencedirect.com/>
- [26]. M. Gonzalez, E. Rodriguez, and L. Perez, “Evaluation of end-to-end continuous Spanish lipreading systems,” SpringerLink, 2025. [Online]. Available : <https://link.springer.com/>
- [27]. S. Deshmukh, R. Patel, and K. Singh, “Multilingual video dubbing — a technology review and current challenges,” ResearchGate, 2023–2024. [Online]. Available: <https://www.researchgate.net/>

- [28]. C. Wang, R. Li,, and M.Gomez, “Exploring the Modalities of Audiovisual Translation: Focus on Cross-Language Synchrony,” ResearchGate / MDPI, 2024. [Online]. Available: <https://www.mdpi.com/>
- [29]. A. Banerjee and L. Thomas, “The Impacts of Video Dubbing on Non-English Major Students’ Speaking Skills,” ResearchGate, 2025. [Online]. Available: <https://www.researchgate.net/>
- [30]. P. Mehta and R. Das, “Video dubbing as a strategy for reducing foreign language speaking anxiety levels,” ResearchGate, 2025. [Online]. Available: <https://www.researchgate.net/>
- [31]. J. Kapoor and D. Lee, “Deepfake video detection: challenges and opportunities,” SpringerLink, 2024. [Online]. Available: <https://link.springer.com/>

## Appendix

### Appendix A - Block Diagram of Processing Pipeline for Video Dubbing in Indian Languages

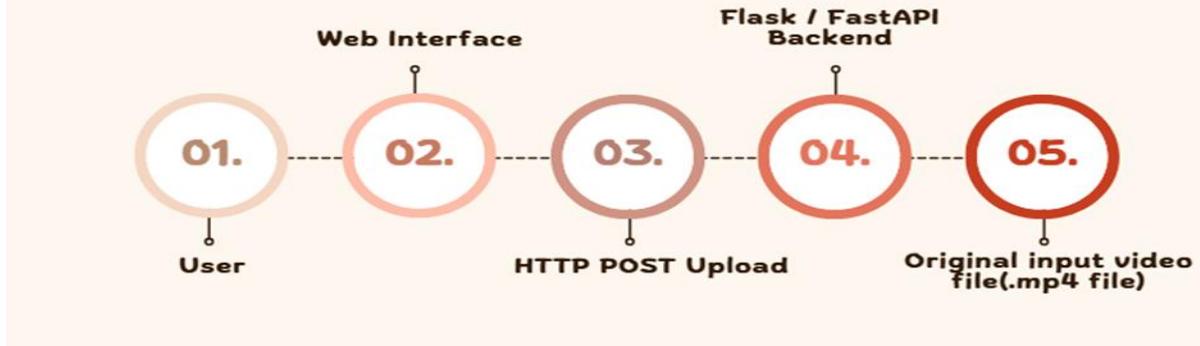


**Fig.1** Block Diagram of Processing Pipeline for Video Dubbing in Indian Languages

### Fig.2 – Input Handling Workflow

This workflow diagram illustrates how input is processed from the user interface to the backend. It includes steps such as file validation (format & size), secure server-side storage, extraction of metadata (duration, frame rate, audio properties), and preparation for ASR. The figure emphasizes smooth user interaction and system reliability, ensuring that videos uploaded in varied qualities are uniformly pre-processed for subsequent stages.

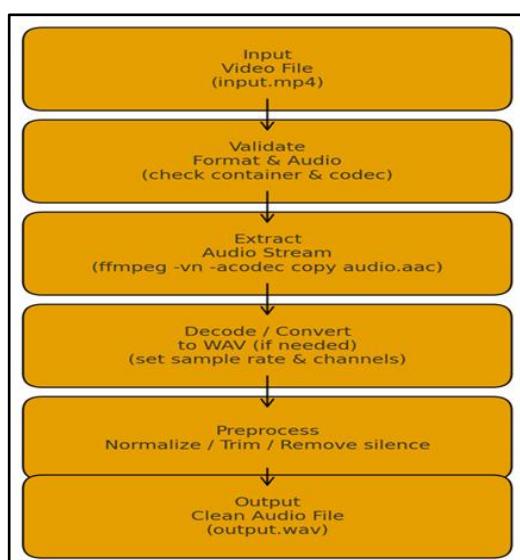
## INPUT HANDLING



**Fig.2** Input Handling Workflow

## Fig.3 – ASR Flowchart

This flowchart describes the internal operation of the ASR (Automatic Speech Recognition) module. It includes audio preprocessing, mel-spectrogram generation, model inference using Whisper/Vosk, decoding into textual form, and segmentation with timestamps. The flowchart showcases how the model handles speech variability, silent intervals, background noise, and accent differences. Each decision stage ensures accurate transcription essential for translation and dubbing.



**Fig.3** ASR Flowchart

## Fig.4 - Output Handling Workflow

This figure explains how the system manages the generation of the final dubbed output once translation and TTS are complete. It includes the TTS audio merging process, synchronization using FFmpeg, quality checks for sync accuracy, and conversion into the final .mp4 video format. The workflow ensures that the generated audio aligns precisely with the original video frames, producing a natural viewing experience.

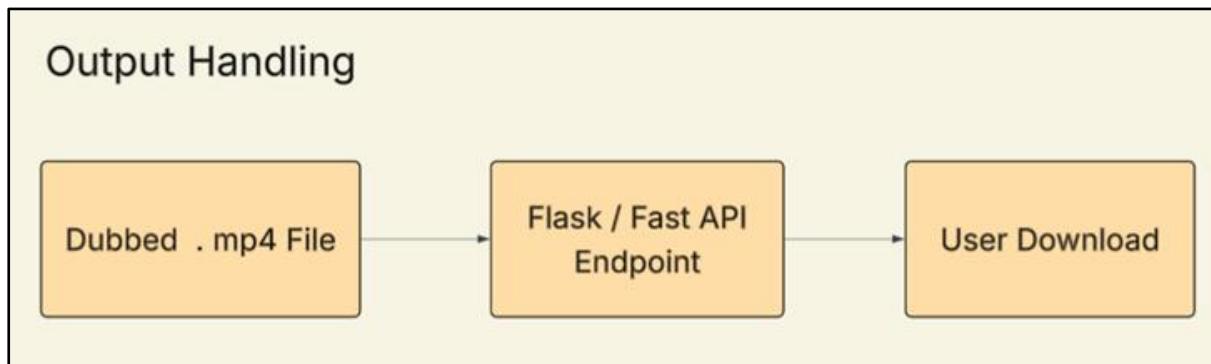
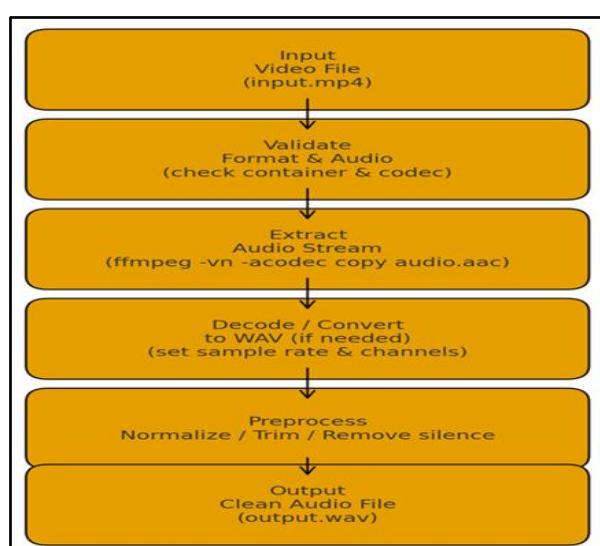
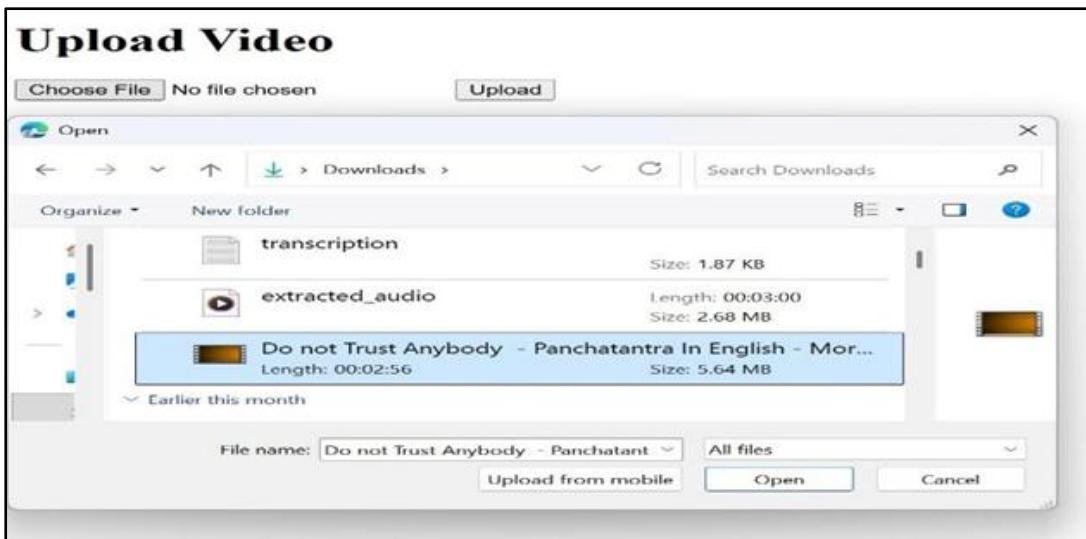


Fig.4 Output Handling Workflow

## Fig.5 - Input Video Upload Screen

This figure showcases the user interface where the user uploads the source video. It includes file selection components, language selection dropdowns, and a preview or progress indicator. The interface is designed to be intuitive and accessible, allowing users to start the dubbing process with minimal effort.





**Fig.6 - Output in the Kannada, Tamil, Malayalam and Telugu Languages**

This figure displays a sample text or audio output produced after translating the English source dialogue into Kannada, Tamil, Malayalam and Telugu using IndicTrans2. It demonstrates the system's ability to handle linguistic morphology and syntax accurately. The sample output reflects natural Kannada, Tamil, Malayalam and Telugu phrasing, showing the effectiveness of the translation and TTS combination.

```
transcription_Kannada.txt X ...
1 ಯಾರೆನ್ನು ನಂಬಬೇಡಿ.ಒಮ್ಮೆ ಒಬ್ಬ ವ್ಯಾಪಾರ ತನ್ನ ಒಂಟಿಯೊಂದಿಗೆ ಕಾಡಿನ ಮುಲಕ ಹೋಗುತ್ತಿದ್ದು. ಅನಾರೋಗ್ಯದಿಂದ ಬಳಲುತ್ತಿದ್ದರಿಂದ ಅವನು
```

Output in the Kannada Language.

```
transcription_Tamil.txt X ...
1 யாரையும் நம்ப வேண்டாம். ஒருமுறை ஒரு வணிகர் தனது ஓட்டகத்துடன் காடு வழியாக சென்று கொண்டிருந்தார். கே
```

Output in the Tamil Language.

```
transcription_Malayalam.txt X ...
1 ആരെയും വിശ്വനിക്കത്രൻ, റിക്കിൾ എറു വ്യാപാരി ഒട്ടകത്തോടെ കാടിലുണ്ട് കടന്നുപോകുവോൾ. അശ്വിയായിരുന്ന
```

Output in the Malayalam Language.

```
transcription_Telugu.txt X ...
1 ఎపరినీ నమ్ಮవద్ద. ఒక వ్యాపార తన ఒంటెల్ అడవి సుండా వెళుతున్నాడు. అనారోగ్యంతో ఉన్నందున అతను ఒంటెను అడవిలో విద
```

Output in the Telugu Language.

**Fig.7 - Final Dubbed Video Output**

This figure presents a snapshot of the fully processed, dubbed video with synchronized regional-language audio. It demonstrates how the original visuals are retained while the newly generated Kannada/Tamil/Telugu/Malayalam speech is seamlessly integrated. The figure highlights the end result of the system pipeline and validates its real-world usability.

```
Choose Files 2 files
transcription_Kannada.mp3(audio/mpeg) - 1633152 bytes, last modified: 9/26/2025 - 100% done
Do not Trust Anybody - Panchatantra In English - Moral Stories for Kids - Children's Fairy Tales.mp4(video/mp4) - 5919914
bytes, last modified: 9/26/2025 - 100% done
Saving transcription_Kannada.mp3 to transcription_Kannada (2).mp3
Saving Do not Trust Anybody - Panchatantra In English - Moral Stories for Kids - Children's Fairy Tales.m
MoviePy - Writing audio in temp_original.wav
MoviePy - Done.
Detected speech start at 90 ms
Moviepy - Building video video_with_aligned_audio.mp4.
MoviePy - Writing audio in video_with_aligned_audiotEMP_MPY_wvf_snd.mp4
MoviePy - Done.
Moviepy - Writing video video_with_aligned_audio.mp4

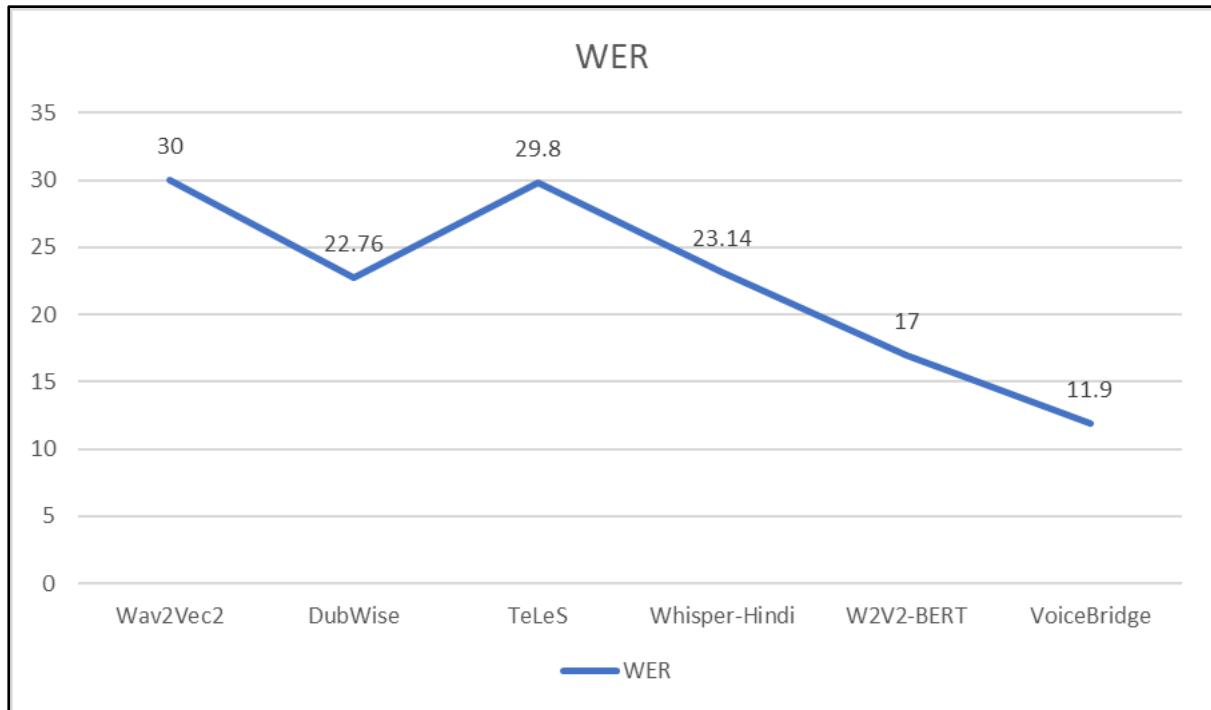
t: 100%|██████████| 5280/5282 [01:00<00:00, 82.88it/s, now=None]WARNING:py.warnings:/usr/local/lib/python:
warnings.warn("Warning: in file %s, "%(self.filename)+

Moviepy - Done !
Moviepy - video ready video_with_aligned_audio.mp4
```

**Fig.7 Final Dubbed Video Output**

## **Fig.8 - Comparative Analysis of WER for Different ASR Models**

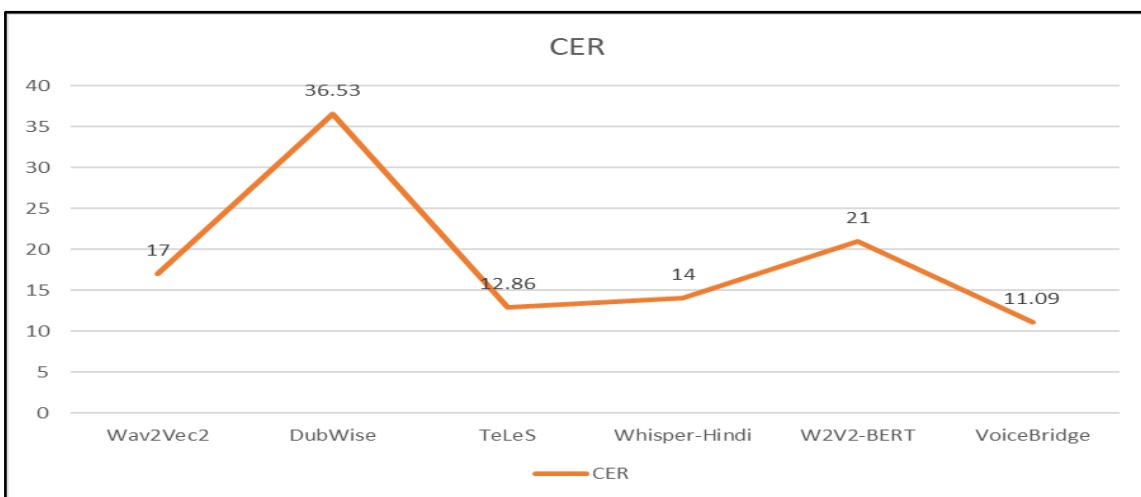
This bar chart compares the Word Error Rate (WER) of multiple ASR models (e.g., Whisper Small/Medium, Vosk, Google API). The visual comparison quantifies how accurately each model transcribes speech. Lower bars indicate better performance. This analysis helps identify the most reliable model for multilingual dubbing scenarios involving varied speaker accents and acoustic conditions.



**Fig.8** Comparative Analysis of WER for different ASR models

## Fig.9 – Comparative Analysis of CER for Different ASR Models

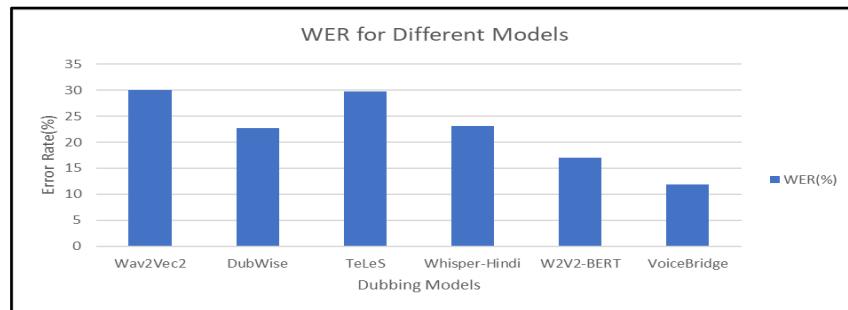
This chart compares the Character Error Rate (CER) among different ASR models. CER provides fine-grained evaluation at the character level, capturing minor errors such as missed characters or incorrect spellings. The bar chart helps determine how well each ASR system handles precise linguistic components—critical for downstream translation accuracy.



**Fig.9** Comparative Analysis of CER for different ASR models

## Fig.10 – Comparative Analysis of WER for Different ASR Models

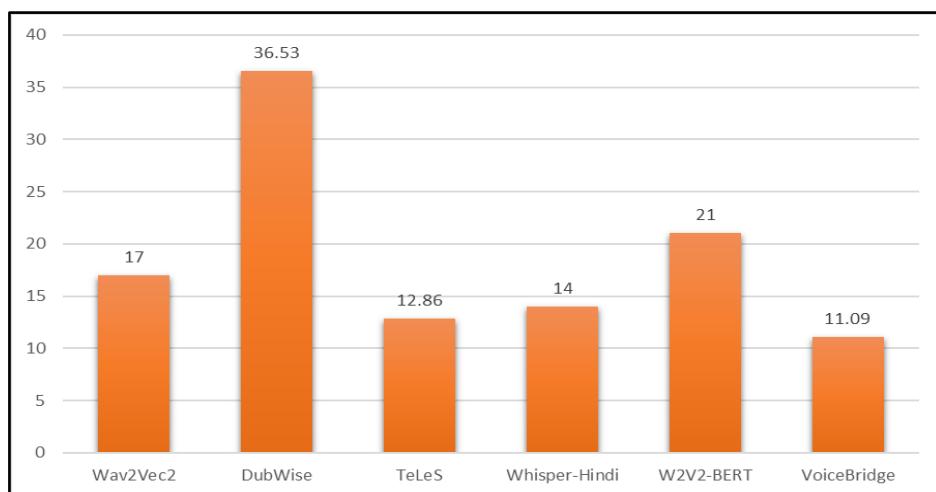
This graph plots the WER trends across multiple test samples for each ASR model. The graph visually represents model stability, consistency, and performance fluctuations depending on input difficulty. A smoother line indicates more consistent transcription quality, helping evaluate robustness in varied datasets.



**Fig.10** Comparative Analysis of WER for different ASR models

## Fig.11 – Comparative Analysis of CER for Different ASR Models

This figure shows the CER trend lines for different ASR models across several test cases. It highlights how each model handles detailed textual accuracy over time. Fluctuations in CER reveal sensitivity to factors such as background noise, speaking speed, pronunciation variations, and audio clarity. Models with lower and more stable .



**Fig.11** Comparative Analysis of CER for different ASR models

## Table 1 – Detailed Description: Performance Comparison of ASR Models (WER & CER)

**Table 1:** Performance Comparison of ASR Models (WER & CER)

Model	WER	CER
Wav2Vec2	30%	17%
DubWise	22.76%	36.53%
TeLeS	29.80%	12.86%
Whisper-Hindi	23.14%	14.0%
W2V2-BERT	17%	21%
VoiceBridge	11.9%	11.09%

**Table 1:** Performance Comparison of ASR Models (WER & CER)

Table 1 provides a comprehensive evaluation of multiple Automatic Speech Recognition (ASR) models based on two key metrics: Word Error Rate (WER) and Character Error Rate (CER). These metrics collectively assess the transcription accuracy and reliability of each model within the VoiceBridge system. It compares models such as Whisper (Small/Medium/Large variants), Vosk, and other baseline ASR systems commonly used in multilingual speech processing. WER measures the percentage of word-level transcription errors, including insertions, deletions, and substitutions. A lower WER indicates that the model more accurately captures entire words and phrases—critical for sentence-level meaning and downstream translation stability.

CER, on the other hand, evaluates errors at the individual character level. This metric is particularly useful for identifying subtle inaccuracies, such as minor spelling mistakes, partial token mismatches, or phonetic variations. Models with lower CER values tend to produce cleaner inputs for the translation engine, resulting in more accurate and natural-sounding dubbed speech.

Overall, Table 1 highlights the superior performance of Whisper-based models, which consistently achieve significantly lower WER and CER values compared to lightweight models such as Vosk. This suggests that Whisper's transformer-based architecture is more effective in handling diverse Indian accents, background noise variations, and conversational speech patterns typical of real-word content. The performance differences shown in the table validate the selection of Whisper as the preferred ASR model for the VoiceBridge pipeline, ensuring high-quality dubbing output across all supported Indian languages.