

# Enabling Multilingual Communication: Automated Lip-synchronization Dubbing for Albanian Videos

ALDO DIKU

UNIVERSITY OF NEW YORK TIRANA

JULY 2025

This is to certify that I have read this project and that, in my opinion, it is fully adequate, in scope and quality, as a thesis for the degree of Bachelor of Arts in Computer Science.

(Title and Name)

(Project Advisor)

\_\_\_\_\_

\_\_\_\_\_

This is to confirm that this thesis complies with all the standards set by the Department of Computer Science of University of New York Tirana.

Date:

Seal/Signature:

I hereby declare that all information in this document has been obtained and presented in accordance with academic rules and ethical conduct. I also declare that, as required by these rules and conduct, I have fully cited and referenced all material and results that are not original to this work.

First Name, Last Name:

Signature:

## ABSTRACT

ENABLING MULTILINGUAL COMMUNICATION: AUTOMATED  
LIP-SYNCHRONIZATION DUBBING FOR ALBANIAN VIDEOS

Diku, Aldo.

BA. in Computer Science

Thesis Advisor: Prof. Miralda Çuka

June 2025, 20 pages

Using available methods and tools (which i will mention when i gather all the tools and methods) to create lip-synced dubbed videos from an albanian video input and outputting a video of the same speaker in another language with the lips moving according to that languages movements.

Keywords: Lip-sync, dubbing, albanian language

First and foremost, I would like to express my sincere gratitude to my advisor, Professor Miralda Çuka, for their invaluable guidance, unwavering support, and insightful feedback throughout this project. Their expertise and encouragement were instrumental in shaping this work.

Finally, I am grateful for the support of my family and friends. Something something a little bit longer and better made.

# Contents

<b>1</b>	<b>Abstract</b>	<b>3</b>
<b>2</b>	<b>Acknowledgements</b>	<b>4</b>
<b>3</b>	<b>Introduction</b>	<b>6</b>
3.1	Background of the study . . . . .	6
3.2	How lip synchronization works . . . . .	6
3.3	Methodology . . . . .	7
3.4	Work Steps Overview . . . . .	9
<b>4</b>	<b>Transcription</b>	<b>9</b>
4.1	Tests . . . . .	9
4.2	Audio transcription . . . . .	10
4.3	Voice Cloning . . . . .	13
<b>5</b>	<b>Lip synch model</b>	<b>13</b>
<b>6</b>	<b>Future improvements</b>	<b>14</b>
6.1	Questions raised . . . . .	14

## 3 Introduction

### 3.1 Background of the study

The increasing globalisation of online content consumption presents a significant challenge for video creators aiming to reach diverse audiences, thereby creating the need to produce multilingual versions of their work. Being able to make only one video and have it transform into multiple languages will help creators but also the target audience as well. Cutting down on the time and money it takes to make two or more videos with the same subject but in different languages. One such case is a youtube creator named der8auer from Germany, for the same topic he creates two videos, one in english and one in german. Sometimes his videos are quite long and having to do them again in another language will tire you out. Having a model that can take your video, translate it into another language and synchronize your lips to the movement of the language of your choosing would be a tool in your arsenal. Some of the usages of this lip-sync dubbing technology are: video dubbing and translation, real-time Face-to-Face translation that will be possible in the future and multilingual communication, gaming and virtual environments reduced cost and labor and time in this case, entertainment and content creation, speech recognition and lip-reading. Film, entertainment and media production. Education and training videos, especially in cases where there is a diverse workforce. It also has promising applications in areas for speech therapy, language learning and assistive communication devices for individuals with hearing impairments [3]. Advancements in Machine Learning (ML) have made it possible to have Automatic Speech Recognition (ASR)

### 3.2 How lip synchronization works

Phonemes are the distinct units of sound in speech, and visemes are visual representations of phonemes. Think phonemes as the audio and visemes as the video. Visemes serve to map sounds (phonemes) to corresponding mouth positions or shapes. The goal is to go from the audible to the visual (visemes). Visemes group similar-looking mouth shapes, which reduces the complexity required for animation. Relating the extracted audio features to the corresponding visual mouth shapes and facial expression. Mapping phonemes down to a smaller number of visemes gives artists fewer expressions to pose. The history of lip-syncing comes from animation and in particular two-dimensional cartoons where artists had to make facial movement according to the sounds of speech [8]. Lip-synchronization is a field that has been studied since 1994 [10], where Automatic lip-sync (ALS) could be used for animating cartoons realistically making the mouth movement more smoother and have it a more natural feel and could also help in aiding people with hearing disability. Later in 1997 a new method was proposed where they used Fast Fourier Transform (FFT) for speech signal analysis, statistical measures called "moments" are used to describe the shape of the FFT. Mouth parameter measurements are used to measure the jaw position, the height of the maximum vertical opening between the lips and horizontal opening of the lips. This method did not require phonemic analysis or prior knowledge of the speech content. ObamaNet was an innovative architecture released in 2017, it used a time-delayed Long Short-Term Memory (LSTM) network to produce synchronized lip-synch videos. The network learned to map raw audio feature to corresponding mouth shapes based on hours of high-quality videos of a specific target person [23]. Convolutional Neural Networks (CNN) played an important role in lip-synchronization as it helped in feature extraction for

visual features, and image processing and generation. CNNs are often used with other models like LSTM blocks, or encoder and discriminator networks based on a Generative Adversarial Network (GAN). A study by Pawar et al. (2024) focused on Marathi, an Indian language which has a shortage of datasets, similar to that of the Albanian language [17]. Integrating CNNs with adversarial training and including RNN-based architectures, specifically LSTM for lip-synch generation, Song et al. (2019), Sadoughi et al. (2021) and Li et al (2021) also used a combination of CNN, RNN, GANs and LSTMs [7] [13]. Another similar study by Exarchos et al. (2024) for lip-reading in Greek language tackled the scarcity of datasets available. The study proposed a combination of 3D CNNs and LSTM networks for word recognition from lip movements.

### 3.3 Methodology

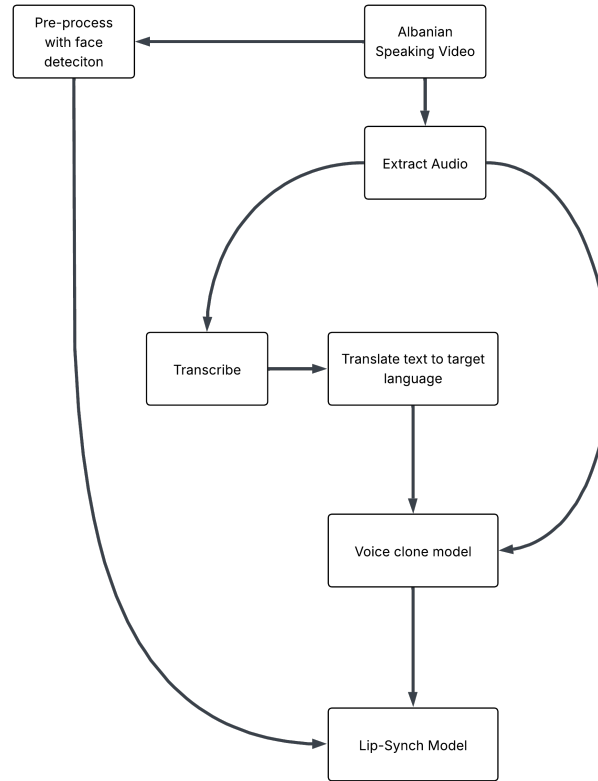
A key limmitation encountered was the lack of suitable Albanian video datasets, specifically open-license content with direct-to-camera speech, reflecting the broader challange of working with low-resource languages in multimedia applications. To address the first issue, custom video contennt was created featuring self-recorded footage in a one by one aspect ratio and lowest camera setting to save. The first step for this project was the data creation, seeing as there is not a lot of free and publicly available free-use content I can use we had to create the data myself. Videos were recorded at 1440x1440 resolution (1:1 aspect ratio) at 30 frames per second, with durations ranging from 10 seconds up till 5 minutes. The study at (citation needed for the optimal length of time of the videos) found that the optimal length time was for videos with one word in them. To adapt to this, a script was made to create small chunks of videos to be closer to that optimal format. Videos could have been created by following and reading a script thereby avoiding the transciptiton step, however we want to create a workflow that will take any video of an albanian speaker and turn it into a dub lip synched video. **This step should be removed if not implemented at the end. Also need to add the study for the dataset with one word that showed better results, I think it was Wav2Lip**

Audio was then extracted from the videos using ffmpeg package, Pulse Code Modulation, Signed 16-bit, Little Endian (pcm\_s16le), a lossless, uncompressed audio codec format, 16000hz sample rate and monophonic, meaning single audio channel. which will be used for transcription and for the base of the voice cloning tool before inputing it into the final lip synchronization model. For the audio extraction ffmpeg library was used, and audio was extracted in a loss-less codec, 16000hz sample rate and signed 16-bit Pulse Code Modulation (PCM), per Google Speech-to-text API documentation [2].

The audio gets sent through an Automatic Speech Recognition model, converting spoken language into written text. and gets saved in a text file, each line contains the name of the file transcribed and a colon separates it from the text. After several tests of three models, Google Speech to Text Chirp 2, Kushtrims and Neura, the latter one was selected as being the most accurate [2] [11] [15].

Transcribed text is translated to a target language that is supported by the voice cloning model. Using deep-translator python package and utilizing google translate as its core. Audio of the same file is used as reference to make a voice clone using OpenVoice, in the target language, a check is made if the cloned output is longer or shorter than the reference audio [18]. Another clone is made with the speed of the audio modified to be as close as possible to the original. OpenVoice offers emotion and accent cloning offering flexible voice control. The translated text





**Figure 1:** Work Flow

will be used as the base of what the voice clone should output. OpenVoices Zero-shot cross-lingual voice cloning can generate speech that can be not present in the multi-lingual training dataset used by the model. There is "native" support and generalisation of voice clones. The native support is offered for languages like english, spanish, french, chinese, japanese and korean in the second version of OpenVoice.

The reference video is processed to be cropped in a 512x512, 25 frames per second format with the face centered, this was made using a face detection model which would detect faces in every frame of the video, from there we take the center and crop accordingly. Googles mediapipe framework was used for the face detection task [4]. This format is a pre-requisite before inputting videos into the final lip-synching model. Pre-emptive tests were done on the videos to see if the face detector worked as intended and could detect a face in every frame of the videos.

After translating, the translated text and the original audio are used to make the voice clone. We need the original audio of the transcript because of emotions displayed in the audio can be used by the voice cloning tool.

Cloned voice and the original video from where the cloned voice was based on are then used to train the lip synchronization model.

**One of the challanges of creating a lip synchronization in videos is the need for the lip movements to accurately align with a specified target speech segment, especially in multilingual and unconstrained environment. Visual data falling out of sync with updated audio and inaccurate lip movements in target videos.**

CNN and Generative Adversial Networks (GANs) to create the lip movements that sync up,

this combined with a Discriminator Network which would try and detect the GANs fake lip movements and real ones, pushing each other to get better.

### 3.4 Work Steps Overview

From the video we need to get the transcription in albanian, for this step several models were tried and tested although only the google cloud API were the best performing one. The next and easiest step to implement is translating the transcription to the target language, translation has come a long way and there are many tools which can achieve high accuracy, we decided to use google translate. The following task is to take the translated transcript and use a voice cloning tool/model to make the voice dubbing. **There are several models available which have not been tested by me yet.** After this we have to use object tracking to track the mouth area in the video to then use the deep learning model for lip frame creation. **This depends heavily on the model the will be applied.**

## 4 Transcription

### 4.1 Tests

From early review there was very little research done on the subject of Automatic Speech Recognition (ASR) for Albanian language, (Florjan citation not sure if it should be applied) Albanian-ASR project in github was a model trained with 40 hours of high quality audio data, using DeepSpeech architecture and additional mechanisms such as Attention mechanism and context-dependent phoneme modeling to enhance the accuracy of the Albanian language ASR system. The low ammount of hours led the model to have a low accuracy of 46.3% which cannot be used in real-world applications.

Whisper made by OpenAI had support for Albanian language but its transcription capabilities were also not able to perform much better on Albanian language [16]. Trying to fine tune the model with the 40 hours of albanian audio provided by the Albanian-ASR proved challanging when trying to fine-tune their largest whisper-large-v3 model (1.55 Billion parameters). With the current hardware being an Nvidia 3080ti 12GB graphics card, it takes 177 hours for the largest model to train. Only when we moved on to the small model (244 Million parameters) was the training time less than 24 hours. The accuracy of the whisper-small model with the fine-tuned albanian audio dataset also did not have enough accuracy for real world usage.

Wav2Vec2 was also considered as an option and fine-tuned with the 40 hours of audio from Florjans datasets. However upon testing, the model's accuracy was not suitable for the project. Google Cloud offers a Speech-to-Text API that has support for albanian language [2]. Their model Chirp 2 seemed exceptional when using their transcription service using the user interface, but there were significant problems implementing their API for it to have the same performance when transcribing being the same as the performance from their user interface. The API service gave a confidence score that was very usefull to use as something above 0.9 was very accuracy while you could ignore the samples below that threshold. Another usefull feature of the speech-to-text API was the individual word confidence score and timing of when each word started and ended, similar to a SubRip Subtitle file (SRT file).

Kushtrim/whisper-large-v3-turbo-shqip is a model in huggingface made by Kushtrim Vioska

also based on the whisper model. It was fine-tuned with 200 hours of diverse and well annotated audio, included was gegh dialect. As of now this model was not available to host on local machines but could be used in HuggingFace spaces where you could connect it with a Gradio API. The speed of audio transcription changed based on the size of the file and the availability of the resources as it was open to everyone. The accuracy was incredibly good, (needs some short testing, with actual transcripts being compared). Unfortunately there was no confidence score or word timings.

Kushtrim seems to have updated the model just recently, adding another 100 more hours of audio making the total to "300 plus" in the Huggingface space.

As the better models, Kushtrims and Neura transcribed the same audio files. Scores were compared and Neura had the best performance. With a WER score of 0.0 while Kushtrim managed to get a WER score of 0.1193. The difference is very small and the performance of both models is very good, although Kushtrims model, being only in HuggingFace space, is more volatile because of the amount of requests it gets and the speed of the transcription. The model has gotten updated recently and now has 300 plus hours of audio.

## 4.2 Audio transcription

The first step of the project is extracting audio from the videos and this was done using the open-source ffmpeg package. Going by the suggestions for optimizing transcription quality through googles speech-to-text API documentation [2], best practices, the audio was sampled at 16000 Hz, converted into a signed-integer bitrate, and in a lossless pcm\_s16le format. One group of audio files were made in a lossless format, signed 16bit PCM (linear16), 16000hz sampling rate. While the other folder, was made with the same codec but with 48000hz sampling rate.

There is no open-source ASR model for Albanian language, and the amount of Albanian audio dataset is very low. Florijan (2023) managed to collect 40 hours of diverse accents, dialects and speech styles of Albanian audio, the purpose was to exceed an initial accuracy target of 45% [19]. This data was filtered out of 133 hours available, going through validation accuracy above 0.9 used for training. The base model used for this was Deep Speech 2 by Baidou [5], consisting of a combination of CNNs and RNNs. The model did managed to reach an accuracy of 46.3%, passing the threshold but the model is not suitable for real world application. To test if there was another better model suitable for the dataset used by Florijan we used the Whisper model by OpenAI [16]. At first we tested the biggest model out of the box to check its performance, There are seven distinct whisper models whos difference lies in the parameter size and support for english only. Whisper is trained on 680,000 hours of labeled audio data, 438,218 hours of English language only, 117,000 hours covering 96 languages, Albanian included, though it does not mention the amount of hours for each language. Audio was sampled at 16000hz in mono, similar to that of Googles documentation. Its English ASR performance was very close to that of a human. Models range from tiny (39 million parameters) to large-v3 (1550 million parameters) these models were tested, attempts to fine-tune the large-v3 model with the dataset from Florijan we fruitless as the training time for such large parameter model was over 177 hours on a RTX 3080Ti 12GB Nvidia graphics card and 32GB of RAM. All models were tested from ones with the biggest parameter to the smallest, the only model that we could feasibly train was the whisper-small with a little bit more than 8 hours. ASR systems are evaluated using a

Word Error Rate (WER) score, measuring the number of errors (substitutions, insertions and deletions) needed to transform the output of the ASR system into the reference transcript or ground truth. A lower WER score indicates a better performance in 1 we can see the formula used for the calculation. The fine-tuned whisper-small model achieved a 61.52% WER, not usable in any application. Testing the accuracy of the whisper-large-v2 model as a baseline reference without any fine-tuning resulted in a 55.1% Word Error Rate (WER) accuracy, the larger whisper-large-v3 model did not have much change and managed to score a 49.45% WER, both of these not up to par. Facebooks Wav2Vec2 baseline model with 960 hours of training data was tested and managed to achieve a 97.74% WER [1].

$$WER = \frac{S + D + I}{N} = \frac{S + D + I}{S + D + C} \quad (1)$$

Where:

- S is the amount of substitutions,
- D is the amount of deletions,
- I is the amount of insertions,
- C is the amount of correct words,
- N is the amount of correct words.

From this point onward the search continued into closed systems, offered commercially. Googles Speech to Text API offered support for Albanian language and had free credits for first time user. Their system gives a confidence score for each transcription section and confidence score in individual words aswell. After transcription, a subtitle file format was also available to check the timings of when each word or sentence is spoken, this is an incredibly usefull feature as we could use it to cut the videos into smaller chunks as Pawar et al. (2024) experienced accelerated training process and improved performance when utilising a "vVISWA" dataset composed of isolated words [17]. In the first test Googles Speech to Text API, using Chirp 2 as a model, managed to achieve a WER score of 10.7%. Kushtrims model based on the whisper-large-v3-turbo and fine-tuned with 200 hours of curated Albanian audio with diverse dialect, hosted on the Huggingface Spaces, the model weights were not opensource, nor was the dataset used for the training [11]. It did not offer a confidence score nor a SRT file with audio timings. The model could be accessed for inference using the API available from Gradio, which is how the transcriptions were processed, however the model was quite slow in times as more than one individual would try to transcribe audio at the same time, increasing the transcription time a lot, in some cases taking more than three hours for twelve seconds of audio. At a later point the model was updated with an additional 100 hours of Albanian audio. Both versions of the model were evaluated, in table 1 we can see the scores for all models, and Kushtrims managed to get a 16.46% WER. The other ASR model tested was Neuras [15], a commercially available product, we were granted access and free credits. Their website advertises the model as having a 93% accuracy, which a bit of an underestimation on their part as it scored a 3.09% WER which is equivalent to 96.91% accuracy. The first set of audio files contained casual speech and did not have language that would be deemed difficult to understand or hear, testing was done for ten files, comparing the ground truth to the hypothesis file.

WER as a metric has its limitations, from the formula we can see that the substitutions, deletions and insertions have the same weight. A minor misspelling might have the same weight as a word that completely changes the meaning of the sentence. It also doesn't take in consideration a mistake in letters inside a word. For example, the word "duket" was incorrectly transcribed as "duhet", in this case the whole word is a mistake. All unnecessary elements are removed, punctuation, making numbers into 123. The core of the WER calculation is finding the minimum number of edits needed to transform the hypothesis (transcription) text into the reference (ground truth) text. This is typically done using a dynamic programming algorithm, specifically the Levenshtein distance algorithm, applied at the word level.

**Table 1:** Error Analysis for Challenging Albanian Scripts

File ID	Neura			Kushtrim			Google API		
	S	D	I	S	D	I	S	D	I
File 1	0	0	0	2	0	0	1	0	0
File 2	0	0	0	8	0	0	7	0	0
File 3	0	0	0	27	4	3	14	17	1
File 4	27	7	4	53	15	24	25	7	3
File 5	2	0	1	31	3	7	18	7	6
File 6	0	0	0	6	0	0	2	0	0
File 7	0	0	0	12	2	6	13	2	6
File 8	0	0	0	6	0	1	4	0	0
File 9	2	0	0	12	0	3	11	0	5
File 10	2	0	0	10	3	2	5	2	0
<b>Total Errors</b>	<b>33</b>	<b>7</b>	<b>5</b>	<b>167</b>	<b>27</b>	<b>46</b>	<b>100</b>	<b>35</b>	<b>21</b>
<b>Overall WER</b>	<b>3.09%</b>			<b>16.46%</b>			<b>10.70%</b>		

A second test was made reading along a script made with the purpose of being a difficult dialog. Five new videos were created, audio extracted and we tested Googles Chirp 2, Kushtrims and Neuras model as having the best results for a further analysis. Table 2 shows the results from the tested audio, very different results shown from the first test, Neuras accuracy dropped close to their advertised numbers, Kushtrims model was the one that moved less, although we have to keep in mind that this was the model with 300+ hours of training data, despite this the models performance decreased. An anomaly was detected with Kushtrims model output because the number of Insertions is 52, in the transcription there were added words "Piii" after the end of the audio file which was not present in the other models. Removing these anomalies the models performance improves to 13.01%, the number of insertions dropped from 52 to 30. Googles Chirp 2 had the best results, achieving a 5.17% WER score, this could be the result of having a larger and more diverse dataset compared to Neuras model which might have more audio content with everyday speech rather than scientific language.

Implementation of Googles API was the most challenging of the three models, as there were many different configurations, and the results varied when accessing it through the API. While Google's web UI offered more consistent transcriptions, the process required manual effort.

**Table 2:** Second Analysis WER score

File ID	Neura			Kushtrim			Google API		
	S	D	I	S	D	I	S	D	I
File 1	3	0	1	4	0	0	4	0	0
File 2	8	0	1	4	0	0	4	0	0
File 3	5	0	0	6	2	52	2	17	1
File 4	5	1	1	4	1	1	5	1	1
File 5	5	0	1	6	0	1	2	0	2
<b>Total Errors</b>	<b>26</b>	<b>1</b>	<b>4</b>	<b>24</b>	<b>3</b>	<b>54</b>	<b>21</b>	<b>1</b>	<b>3</b>
<b>Overall WER</b>	<b>7.08%</b>			<b>18.49%</b>			<b>5.17%</b>		

### 4.3 Voice Cloning

Should add a bit of history for the voice cloning, different models, architectures, its uses etc.

Voice cloning is a subtask of speech synthesis, using deep learning to create speech imitating a specific voice. Prioritizing the preservation of the identity of the target speaker. Modern systems use a combination of speaker embedder (or encoder) to capture unique voice characteristics, a synthesizer to predict acoustic representations from text and the speaker embedding, and a cocoder to convert these representations into speech.

The NAUTILUS system from [14] has the ability to clone untranscribed speech, using a small amount of untranscribed speech, cited as about five minutes. Experiments comparing performance with different amounts of adaptation data show that Mean Opinion Score (MOS) for naturalness and similarity can improve significantly from five seconds to 15 seconds of untranscribed audio, plateauing or showing only minor improvement beyond that (e.g. 30 or 60 seconds) [21].

Noisy audio can negatively impact the quality of generated speech, a real-time voice cloning system notes that generated speech can have low intelligibility due to artifacts, such as murmurs, hums, or noisy audio [6].

Because of different training and test datasets and metrics used, MOS for example is a subjective measurement. A paper in 5th of May 2025 was released by [22] which created a foundation model for real-time autonomous interaction and voice role-play.

\*Show differences between intra-lingual and cross-lingual models in voice-cloning research\*.

Using the translated text into the language of our choice we can then try and use a voice cloning model to make the voice dubbing. A lot of models were available for this task, Nari Labs released Dia, in 22nd April 2025, coqui-ai, which a lot of other models are based on.

Fine-tuning the voice cloning models since they might respond better to more of my voice.

## 5 Lip synch model

Before deep learning lip-sync was achieved using viseme-based and rule-based system, mapping phonemes to corresponding visual mouth shapes. These methods were more commonly

used in cartoons as they lacked the realism and fluidity of natural speech [9] [12] [20].

## 6 Future improvements

### 6.1 Questions raised

List of questions:

- How do we decide if we are going to use a transcription or not?
- How do we know if the transcription is correct programatically?
- How should the data be prepared for the lip-synchronization model?
- Video format, 30 frames per second or 25, quality of the video, resolution, etc.
- Audio format, mp3, lossless, bitrate, etc.
- How does the duration of video chunks used during preprocessing impact the performance of the final machine learning model, and what chunk length yields the best trade-off between information retention and computational efficiency?
- How are voice cloning models evaluated, as it feels like it should be done with a human evaluation and that is faulty most of the time?
- Which would be the best model for the lip-synchronization task based on the data we have?
- Best way to handle differences between original audio timing and the voice cloned sample timing. (slow down the voice? speedup)
- Is an LLM better than a traditional tool for translating? Neural Machine Translation (NMT) vs Large Language Models.

AS for the optimal length of the video chunks, the study done by [17] faced challenges with lengthy continuous speech, including time-consuming merging issues and potential overfitting. This was addressed by using the "vVISWA" dataset which contains isolated words or independent speech. This improved the performance of the model a lot and reduced overfitting.

List of challenges:

- Low resources for the datasets in albanian, in video and audio.
- No open source models for the transcription of the audio from albanian.
- Voice cloning models being evaluated using a human evaluation and that is biased.
- Storage of data, both in the cloud and locally.
- Computing power, GPU, TPU, etc.
- Training time.



- If the spoken albanian is with a heavy dialect, an Albanian ASR trained with large number dialects is needed
- If the cutting of videos is not done correctly, it is possible to cut a word in half
- Translating the transcribed albanian dialect. How possible it is? NLP script transformation + translate\*
- Voice cloning emotions.
- Video chunking first then transcription, or transcription (have to keep timings) then chunking - might be more possible and more cost effective than the NLP to find the problems in the language/chunked words but not for dialects.
- Could probably use the timings from google speech-to-text api (or any timings) to make the cuts since it offers that. Making the better option to be transcript first then cut/chunk second.
- Multiple Ethical challenges
- NLP + basic quantized LLM for better albanian understanding and translation?
- 

List of improvements:

- Better dataset, quantity and quality, diversity (age, gender, moustache, skin color, lip color- lipstick, ethnicity, also in this case, dialect)
- Separation of dataset into single word video for better model performance.
- Albanian ASR open-source model, google speech-to-text API is paid service, and Kushtrims uses huggingface spaces which is slow
- Comparison between Google Speech-to-text API and Kushtrim ASR
- Addition on an NLP to handle transcription errors
- NLP for dialect speech, to transform it into something that can be translated, maybe LLM is better in this case.
- Better voice cloning. More emotional expression
- Better face reproduction from the emotions of the voice clone.
- Real time lip-synch dubbing
- Separating the dataset into two groups the new dataset with videos close to one word per video and see the results of that and the original dataset with longer videos.
- NLP can also be utilized in the audio length check, if we have too much of a difference in audio length, we could apply NLP-s to change the sentence into something shorter or longer but keeping the sentence about the same.



List of things to not forget to add:

- Mention confidence score in google api, it also offered timings like SRT file
- Kushtrim first 63 videos were on 200 hours of dataset, and the next 5 "hard" videos were on 300 hours. There can be a slight comparison on the accuracy. Should retest it.
- Write about the changes in translation from NTL to LLM being better, keeping context etc.
- Mention the fact that there is no open-source Albanian ASR model, no large open-source free-use dataset either.
- 

## References

- [1] Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations, 2020.
- [2] Google Cloud. *Speech to Text API*.
- [3] Themis Exarchos, Georgios N. Dimitrakopoulos, Aristidis G. Vrahatis, Georgios Chrysosvitiotis, Zoi Zachou, and Efthymios Kyrodimos. Lip-reading advancements: A 3d convolutional neural network/long short-term memory fusion for precise word recognition. *BioMedInformatics*, 4(1):410–422, 2024.
- [4] Google. MediaPipe Face Detector — Google for Developers. [https://ai.google.dev/edge/mediapipe/solutions/vision/face\\_detector](https://ai.google.dev/edge/mediapipe/solutions/vision/face_detector), 2024. Accessed: 2025-06-11.
- [5] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, and Andrew Y. Ng. Deep speech: Scaling up end-to-end speech recognition, 2014.
- [6] Weixin Hu and Xianyou Zhu. A real-time voice cloning system with multiple algorithms for speech quality improvement. *PLOS ONE*, 18(4):1–14, 04 2023.
- [7] Diqiong Jiang, Jian Chang, Lihua You, Shaojun Bian, Robert Kosk, and Greg Maguire. Audio-driven facial animation with deep learning: A survey. *Information*, 15(11), 2024.
- [8] Arien Kock Jonathan Gratch. An evaluation of automatic lip-syncing methods for game environments. Technical report, University of Twente, Department of Computer Science, University of SOurthern California, Institute for Creative Technologies, 2022.
- [9] Jonathan Kock, Arien ; Gratch. An evaluation of automatic lip-syncing methods for game environments. Technical report, UNIVERSITY OF SOUTHERN CALIFORNIA LOS ANGELES, 2005.

- [10] B.E. Koster, R.D. Rodman, and D. Bitzer. Automated lip-sync: direct translation of speech-sound to mouth-shape. In *Proceedings of 1994 28th Asilomar Conference on Signals, Systems and Computers*, volume 1, pages 583–586 vol.1, 1994.
- [11] Kushtrim. Whisper large v3 turbo shqip. <https://huggingface.co/spaces/Kushtrim/whisper-large-v3-turbo-shqip>, 2025.
- [12] John Lewis. Automated lip-sync: Background and techniques. *The Journal of Visualization and Computer Animation*, 2(4):118–122, 1991.
- [13] Xiaohong Li, Xiang Wang, Kai Wang, and Shiguo Lian. A novel speech-driven lip-sync model with cnn and lstm. In *2021 14th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI)*, page 1–6. IEEE, October 2021.
- [14] Hieu-Thi Luong and Junichi Yamagishi. Nautilus: a versatile voice cloning system, 2020.
- [15] Neura. Neura.al. <https://neura.al/>, 2025.
- [16] OpenAI. Robust speech recognition via large-scale weak supervision. Technical report, OpenAI, 2022.
- [17] Diksha Pawar, Prashant Borde, and Pravin Yannawar. Generating dynamic lip-syncing using target audio in a multimedia environment. *Natural Language Processing Journal*, 8:100084, 2024.
- [18] Zengyi Qin, Wenliang Zhao, Xumin Yu, and Xin Sun. Openvoice: Versatile instant voice cloning. *arXiv preprint arXiv:2312.01479*, 2023.
- [19] Florijan Qosja. Development of an albanian language transcriber using artificial intelligence. <https://github.com/florijanqosja/Albanian-ASR>, 2023.
- [20] Amirkia Rafiei Oskoei, Mehmet S. Aktaş, and Mustafa Keleş. Seeing the sound: Multilingual lip sync for real-time face-to-face translation. *Computers*, 14(1), 2025.
- [21] Tasnima Sadekova, Vladimir Gogoryan, Ivan Vovk, Vadim Popov, Mikhail Kudinov, and Jiansheng Wei. A unified system for voice cloning and voice conversion through diffusion probabilistic modeling. In *Interspeech 2022*, pages 3003–3007, 2022.
- [22] Yemin Shi, Yu Shu, Siwei Dong, Guangyi Liu, Jaward Sesay, Jingwen Li, and Zhiting Hu. Voila: Voice-language foundation models for real-time autonomous interaction and voice role-play, 2025.
- [23] Supasorn Suwajanakorn, Steven M. Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing obama: learning lip sync from audio. *ACM Trans. Graph.*, 36(4), July 2017.