



Artificial Intelligence Department

Faculty of Computers and Artificial Intelligence

Cairo University

Sign Language Translator

Under Supervision of:

Dr. Amin Allam

TA. Mohamed Ramadan

Prepared by:

Samy Gerges Ibrahim 20190235

Reem Ahmed Khalil 20201075

Aya Ashraf Mohamed 20200106

Academic Year

2023/2024

Table of Contents

1. Table of Contents.....	2
2. List of Illustrations.....	3
3. Abstract	4
4. Introduction.....	4
4.1. Problem	4
4.2. Objective.....	5
4.3. Solution	5
4.4. Scope	6
4.5. Methodology	6
4.6. Competitors.....	7
5. User Centered Analysis.....	7
5.1. Functional Requirements.....	7
5.2. Non Functional Requirements	8
5.3. Project Stakeholder.....	8
5.4. Time Plan	9
6. System Analysis & Design	10
6.1. Use Case Diagram.....	10
6.2. User Stories & Prototype	10
6.3. Entity Relationship Diagram (ERD).....	18
6.4. Class Diagram.....	19
6.5. Sequence Diagram.....	20
6.5.1.Registration.....	20
6.5.2.Sign to text.....	21
6.5.3.Text to sign.....	22
6.5.4.Sign chatbot.....	23
6.6. System Architecture	24
7. Implementation	24
7.1. Dataset.....	24
7.2. Techniques.....	24
7.2.1.Text to sign model.....	24
7.2.2.Sign to text model.....	27
7.2.3.Chatbot.....	30
8. References.....	31

List of Illustrations

Table 1: competitors

Figure 1: time plan

Figure 2: use case diagram

Figure 3: login interface

Figure 4: missing field interface

Figure 5: sign up interface

Figure 6: choose your service interface

Figure 7: sign to text interface

Figure 8: recognition results interface

Figure 9: chatbot interface

Figure 10: text to sign translation interface

Figure 11: ERD

Figure 12: class diagram

Figure 13: register sequence diagram

Figure 14: sign to text sequence diagram

Figure 15: text to sign sequence diagram

Figure 16: sign chatbot sequence diagram

Figure 17: system architecture

Table 2: WLASL dataset description

Figure 18: Previous Video architectures

Figure 19: VKNet architecture

Figure 20: head network architecture

Figure 21: Two stream 3D ConvNet

Figure 22: the Inflated Inception-V1 architecture (left) and its detailed inception submodule (right)

Table 3: comparison with previous works on WLASL

Abstract

Sign language translators play a crucial role in enhancing communication accessibility for individuals who use sign languages. However, these applications face several challenges in meeting the expectations and needs of their users.

One primary challenge revolves around ensuring fair and equal accessibility for deaf individuals in today's programs and services. Additionally, users find it quite hard to find one translator that translates both ways with accurate and fast translations.

Sign language translator solves those problems by providing a user friendly interface that translates between sign language and text/voice bi-directionally, which makes the translation process easy and fast both ways. Additionally, Sign language translator includes a chatbot that is designed especially for deaf users to access its services with their own native language to make them feel more included.

Introduction

Problem

In a world where communication is predominantly verbal, deaf individuals face barriers in interacting seamlessly with the hearing population, and vice versa. This project addresses this challenge by introducing a revolutionary communication platform designed to foster inclusivity. The main goal is to make communication more fluid and natural between deaf and hearing individuals. By seamlessly integrating text, spoken language, and sign language, this technology fosters an environment where conversations flow effortlessly, creating a more connected and inclusive world for everyone.

The incorporation of a chatbot to facilitate communication in sign language serves a profound purpose. By allowing deaf individuals to interact with a chatbot using their native sign language and receiving responses in the same medium, we aim to convey a powerful message of inclusivity and recognition.

This approach not only breaks down communication barriers but also reinforces the idea that every individual, regardless of their hearing abilities, is valued and considered in the realm of technology. Beyond empowering the deaf community, it's about making technology a tool that enhances communication, ensuring that it is not just accessible but genuinely responsive to the diverse needs of its users.

Objective

By the end of this project a sign language translation program would be designed and developed based on computer vision, NLP and machine learning models to address the current challenges facing users such as the difficult communication with sign language speakers and sign language speakers not being included in today's services.

The platform will be designed with a well-structured interface to allow users to easily navigate between different services such as bi-directional sign language translations and communication with sign language chatbot. By addressing these main points, the platform aims to improve the overall communication experience for users and become a competitive player in the industry.

Solution

The solution to the communication barriers between deaf and hearing individuals lies in the creation of an innovative communication platform. This platform seamlessly integrates text, spoken language, and sign language to make interactions more fluid and natural. By providing a user-friendly interface that accommodates multiple communication modes, the technology ensures that conversations can effortlessly flow between deaf and hearing users.

The integration of a chatbot, a sophisticated language model, to enable communication in sign language. Deaf users can express themselves using their native sign language, and the chatbot responds in kind.

Additionally, the platform's adaptability allows hearing users to engage in conversations through text or spoken language, promoting a two-way communication channel that caters to the diverse needs of both communities.

This comprehensive solution is designed not just to address a problem but to transform the way we perceive and experience communication, fostering a world where everyone can connect, communicate, and be heard.

Scope

The scope of this project is to create a Sign Language Translator (SLT) application, enabling bidirectional translation between sign language and written/spoken language. The system will also create a chatbot interface which will be designed to recognize the sign language queries and respond with an avatar using sign language.

The development of the system will be thoroughly documented, starting with an introduction to the stakeholders and user analysis. This will be followed by a description of the system analysis and design process, including the identification of functional and non-functional requirements and the creation of design diagrams. The tools used in the development process will also be outlined, and detailed explanations of the main features of the platform and the models used to implement them will also be provided.

Methodology

For the development of a Sign Language Translator (SLT), an iterative and agile methodology is suitable. Given the dynamic nature of language, the need for constant user feedback, and the potential for evolving requirements, an agile approach allows for flexibility and responsiveness throughout the development process.

- 1- Requirements Gathering: Initial requirements are gathered, but the agile process allows for ongoing refinement and addition of requirements as the project progresses.
- 2- Analysis and Design: These phases are concurrent and iterative. As new requirements emerge or existing ones evolve, the design is adjusted accordingly, allowing for continuous improvement.

- 3- **Implementation:** Incremental development of the SLT occurs, with regular delivery of functional components. This allows stakeholders to see tangible progress and provide feedback early in the process.
- 4- **Testing:** Continuous testing is integrated throughout the development cycle, ensuring that each increment is thoroughly tested. This agile approach promotes early detection and resolution of issues.
- 5- **Feedback and Iteration:** Regular reviews and feedback sessions with users and stakeholders drive iterative improvements. Changes can be made swiftly, ensuring that the SLT aligns closely with user needs.
- 6- **Maintenance:** Ongoing maintenance and enhancements can be seamlessly integrated into the development process, addressing bugs and incorporating user-driven modifications without disrupting the overall project timeline.

Competitors

App	Free	Fast	Multi-language	Avatar	Text to sign	Sign to text	Sentences	Chatbot
Pocket Sign	✓	✓	✗	✗	✓	✗	✗	✗
Sign Translate	✓	✓	✓	✗	✓	✗	✓	✗
Mimix	✓	✓	✗	✓	✓	✗	✓	✗
Prodeaf	✓	✓	✓	✓	✓	✗	✓	✗
Hand Talk	✓	✓	✓	✓	✓	✗	✓	✗
Motionsavy	✗	✓	✗	✗	✓	✓	✓	✗
Signfy (ours)	✓	✓	✗	✓	✓	✓	✓	✓

Table 1: competitors

User Centered Analysis

Functional Requirements

- 1- The customer creates an account by entering his first name, last name, email, username and password.
- 2- The user will be able to translate his sign language to text/voice by entering his sign language gestures as a live video and getting a text which can be heard as a voice.

- 3- The user will be able to translate his text/voice to sign language by entering his text/voice and getting an avatar the translated as sign language
- 4- The user will be able to communicate with a chatbot interface that understands and responds to inquiries made in sign language.

Non Functional Requirements

- 1- Usability: we aim to have a user-friendly system with easy-to-use features that facilitate the communication between Deaf and speaking user and using a chatbot by deaf user to ensure a maximum usability we are carrying on a solid development process from wire framing to the final deliverable.
- 2- Compatibility: our system will be provided as an application.
- 3- Reliability: the system should be reliable enough to do all its functions smoothly and operate in a defined environment. Overall user experience would be consistent and predictable.
- 4- Performance: our system is descent response time despite the highly complex tasks it does.

Project Stakeholder

- 1- Stakeholders are all individuals or groups who have an active stake in the project and can potentially impact, either positively or negatively, its development during its lifecycle.
- 2- **Non-Verbal Individuals:** People who can't speak due to various reasons, such as speech impediments, medical conditions, or other communication challenges. They are direct beneficiaries of the program, relying on it to express themselves effectively using sign language.
- 3- **Verbal Individuals:** Individuals who can speak but may not have proficiency in sign language. They use the program to translate their spoken words into sign language, facilitating communication with non-verbal individuals. This group includes family members, friends, and professionals interacting with the non-verbal community.

- 4- **Developers:** The technical team responsible for developing and maintaining the software are vital stakeholders. Their expertise ensures the program's functionality, security, and continuous improvement to meet the evolving needs of users.

Time plan

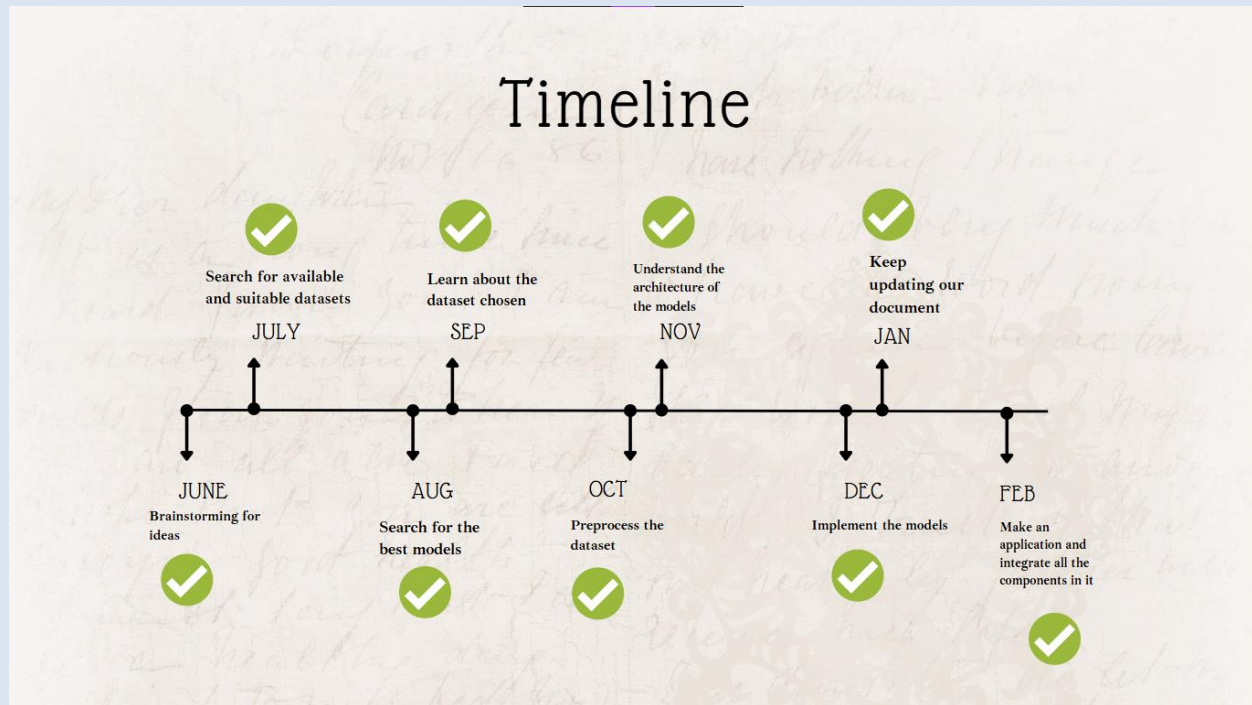


Figure 1: timeline

System Analysis & Design

Use case diagram

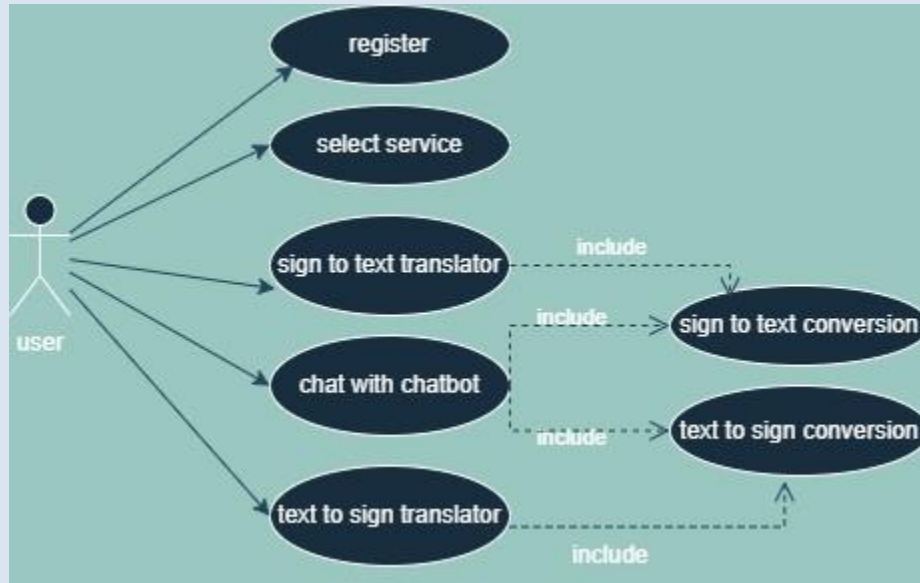


Figure 2: use case diagram

User Stories & Prototype

➤ User Story #1

User Story ID	US #1
User Story Name	Registration / login
Actors	User
Description	As a customer I like to be able to have an account So that my data will be saved in system
Acceptance Criteria	Given I'm logged-out and I'm on the Sign-In page When I fill in the "Username" and "Password" fields with my authentication credentials and I click the Sign-In button in case I already have an account or I register instead. Then the system signs me in.

Normal Scenario

Actor Action	System Response
1-User Enters his Username and Password. 2-Click Sign in	
	3- System Verifies user data 4- System redirects to home page of user
5-Click Sign up	
	6-System displays a form for the user to fill
7- User enters [First name, Last name, Email, Username and Password]	
	8-System verifies the input data 9- System redirects to home page of user

Exceptional Scenario

Actor Action	System Response
1-User Enters Username and Password. 2- Clicks Sign in	
	3-Credentials is invalid 4- System rejects signing in and displays an error message
5-Clicks Sign up	
	6-System displays a form for the user to fill
7- User enters [First name, Last name, Email, Username and Password]	
	8-Form is incomplete 9- System rejects registration and displays an error message with the empty fields or invalid data.

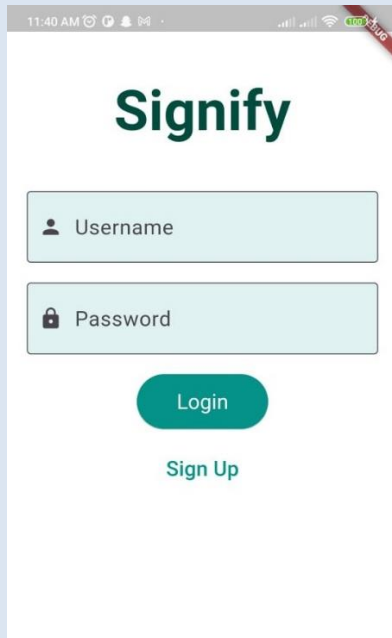


Figure 3: login interface

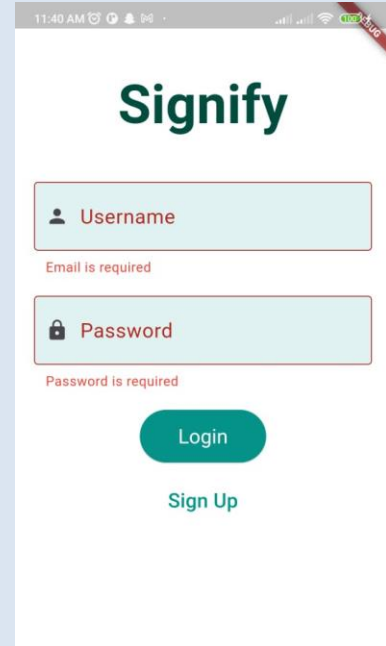


Figure 4: missing field interface

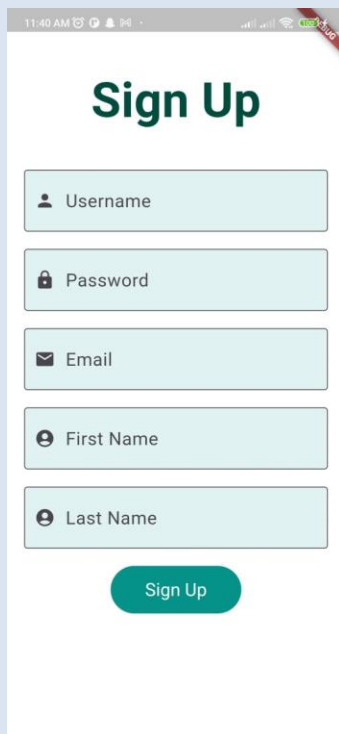


Figure 5: sign up interface

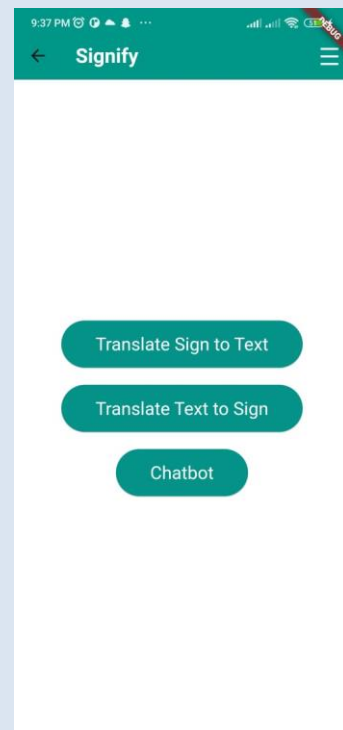


Figure 6: choose your service interface

➤ **User Story #2:**

User Story ID	US #2
User Story Name	Sign Language to text/voice Translation
Actors	User
Description	As a user, I want to utilize the Sign Language Translator (SLT) for accurate recognition from sign language to written/spoken language
Acceptance Criteria	<p>Given I want my signs to be translated</p> <p>When I choose the service of sign to text translation</p> <p>Then the system accurately translates and displays the corresponding sign language to text or voice</p>

Normal Scenario

Actor Action	System Response
1-user registers	
	2-system shows three options for user one for translating sign language to text/voice and one for translating text/voice to sign language and one for chatting with chatbot
3-user chooses the service of translation 4-user inputs his sign language gestures	
	5-system translates the signs into text with an option to hear it as voice
6-user press the voice option	
	7-system plays the voice translation

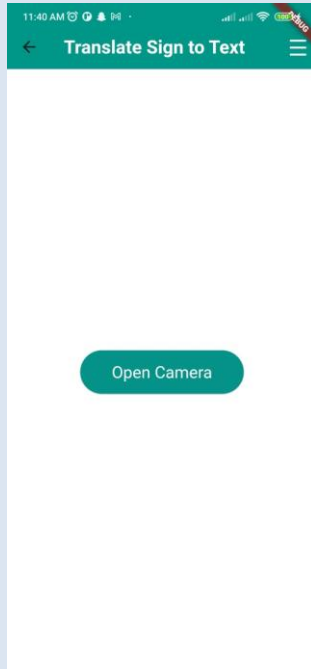


Figure 7: sign to text interface

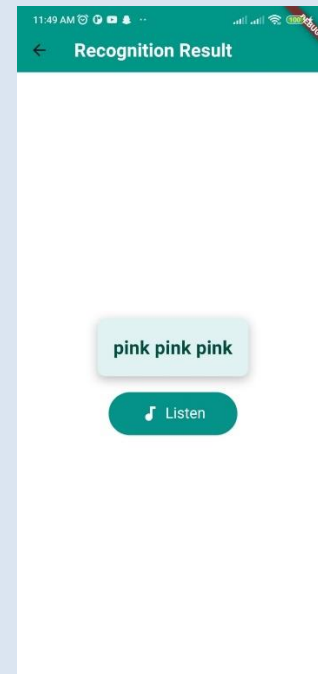


Figure 8: recognition results interface

➤ User Story #3:

User Story ID	US #3
User Story Name	Chatting with Sign Language chatbot
Actors	User
Description	As a user, I want to engage in a chat with the Sign Language chatbot, using sign language for inquiries and receiving responses in sign language.
Acceptance Criteria	<p>Given I want to chat with a sign language chatbot</p> <p>When I choose the chat service with the Sign Language chatbot</p> <p>Then the system provides an interface for inputting sign language gestures and outputs responses for my questions in sign language.</p>

Normal Scenario

Actor Action	System Response
1-user registers	
	2-system shows three options for user one for translating sign language to text and one for translating text to sign language and one for using the chat bot
3-user chooses the service of chatbot 4-user inputs his sign language gestures	
	5-system respond to the user question with an avatar doing sign language

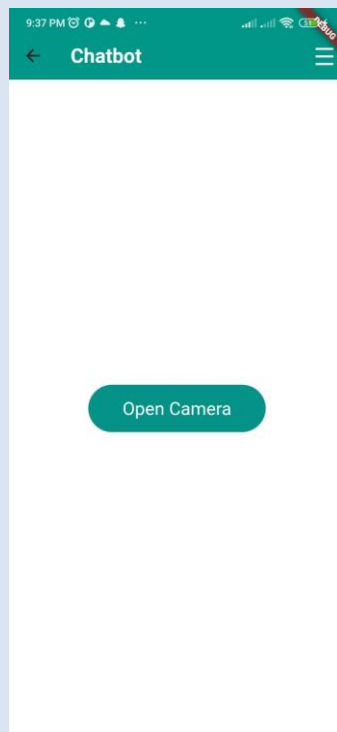


Figure 9: chatbot interface

➤ **User Story #4:**

User Story ID	US #4
User Story Name	text to sign language Translation
Actors	Speaking User
Description	As a speaking user, I want to utilize the Sign Language Translator (SLT) for accurate translation from written to sign language
Acceptance Criteria	<p>Given I want my text to be translated to sign language</p> <p>When I choose the service of translation</p> <p>Then the system accurately translates and displays the corresponding text or voice to sign language</p>

Normal Scenario

Actor Action	System Response
1-user registers	
	2-system shows three options for user one for translating sign language to text and one for translating text to sign language and one for using chatbot
3-user chooses the service of translation text to sign language	
	4-system has an input window
5-user inputs his text	
	6-system translates the text to sign language and display the translation as an avatar

Exceptional Scenario

Actor Action	System Response
1-user registers	

	2-system shows three options for user one for translating sign language to text and one for translating text to sign language and one for using chatbot
3-user chooses the service of translation text to sign language	
	4-system has an input window
5-user inputs his text	
	6-system cannot find the text in the dictionary so it splits the text and display an avatar doing each letter of the text

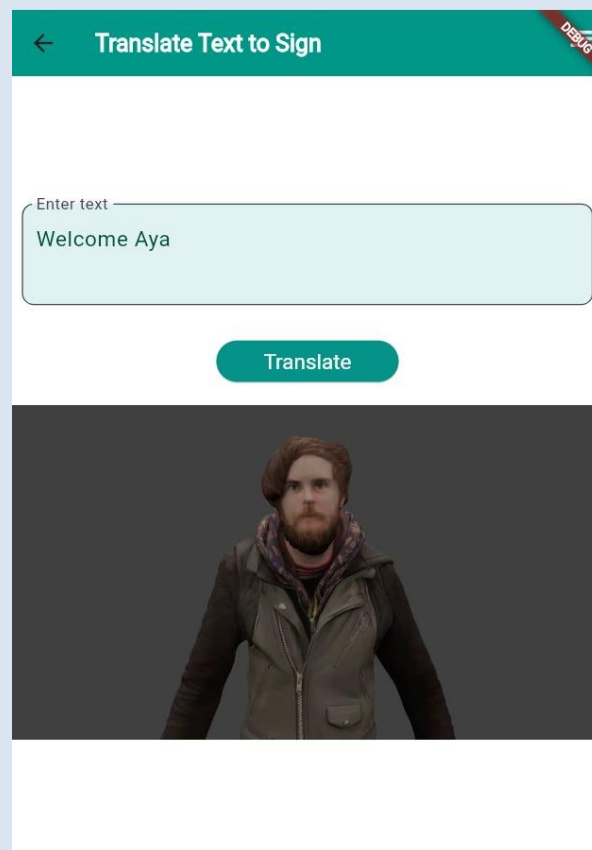


Figure 10: text to sign translation interface

Entity Relationship Diagram (ERD)

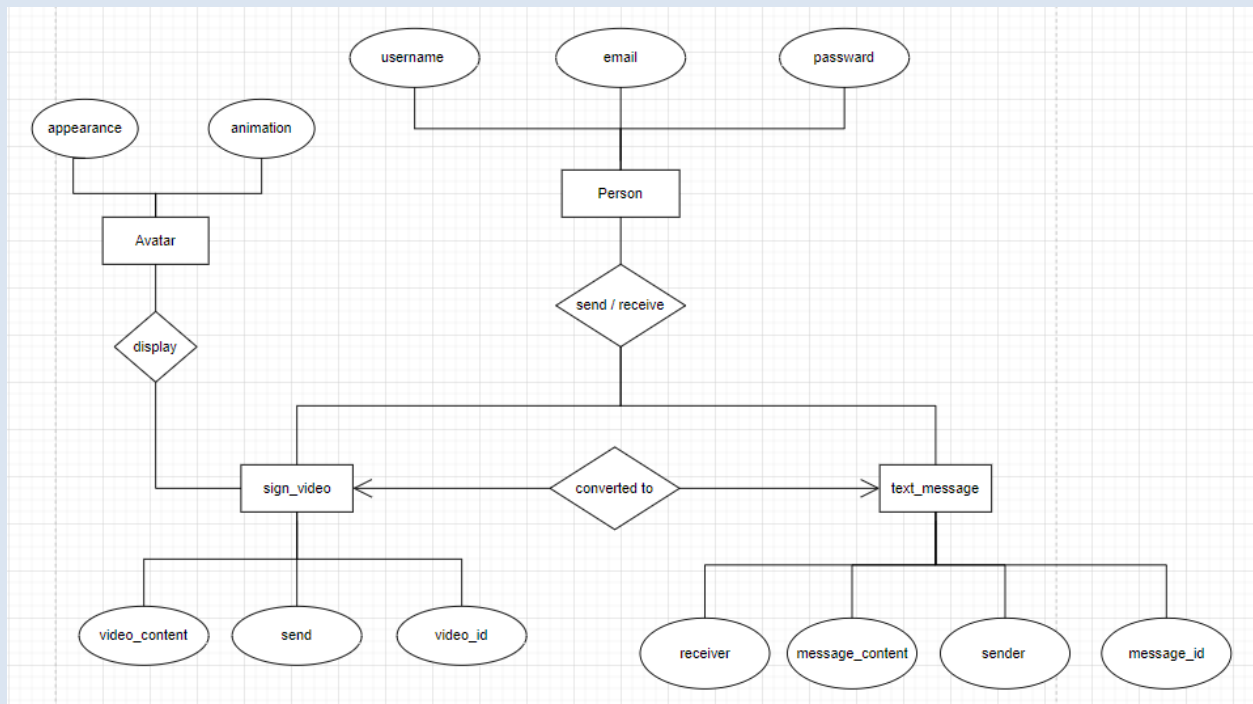


Figure 11: ERD

Class Diagram

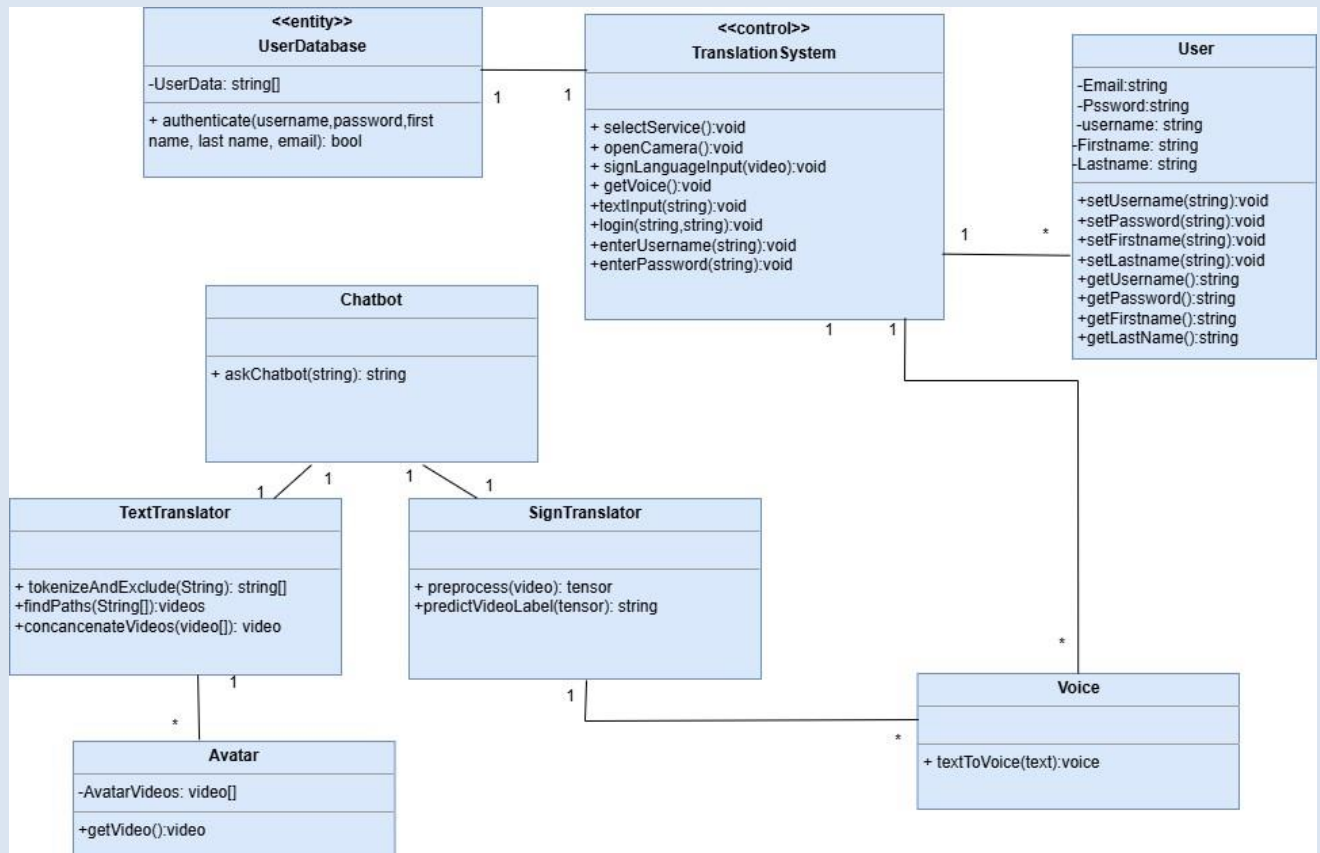


Figure 12: class diagram

Sequence Diagrams

1- Registration

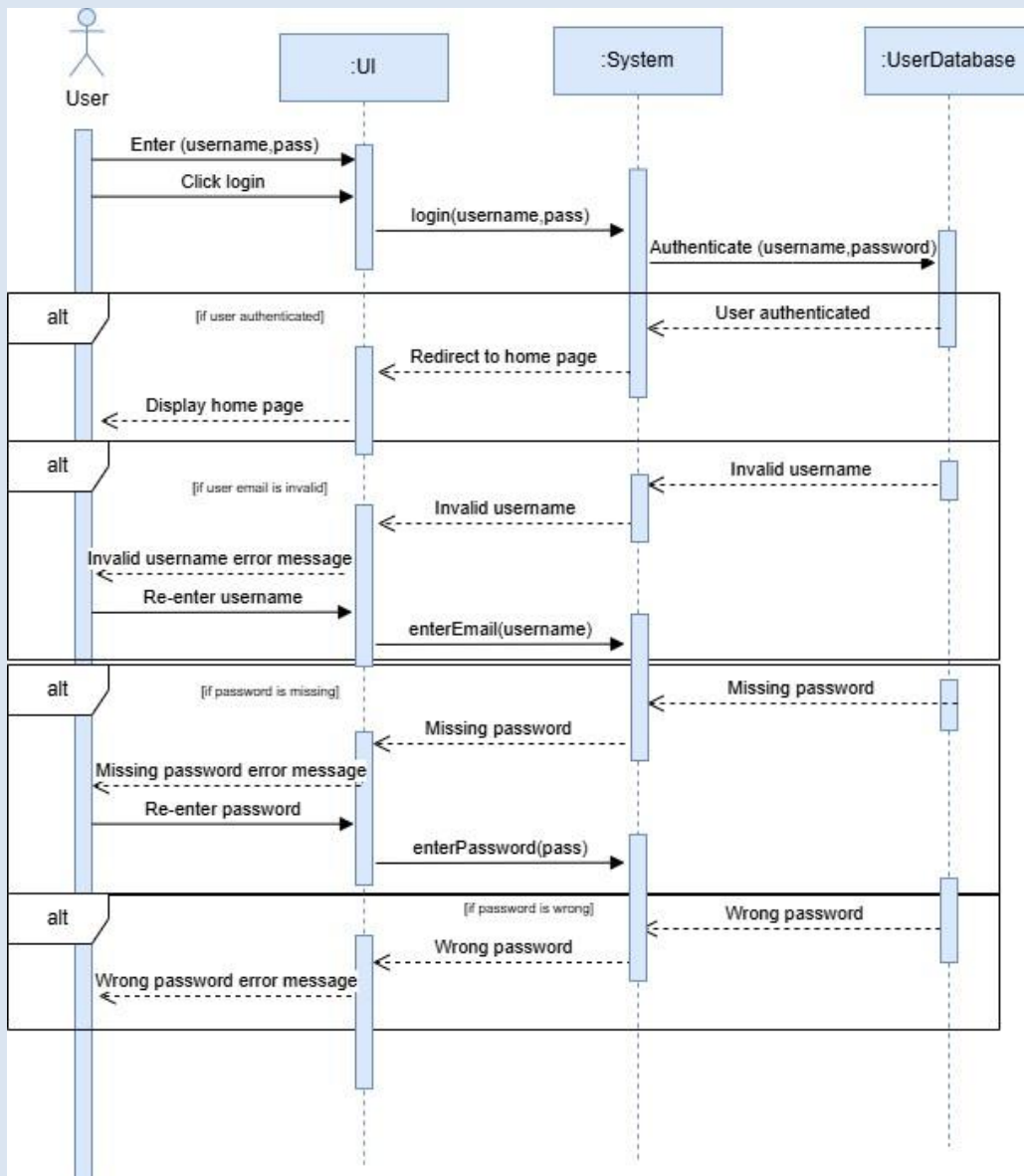


Figure 13: register sequence diagram

2- Sign to text

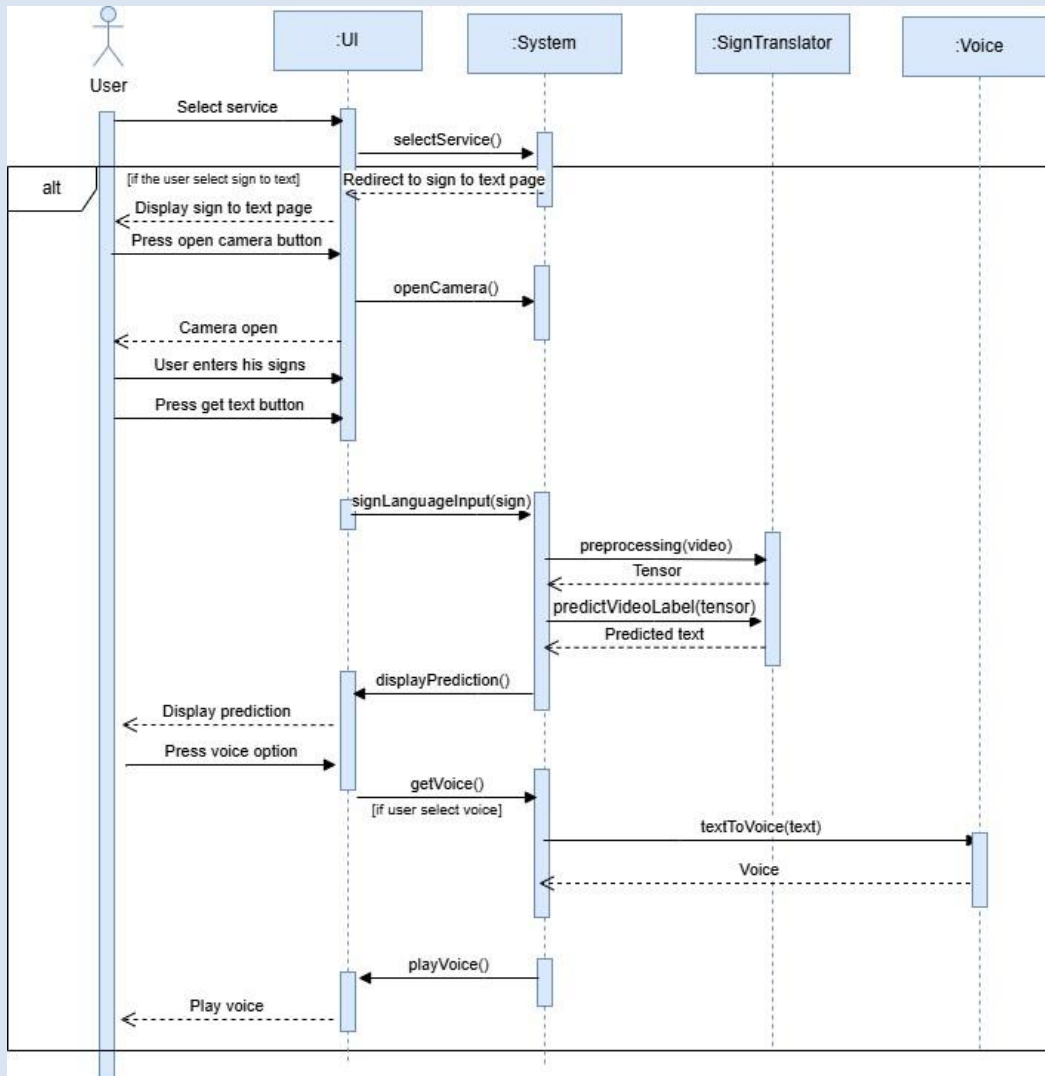


Figure 14: sign to text sequence

3- Text to sign

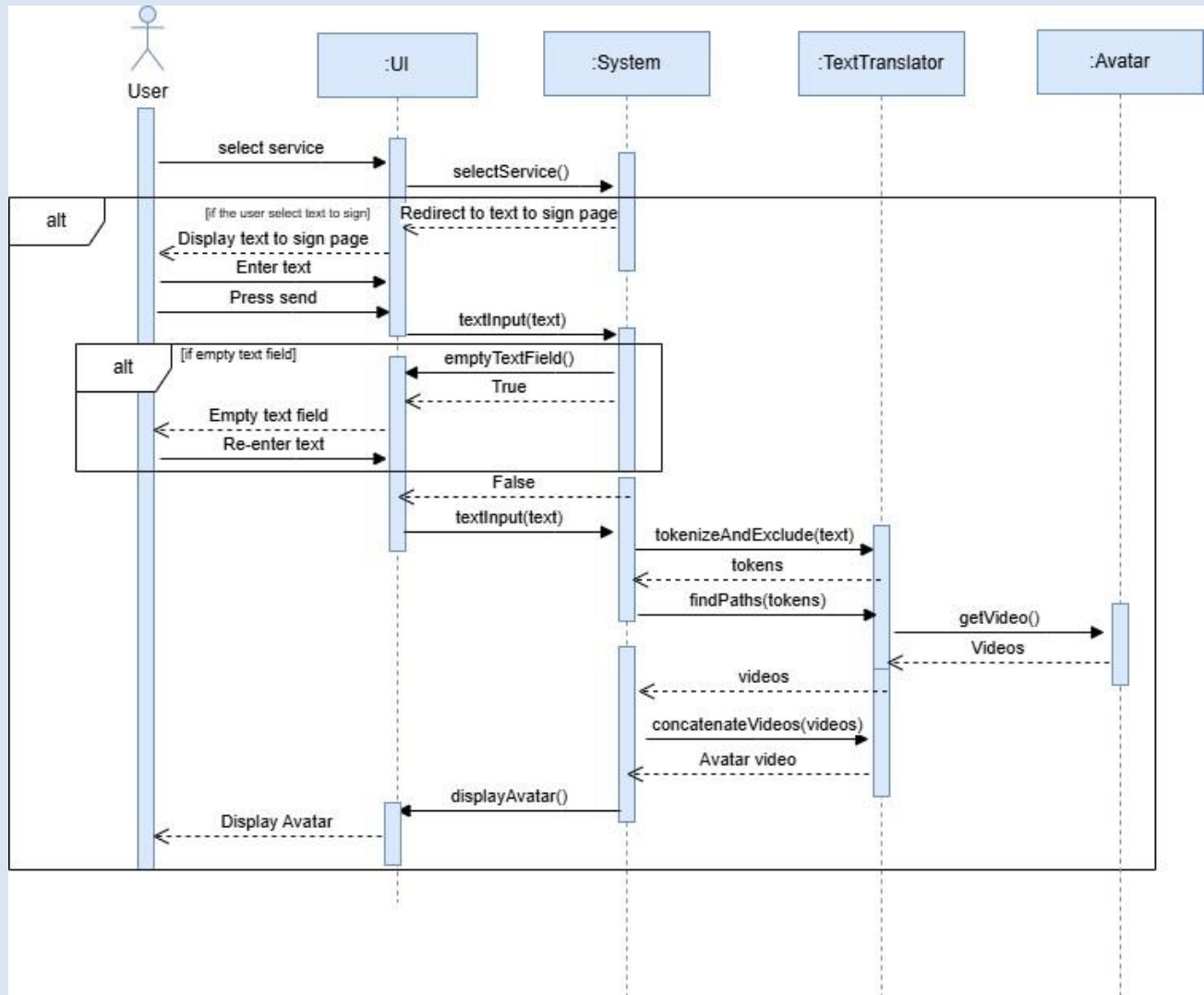


Figure 15: text to sign sequence

4- Sign chatbot

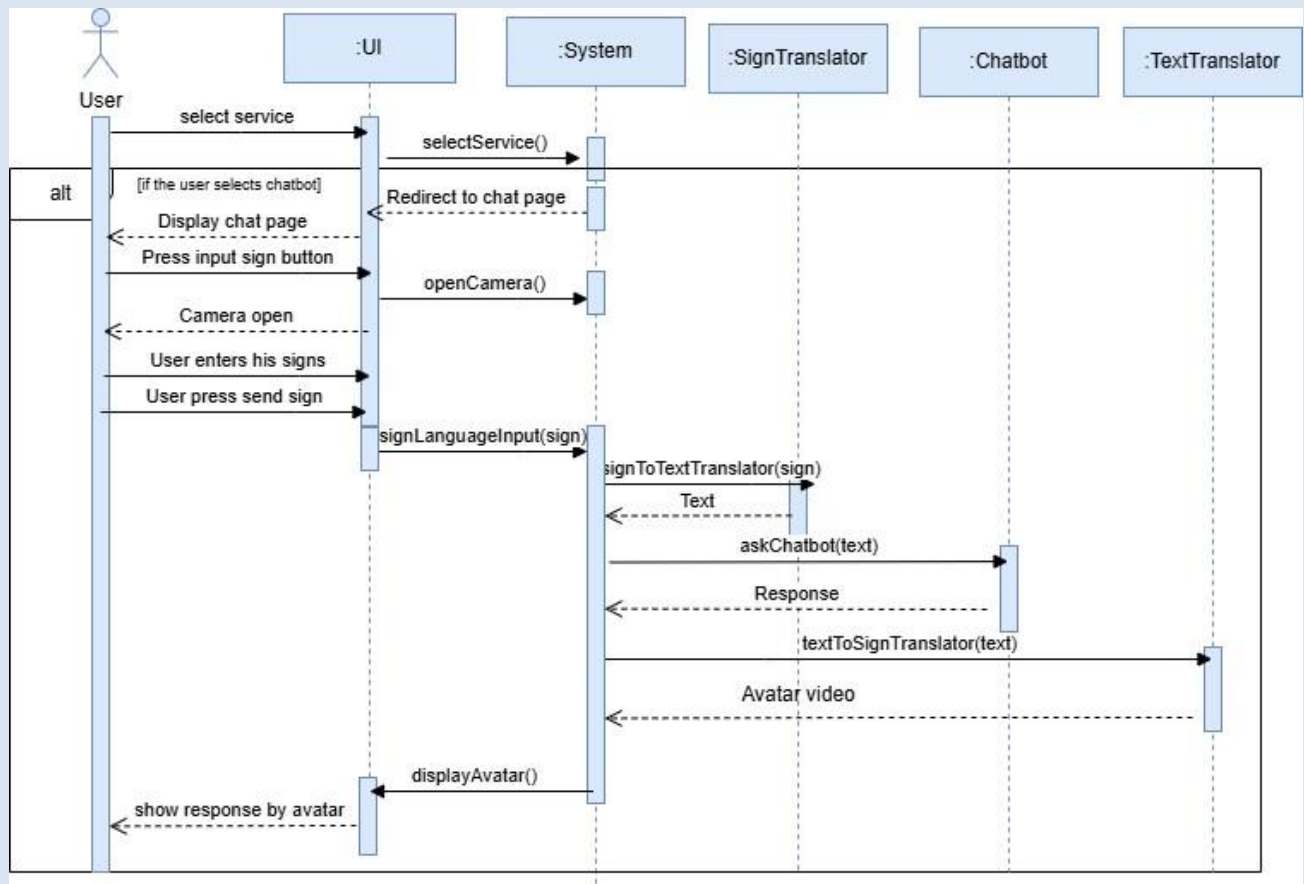


Figure 16: sign chatbot sequence

System Architecture

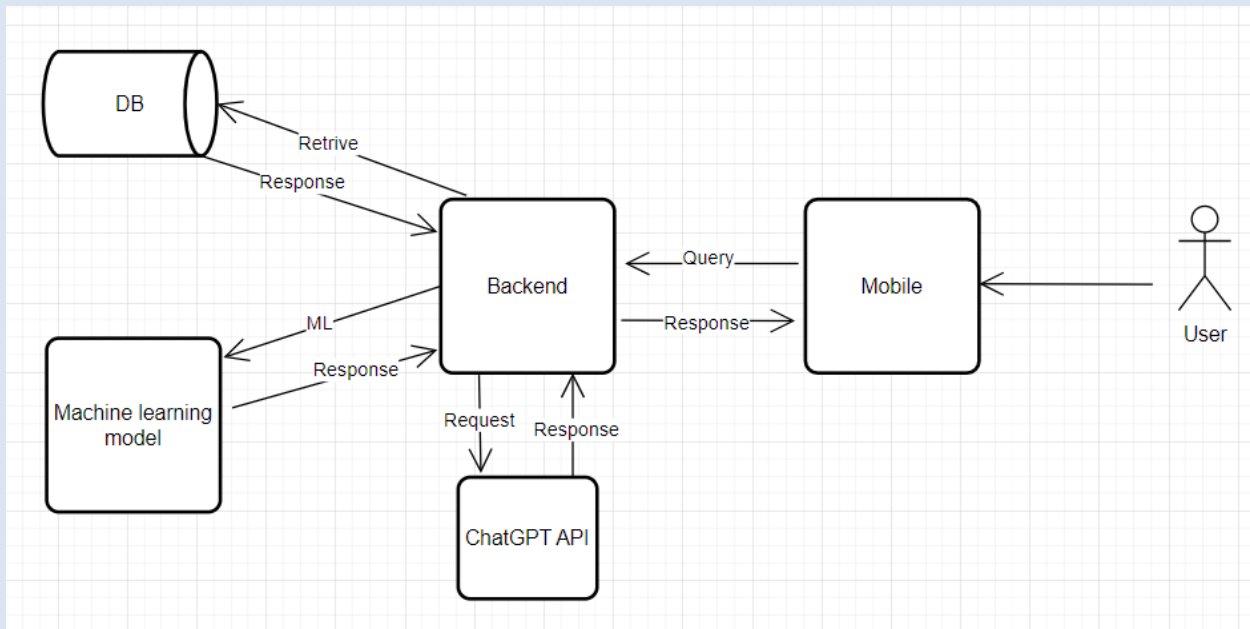


Figure 17: system architecture

Implementation

Dataset

DATASET	# SIGNERS	# VIDEOS	# GLOSSES
WLASL-100	97	2,038	100
WLASL-300	109	5,117	300
WLASL-1000	116	13,168	1,000
WLASL-2000	119	21,083	2,000

Table 2: WLASL dataset description

WLASL is an ASL dataset with a vocabulary of 2,000 words and 21,083 samples. It releases WLASL100 and WLASL300 as its subsets. WLASL is collected from Web videos and bring new challenges due to unconstrained real-life recording conditions.

We utilize the WLASL100 dataset for both our text-to-sign and sign-to-text translators.

Techniques

Text to sign model

A majority of prior works on Spoken2Sign translation focused on expressing translation outcomes through keypoints. However, the keypoint

representations often pose interpretable challenges for signers. With the evolution of generative models, several studies have employed these keypoints to animate sign images, subsequently creating a sign video. However, the output video format is prone to blurriness and visual distortions. In this work, we introduce an innovative method for Spoken2Sign translation by utilizing an avatar to represent the translation results.

Our Spoken2Sign translation pipeline primarily consists of three components: 1) a Text2Gloss translator that translates the input text into a gloss sequence; 2) an avatar producer, which produces the animated sign avatars; and 3) a sign connector, which stitches adjacent signs together.

1- Text2Gloss Translator:

In the text gloss generator, we split the input sentence into words, do some preprocessing to these words such as removing punctuations and correcting grammatical errors, then we check if a word doesn't exist in the dataset we split it into letters, and then for each word (gloss) we get its corresponding video from the dataset.

2- Avatar producer:

In the avatar producer, we utilize the videos to generate the movements for the avatar. This process is facilitated using Blender.

3- Sign connector:

In the sign connector, we take the avatars coming from the avatar producer and stitch them together into one.

Sign to text model

The framework used is NLA-SLR and it mainly consists of three parts: 1) data pre-processing which generates video-keypoint pairs as network inputs; 2) a video-keypoint network (VKNet) which takes video-keypoint pairs of various temporal receptive fields as inputs for vision feature extraction; 3) a head network containing a language-aware label smoothing branch and an inter-modality mixup branch.

1- Data pre-processing

- Input videos V are cropped and represented as a tensor with dimensions $T * H * W * C$ where number of frames $T = 64$, spatial resolution $H = W = 224$, number of channels $C = 3$
- Keypoints in these videos are estimated using HRNet trained on COCO dataset resulting 63 keypoint for each video 42 out of them are for hands.

- Those keypoints are represented in a heatmap with dimensions $H * W * K$ where $H = W = 112$, $K = 63$. Elements on the heatmaps are calculated using a Gaussian function. We get those heatmaps for all videos then we stack them together in K with dimensions $T * H * W * K$.
- Original 64-frame video (V_{64}) paired with its corresponding keypoint heatmap sequence (K_{64}) are cropped into (V_{32}) and (K_{32}) then both pairs are fed into VKNet leading to more informative and discriminative features.

2- VKNet

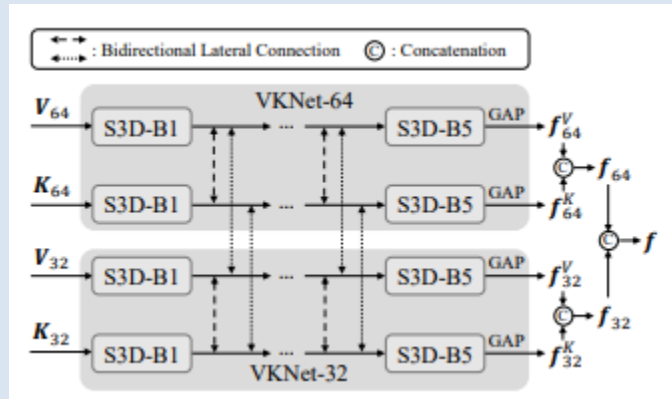


Figure 18: VKNet architecture

- VKNet is composed of two sub-networks, namely VKNet-32 and VKNet-64, which take (V_{32} , K_{32}) and (V_{64} , K_{64}) as inputs, respectively
- Both have the same architecture which consists of a video encoder and a keypoints encoder, since we use heatmaps for keypoints representation, we use a CNN for extracting the feature out of it.
- The CNN used is S3D due to its excellent accuracy-speed trade-off
- The features coming from both VK-32 and VK-64 are concatenated resulting the final features F .

3- Head network

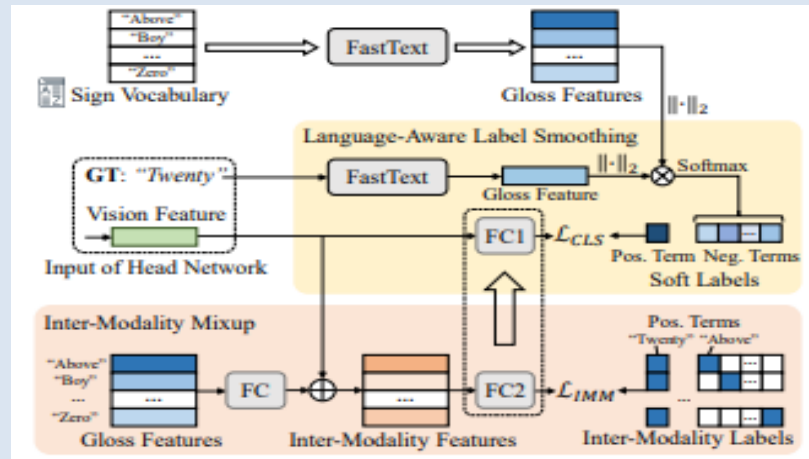


Figure 19: head network architecture

- The head network consists of two parts a language aware smoothing and an intermodality mixup
- The language aware smoothing solves the problem of two signs that look the same and have similar semantic meanings by getting the glosses features and computing similarities between those glosses then include them in the label calculation
- The inter modality mixup solves the problem of two signs that look the same and have different semantic meanings by adding our original features to glosses features to maximize signs separability in the latent space.

4- Implementation

- When it comes to implementation we could not continue with this model due to its need of high computational powers

Sign to text model (NEW MODEL)

1- Previous work

- There are some papers that used rnns and lstms but the problem with them is that they are sequential, you take your image, push through conv layer, get your lstm, then to go to the next step you need to know the output of the previous step
- LSTM (a): Combines convolutional neural networks (ConvNet) with Long Short-Term Memory (LSTM) networks. The ConvNet processes each frame individually, and the LSTM captures the temporal dependencies across frames.

- 3D-ConvNet (b): Uses 3D convolutional neural networks to process spatiotemporal data directly, where the third dimension captures temporal information.
- Two-Stream (c): Combines two separate ConvNets, one processing the spatial information (image frames) and the other processing temporal information (optical flow).
- 3D-Fused Two-Stream (d): Merges the two-stream approach with 3D convolutions, where spatial and temporal streams are fused together
- Two-Stream 3D-ConvNet (e): Extends the two-stream approach by using 3D-ConvNets in both streams, one for images and one for optical flow.

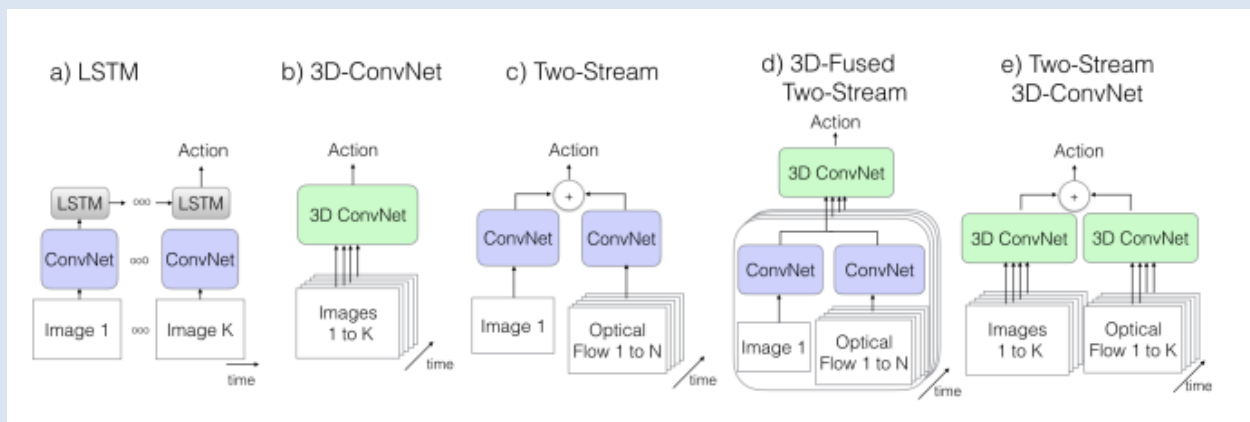


Figure 20: Previous Video architectures. K stands for the total number of frames in a video, whereas N stands for a subset of

2- Preprocessing

- We first load RGB and optical flow frames from each video
- We crop the frames of each video to 224*224 (RandomCrop for train data and CenterCrop for test data)
- Convert the video into tensors of shape (channels*frames*height*width)

3- Model Architecture

- Inflation from 2D to 3D: we took existing successful pretrained 2D imageNet model and adapted its weights to handle videos or sequences (3D data) without redesigning from scratch. We start by taking the pretrained 2D filters, copying them and inflating them by making them cubic $N*N*N$ rather than square then dividing each one by N . Having those weights gives us smart initial weights to give our model to save time compared to creating new models from scratch.

- Two-Stream Configuration: we used two separate data streams:
 RGB Stream: This stream processes video frames directly to capture spatial features (like objects and textures).
 Optical-Flow Stream: This stream captures temporal patterns (movement and changes over time between frames) to enhance performance.

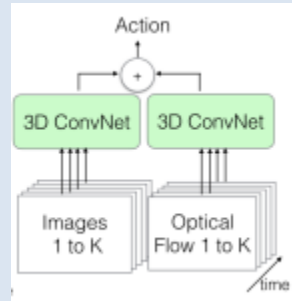


Figure 21: Two stream 3D ConvNet

- The model architecture: when we added our layers we considered the idea of receptive fields (the area each pixel look at to get its information from) this area affect the model ability to detect motion so adding the layers the first convolutional layer and two max pooling layers we do not perform temporal pooling. Then we made receptive field grow over time by using kernels and strides along the temporal dimension.
 Why to start with a low receptive fields then make them higher?
 Because we want to get detailed motion information but we don't want to overfit them.
- When it comes to other layers we start with a convolution layer to extract low level features, then we proceed with other layers including inception modules to extract various features, ending with average pooling to summarize for classification and conv to get the number of classes and finally we predict.
- We use Relu for activation, Softmax for prediction and BatchNorm layers to prevent overfitting

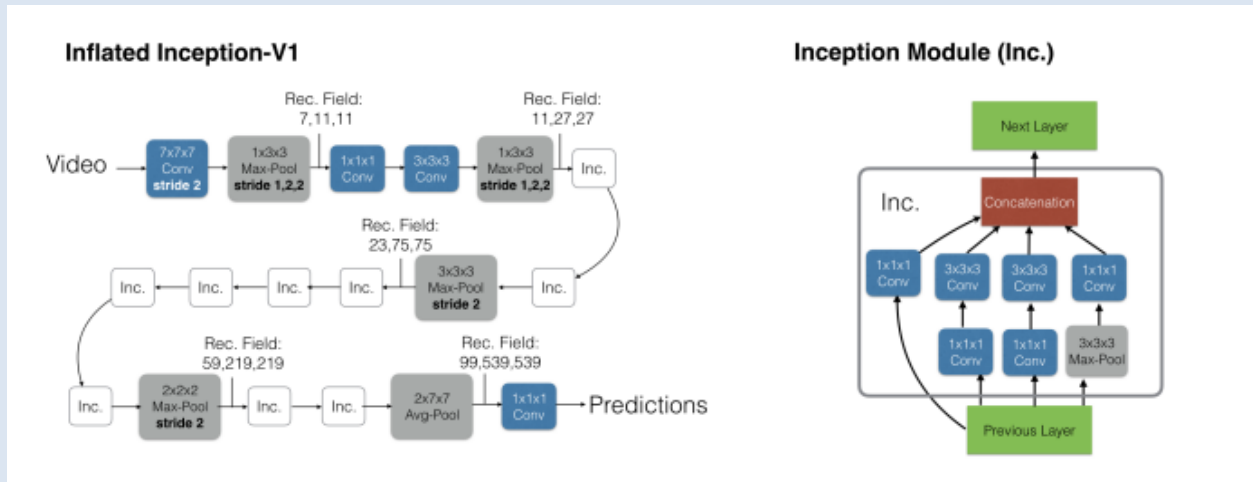


Figure 22: the Inflated Inception-V1 architecture (left) and its detailed inception submodule (right)

4- Implementation

- After training the model on WLASL100 with 896 iteration we got top-1 accuracy 65.89%, top-5 accuracy 84.11%, and top-10 accuracy 89.92% on test data.

Method	WLASL2000				WLASL1000				WLASL300				WLASL100			
	P-I		P-C		P-I		P-C		P-I		P-C		P-I		P-C	
	T1	T5	T1	T5	T1	T5	T1	T5	T1	T5	T1	T5	T1	T5	T1	T5
Pose-GRU									33	64			46	76		
Pose-TGCN	23	51							38	67			55	78		
GCN-BERT									42	71			60	83		
OpenHands	30															
PSLR									42	71			60	83		
I3D	32	57			47	76			56	79			65	84		

Table 3: comparison with previous works on WLASL

Chatbot

- 1- We employ sign-to-text model to take the input sign and translate it into text.
- 2- We send this text to chatbot to generate a response.
- 3- We employ text-to-sign model to translate the text coming from chatbot into sign language represented by an avatar.

References

- [Hand Talk App](#)
- [ASL Sign Language - start learning ASL now with our video lessons \(pocketsign.org\)](#)
- [Sign Translate](#)
- [Mimix Speech to sign language website](#)
- [ProDeaf Translator - Microsoft Apps](#)
- [Hand Talk: your website accessible in ASL](#)
- [MotionSavvy UNI: 1st sign language to voice system | Indiegogo](#)
- [WACV 2020 Open Access Repository \(thecvf.com\)](#)
- [\[2401.04730\] A Simple Baseline for Spoken Language to Sign Language Translation with 3D Avatars \(arxiv.org\)](#)
- [CVPR 2023 Open Access Repository \(thecvf.com\)](#)
- [WACV 2020 Open Access Repository \(thecvf.com\)](#)
- [Pose-Based Sign Language Recognition Using GCN and BERT \(thecvf.com\)](#)
- [\[2110.05877\] OpenHands: Making Sign Language Recognition Accessible with Pose-based Pretrained Models across Languages \(arxiv.org\)](#)
- [WACV 2021 Open Access Repository \(thecvf.com\)](#)
- [CVPR 2017 Open Access Repository \(thecvf.com\)](#)
- [Spatial Temporal Graph Convolutional Networks for Skeleton-Based Action Recognition | Proceedings of the AAAI Conference on Artificial Intelligence](#)
- [WACV 2021 Open Access Repository \(thecvf.com\)](#)
- [BSL-1K: Scaling Up Co-articulated Sign Language Recognition Using Mouthing Cues | SpringerLink](#)
- [Hand-Model-Aware Sign Language Recognition | Proceedings of the AAAI Conference on Artificial Intelligence](#)
- [CVPR 2020 Open Access Repository \(thecvf.com\)](#)
- [\[2302.05075\] BEST: BERT Pre-Training for Sign Language Recognition with Coupling Tokenization \(arxiv.org\)](#)
- [ICCV 2021 Open Access Repository \(thecvf.com\)](#)
- [CVPR 2021 Open Access Repository \(thecvf.com\)](#)
- [\[2110.06161\] Sign Language Recognition via Skeleton-Aware Multi-Model Ensemble \(arxiv.org\)](#)
- [https://arxiv.org/abs/1705.07750](#)