# Sign-Sync

## Software Requirements Specification Document (Demo 2)





| Member | Student Number |
|---|---|
| Michael Stone | u21497682 |
| Matthew Gravette | u23545977 |
| Wessel Johannes van der Walt | u22790919 |
| Jamean Groenewald | u23524121 |
| Stefan Muller | u22498622 |

# Introduction

Sign Sync is a real-time translation system designed to bridge communication between spoken English and American Sign Language (ASL). The system allows users to input spoken or typed English and receive an accurate visual translation in ASL, using a combination of natural language processing, speech recognition, and gesture playback. It is developed as part of the COS301 Capstone Project to assist the Deaf and Hard-of-Hearing community.

# User Characteristics

**Deaf or Hard-of-Hearing users**: Require real-time sign language output from spoken or written English.

**Hearing individuals**: May use the system to learn or communicate via ASL.

**Interpreters**: Professionals who need assistance in translating speech quickly.

**Educators or Students**: Learning environments using ASL as a teaching aid.

**System Administrators**: Maintain the deployed system and monitor performance.

# User Stories

## Deaf or Hard-of-Hearing Users

- View spoken language translated into sign animations.
- Use sign language to communicate back via webcam input.
- Adjust avatar display or gloss format for clarity.
- Access help or FAQs without requiring external assistance.
- Enable accessibility features (e.g., high-contrast mode, larger text, slowed animation).

## Hearing Users

- Speak naturally and have their speech translated into signs.
- Read or hear signed responses translated to text or audio.
- Type messages and see the corresponding ASL gloss and animation.
- Use the system as a teaching or demonstration tool for learning ASL.
- Select different visual styles or avatars for sign output.

## Administrators / Researchers

- Monitor system performance and translation accuracy.
- Gather flagged feedback data for AI retraining.
- Manage user access and customization settings.
- Review usage logs and system health metrics.
- Push updated models or gloss rules into the pipeline.

## Educators / Interpreters

- Demonstrate real-time speech-to-sign translation in classrooms or training.
- Share a live session with multiple learners.
- Export gloss translations for lesson plans.
- Annotate signs or glosses for specific learning contexts.

# Functional Requirements

## R1: Text-to-Sign Translator
**R1.1**: Capture text input from user.
**R1.2**: Translate English text to Sign gloss.
**R1.3**: Search word definition for appropriate sign.
**R1.4**: Display sign through avatar.
**R1.5:** Handle grammar transformations for accurate ASL structure

## R2: Sign-to-Text Translator
**R2.1**: Capture webcam input and extract hand keypoints.
**R2.2**: Classify sign gesture sequences using trained AI models.
**R2.3**: Convert recognized signs to sign gloss.
**R2.4**: Convert sign gloss to English.
**R2.5:** Show real-time visual feedback of detected signs.

## R3: Feedback and AI Improvement System
**R3.1**: Allow users to flag inaccurate translations.
**R3.2**: Log flagged data.
**R3.3**: Retrain AI models periodically using collected data.
**R3.4:** Provide an admin dashboard to review flagged entries and retraining logs

## R4: User Interface
**R4.1**: Display live avatar animations based on translation output.
**R4.2**: Show translated text and/or play voice feedback.
**R4.3**: Offer accessibility options (high-contrast mode, font scaling, voice personas).
**R4.4:** Provide a help menu with tutorials and FAQs.
**R4.5:** Support multi-device layouts (mobile/tablet/desktop).

## R5: User Management and Settings
**R5.1**: Authenticate users.
**R5.2**: Store and retrieve user preferences and settings.
**R5.3**: Allow users to login.
**R5.4**: Allow users to register.
**R5.5:** Allow role-based access for admins, researchers, and general users.

## R6: Speech to Text

**R6.1**: Capture user speech.

**R6.2**: Convert speech to text.

**R6.3**: Display text on screen.

**R6.4:** Send text into ASL gloss converter pipeline.
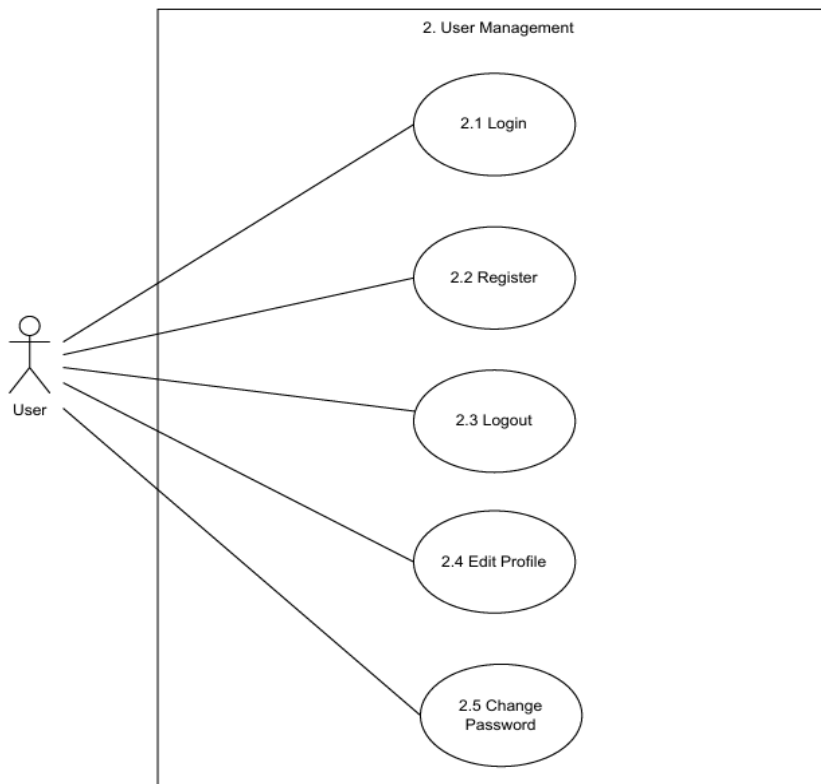

## R7: Text to Speech

**R7.1**: Capture text input by user.

**R7.2**: Convert text to speech.

**R7.3**: Play speech for the user to hear.

**R7.4:** Support voice personalization settings.

# Use Case Diagrams



1.Translator

- 1.1 Convert Sign to text
- 1.2 Convert Text to sign
- 1.3 Convert Speech to text
- 1.4 Convert Text to speech

User



2. User Management

- 2.1 Login
- 2.2 Register
- 2.3 Logout
- 2.4 Edit Profile
- 2.5 Change Password

User

## 3. Feedback system

**User**

3.1 Flag inaccurate translations

## 4. Customisability

**User**

4.1 Manage Preferences

4.2 Customise Avatar

# Domain Model

## User
- userID
- username
- email
- password

## UserPreferences
- userID
- themeMode
- fontSize
- avatarStyle

User 1 — saves > — 1 UserPreferences

User 1 < makes a 0..* SignGesture

User 1 < records 0..* SpeechInput

User 1 types > 0..* TextInput

## SignGesture
- gestureID
- userID
- keypoints
- timestamp

## SpeechInput
- speechID
- userID
- audioData
- transcribedText
- timestamp

## TextInput
- textID
- userID
- inputText
- timestamp

SpeechInput 1 — converted to > — 1 TextInput

SignGesture 1 < classified into 1 SignPrediction

TextInput 1 < translated to 1 GlossTranslation

## SignPrediction
- predictionID
- gestureID
- predictedWord
- confidence

## SignAnimation
- signID
- animation
- keywords[]

## GlossTranslation
- glossID
- sourceText
- glossText

SignAnimation 0..* < matched to 1 GlossTranslation

# Service Contracts (Graded in demo 3)

## Alphabet Translator Service Contract

**CameraFrontend**
+ captureKeypoints(): array
+ requestPrediction(keypoints: array): Text

{performance: 500ms}
{post: prediction delivered}
{post: no user data is kept}
{pre: recipient connected}
{scalability: can sign anything to camera}
{error: return 'No Hand Found' if no keypoints}

**translator_api**
+ predict(request: KeyPointsRequest)
+ normalization(Keypoints: array): array

**model**
+ predict(input: array): string

## Speech-to-Text Service Contract

**API**
+ speech_to_text( : WebSocket) : void
+ get_model() : Model

**KaldiRecognizer** | Import
+ AcceptWaveform( : bytes) : void

{performance: translates an 18 second audio file to text in 3.46 seconds, providing a Real-Time Factor of 0.19 < 1}
{post: text representation of the audio file successfully displayed to the requesting party}
{post: if the audio was in the incorrect format an error message will be sent to the requesting party}
{pre: websocket is open and client has performed CORS handshake}
{pre: server is running and accepting requests}

## Text-to-Sign Service Contract

**TextToSign**

+ showSign( : ) : void

+ processSentence( : string) : void

**TranslatorPage**

{pre: text length > 0}
{pre: server is authenticated and accepting requests}
{post: sign images corresponding to characters displayed}
{post: user able to cycle through images}
{performance: completion in <= 2ms consistently}

# Architectural Requirements

📄 Sign-Sync: Architectural Requirements Document

# Technology Choices

## Frontend Framework

| Framework | Pros | Cons |
|---|---|---|
| Angular | <ul><li>Full-featured MVC framework</li><li>Large enterprise support</li></ul> | <ul><li>Steep learning curve</li><li>Heavy bundle size</li></ul> |
| React | <ul><li>Component-based</li><li>Huge ecosystem and community</li><li>Easy Websocket integration</li></ul> | <ul><li>State management can be difficult</li><li>Setup can be tedious</li></ul> |
| Svelte | <ul><li>Compiles to vanilla JS</li><li>Fast performance</li></ul> | <ul><li>Less enterprise adoption</li><li>Smaller ecosystem</li></ul> |

**Choice:**
React was selected due to its modular structure, vibrant and large ecosystem and ease of integrating real-time features such as websockets. This aligns well with the microservices architecture and enables a maintainable, scalable frontend.

## Backend Language

| Language | Pros | Cons |
|---|---|---|
| Python | <ul><li>Large AI/ML ecosystem</li><li>Simple, readable syntax</li><li>Strong library support</li></ul> | <ul><li>Slower runtime</li><li>Not ideal for multi-threading</li></ul> |
| JavaScript (Node.js) | <ul><li>Full-stack JS</li><li>Large NPM ecosystem</li></ul> | <ul><li>Difficulty in debugging</li><li>Complex async handling</li></ul> |
| Go | <ul><li>Excellent concurrency</li><li>Fast Performance</li></ul> | <ul><li>Limited AI/ML libraries</li></ul> |

**Choice:**

Python was chosen for backend services, especially AI-related modules, due to its excellent support for ML and NLP libraries, such as spaCy and Vosk. While it is not the fastest, its developer productivity and expressiveness make it ideal for rapidly developing and deploying independent services. This aligns perfectly with the microservices architecture.

## API Framework

| Framework | Pros | Cons |
|-----------|------|------|
| ExpressJS | <ul><li>Minimal and flexible</li><li>Well-established</li><li>Fast setup</li></ul> | <ul><li>Requires manual validation</li><li>Not type-safe</li></ul> |
| FastAPI | <ul><li>Fast, async support</li><li>Easy validation with Pydantic</li><li>Auto-generated docs</li></ul> | <ul><li>Lacks some mature integrations</li><li>Still relatively new</li></ul> |
| Flask | <ul><li>Lightweight</li><li>Simple for quicks APIs</li><li>Mature and stable</li></ul> | <ul><li>Not async by default</li><li>Less scalable for real-time</li></ul> |

**Choice:**
FastAPI was chosen as our API framework since it supports our microservice architecture with its async design, fast performance and modular structure. Each microservice can be independently built and deployed using this framework which ensures scalability and maintainability.

## Database

| DB | Pros | Cons |
|---|---|---|
| MongoDB | <ul><li>NoSQL, flexible schema</li><li>Document-oriented (therefore great for JSON data)</li></ul> | <ul><li>Less suitable for relational data</li><li>Data consistency is not always guaranteed</li></ul> |
| PostgreSQL | <ul><li>Strong ACID compliance</li><li>Complex querying</li></ul> | <ul><li>Requires fixed schema</li><li>Slightly more setup for scaling</li></ul> |
| Firebase Realtime DB | <ul><li>Real Time sync</li><li>Easy to use</li><li>Scales well</li></ul> | <ul><li>Less control over backend logic</li><li>No relational structure</li></ul> |

**Choice:**
MongoDB was chosen due to its document-oriented structure which fits well with storing user preferences and data and loosely structured data. It also complements a microservices setup by being easy to scale independently per service.

## Speech Recognition

| Model | Pros | Cons |
| --- | --- | --- |
| Mozilla DeepSpeech | <ul><li>Open source</li><li>Good accuracy</li><li>Active community</li></ul> | <ul><li>Large models</li><li>High resource usage</li></ul> |
| Vosk | <ul><li>Free</li><li>Fast and multilingual</li><li>Real-time</li><li>Raw byte streams</li></ul> | <ul><li>Limited documentation</li><li>Smaller community</li></ul> |
| Google Speech API | <ul><li>Very high accuracy</li><li>Robust language support</li></ul> | <ul><li>Cloud-only</li><li>Latency</li><li>Usage cost</li></ul> |

**Choice:**
We chose Vosk because it runs offline, supports real-time transcription and integrates easily into independent microservices without relying on external APIs. This is crucial for maintaining modularity and reducing latency in a distributed architecture.

## NLP Processing

| Model | Pros | Cons |
| --- | --- | --- |
| spaCy | <ul><li>Lightweight</li><li>Pretrained models</li><li>Easy to integrate</li></ul> | <ul><li>Limited deep semantic analysis</li></ul> |
| NLTK | <ul><li>Rich library for NLP education/research</li></ul> | <ul><li>Slower</li><li>Outdated for production systems</li></ul> |
| HuggingFace Transformers | <ul><li>State-of-the-art models</li><li>flexible</li></ul> | <ul><li>Heavier</li><li>Complex integration</li></ul> |

**Choice:**

spaCy was chosen for its speed and simplicity which is ideal for real-time language processing within our NLP microservice. Its modularity ensures each NLP-related function can scale and update independently in the overall architecture.

## Gesture Recognition

| Model | Pros | Cons |
|---|---|---|
| TensorFlow (TCN) | <ul><li>Great for temporal sequences</li><li>Memory efficient</li></ul> | <ul><li>Steeper learning curve</li><li>Requires model tuning</li></ul> |
| PyTorch (LSTM) | <ul><li>Dynamic graph</li><li>Easy debugging</li></ul> | <ul><li>Slower in production</li><li>Less optimised for mobile</li></ul> |
| MediaPipe | <ul><li>Fast</li><li>Easy gesture pipelines</li></ul> | <ul><li>Limited customisation</li><li>Black-box components</li></ul> |

**Choice:**

TensorFlow with Temporal Convolutional Networks (TCNs) was chosen due to their strong performance in recognising sequences, such as gestures. These models are containerised and deployed as an isolated microservice which aligns well with our architecture's need for scalable, efficient model inference.

## Hand Recognition

| Model | Pros | Cons |
|---|---|---|
| OpenCV | <ul><li>Lightweight</li><li>Cross-platform</li><li>Integrates well with Python</li></ul> | <ul><li>Requires manual tuning</li><li>No built-in hand detection</li></ul> |
| MediaPipe | <ul><li>Fast</li><li>Pretrained hand landmark detection</li></ul> | <ul><li>Harder to customise</li><li>Black-box components</li></ul> |
| OpenPose | <ul><li>Highly accurate for full body/hands</li></ul> | <ul><li>Heavy</li><li>GPU-dependent</li><li>Harder to deploy at scale</li></ul> |

**Choice:**
OpenCV and MediaPipe as they are easy to integrate. They are flexible and lightweight which makes it ideal for our hand recognition microservice. It enables fine-tuned control and, when containerised, it integrates smoothly into the microservices environment without excessive resource demands.

## Hosting

| Service | Pros | Cons |
|---|---|---|
| Amazon Web Services (AWS) | <ul><li>Highly scalable and battle-tested</li><li>Offers free-tier services (EC2, S3, Lambda) suitable for MVP deployments</li><li>Excellent integration with Docker, API Gateways, and CI/CD tools</li></ul> | <ul><li>Complex initial setup</li><li>Steeper learning curve for new developers</li><li>Cost increases quickly beyond the free tier</li></ul> |
| Google Cloud Platform (GCP) | <ul><li>Excellent for containerized deployments (e.g., Cloud Run, GKE)</li><li>Great NLP/AI service integrations if needed in future</li><li>Free-tier credits for students and education teams</li></ul> | <ul><li>Fewer community resources/tutorials compared to AWS</li><li>Region-specific performance may vary</li></ul> |
| Microsoft Azure | <ul><li>Strong enterprise integrations and CI/CD via GitHub Actions</li><li>Azure App Service is simple for deploying Python + React apps</li><li>Offers educational credits for students</li></ul> | <ul><li>Documentation is sometimes inconsistent</li><li>Slightly more expensive for persistent container hosting than GCP</li></ul> |

**Choice:**
Still being discussed with the client.

# Deployment Plan/Model (Demo 3)

The Sign Sync system is designed to run as a containerized, microservice-based application, which allows us to deploy each translation component independently. The deployment model is cloud-ready and optimized for scalability, portability, and modularity.

- **Frontend** (React) — served via Azure App Service or S3 + CloudFront (AWS)

- **Backend Microservices** (FastAPI) — containerized and deployed using Docker on Microsoft Azure or AWS

- **Database** (MongoDB) — hosted via MongoDB Atlas or Azure Cosmos DB with Mongo API

- **WebSocket-enabled Speech Service** — exposed through secure public endpoints

- **CI/CD Pipeline** — automated builds and deployments via GitHub Actions, with deployment scripts targeting Azure or AWS services

# Live Deployed System (Demo 4)

Work towards in demo 3