

Sign Language to Text Conversion Using LSTM Model

Shraddha Pandey¹, Narayan Pandey², Sahil Patel³, Vedant Khangan⁴, Prof.Sonal .R. Jathe⁵

^{1,2,3,4} Final Year Student, B.E, Department of C.S.E, Jawaharlal Darda Institute of Engineering and Technology, Yavatmal, India

⁵ Assistant professor, Department of C.S.E, Jawaharlal Darda Institute of Engineering and Technology, Yavatmal, India

Abstract - Sign language is a crucial medium of communication for individuals with hearing impairments. However, the lack of efficient tools for converting sign language into text limits the accessibility of information for this community. In this paper, we propose a novel approach for sign language to text conversion using Long Short-Term Memory (LSTM) neural networks. LSTM networks are well-suited for sequential data processing, making them an ideal candidate for capturing the temporal dependencies inherent in sign language gestures. We present the architecture of our LSTM model, discuss the datasets preparation, training methodology, and evaluate the performance of our model through comprehensive experiments. The results demonstrate the effectiveness and feasibility of our proposed approach in accurately translating sign language gestures into text, thereby enhancing accessibility for individuals with hearing impairments.

Keywords: sign language, LSTM, neural networks, accessibility, hearing impairments

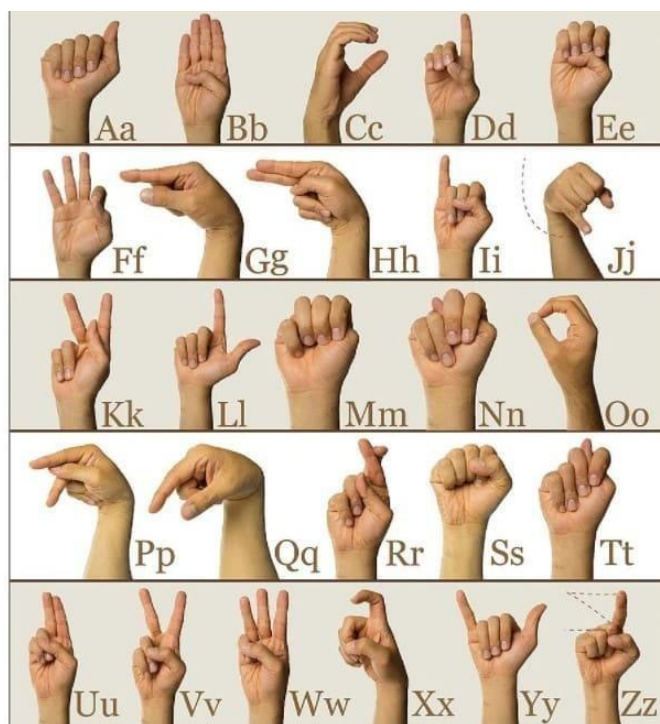
1. INTRODUCTION

Sign language serves as a primary mode of communication for individuals with hearing impairments, allowing them to express themselves and interact with others effectively. However, the lack of efficient tools for converting sign language into text poses significant challenges for accessing information and communication for this community. Traditional methods for sign language recognition often rely on hand-crafted features and rule-based systems, which may not generalize well to diverse sign languages and gestures.

Recent advancements in deep learning, particularly in the field of sequence modeling, have shown promising results in various natural language processing tasks. Long Short-Term Memory (LSTM) networks, a type of recurrent neural network (RNN), are adept at capturing temporal dependencies in sequential data, making them suitable for sign language recognition tasks.

In this paper, we propose a novel approach for sign language to text conversion using LSTM neural networks. We design and implement an LSTM model that learns to map sign language gestures to corresponding textual representations. We present the architecture of our model, discuss the datasets used for training, outline the training methodology, and

evaluate the performance of our approach through comprehensive experiments.



American Sign Language

2. RELATED WORK

Previous research in sign language recognition has explored various approaches, including computer vision-based methods, hand-crafted feature extraction techniques, and deep learning models. Early methods often relied on extracting handcrafted features from video sequences of sign language gestures and using classifiers to recognize gestures. While these methods achieved some success, they often struggled with robustness to variations in lighting conditions, background clutter, and complex hand movements.

More recent approaches have leveraged deep learning techniques, particularly convolutional neural networks (CNNs) and recurrent neural networks (RNNs), for sign language recognition. CNNs have been used for feature extraction from video frames, while RNNs, including LSTM networks, have been employed for sequence modeling and temporal dependency capture. However, most existing approaches

focus on isolated sign recognition rather than continuous sign language translation to text.

3. METHODOLOGY

A) Software Requirements

Operating System: Windows 7 and above.

Language: Python.

Libraries:

a. OpenCV: A robust open-source distribution tailored for real-time applications, particularly adept at managing extensive datasets with computational efficiency. Its primary functions include processing images and videos for tasks such as sign and gesture recognition.

b. TensorFlow: An open-source AI framework renowned for constructing models via data-flow graphs and deploying large-scale neural networks comprising multiple layers.

c. Keras: A high-level neural networks API that works seamlessly with TensorFlow. It's instrumental in constructing neural networks incorporating LSTM layers for sequential data handling.

d. MediaPipe: A versatile framework designed for building ML pipelines to process time-series data like video and audio.

e. Skylearn: Utilized for evaluating system performance using built-in metrics and confusion matrices to determine accuracy.

B) Hardware Requirements

Camera: Good quality, 3MP.

RAM: Minimum 8GB and higher.

GPU: 4GB dedicated.

Processor: Intel Pentium 4 or higher.

HDD: 10GB or higher.

Monitor: 15 inches or 17 inches color monitor.

Input Device: Mouse (Scroll or Optical) / Touchpad.

B) SYSTEM WORKFLOW

Video Capture: Initially, the system captures the video of the person using OpenCV as the input source.

Data Collection: The captured video undergoes data collection using MediaPipe Holistic, which detects facial, pose, and hand landmarks as key points. These key points are extracted and stored in a NumPy array.

Dataset Preparation: The collected data is organized into sequences and formatted into frames of video. The key points extracted from MediaPipe are then pushed into the NumPy array, forming the dataset.

Model Training: The system trains a Long Short-Term Memory (LSTM) deep learning model to recognize sign language gestures. The model architecture consists of three LSTM layers followed by three Dense layers. Training is conducted over 200 epochs with a batch size of 128. The optimization objective is to minimize loss via categorical cross-entropy using the Adam optimizer.

Real-Time Gesture Recognition: After training, the built neural network is deployed for real-time sign language recognition using OpenCV. Gestures captured in the video feed are recognized by the model and displayed as text within a highlighted section.

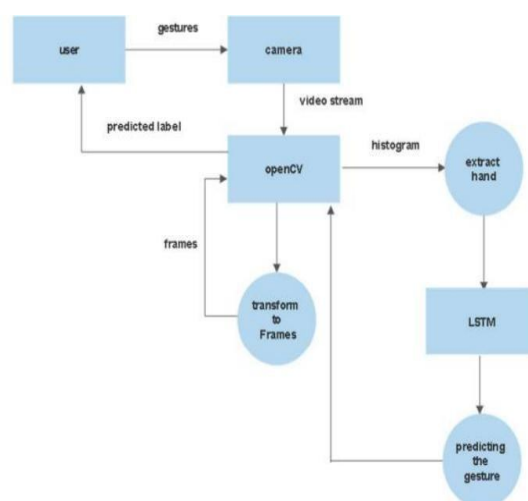


Figure 1: Network architecture of the proposed system.

4. ALGORITHM SPECIFICATION

The step by step procedure to construct this system:

Install and Import Dependencies: Begin by installing necessary libraries and importing required modules in the programming environment.

Detect Key Points with MediaPipe Holistic: Utilize MediaPipe holistic to detect key points such as facial features, body posture, and hand gestures in the captured video.

Extract Key Points: Extract and store the detected key points from the video feed.

