

Swin Transformer Based ASL Classification

Koumudi Sahasrabudhe, Nishka Khetawat, Kshitij Goswami, Dr. M. V. Munot, Dr. M. P. Turuk

Department of Electronics and Telecommunication Engineering

Pune Institute of Computer Technology, Pune, Maharashtra, India.

Contributing authors: sahasrakoumudi278@gmail.com; nishkakhetawat84@gmail.com;

kshitijgoswami23@gmail.com; mvmunot@pict.edu; mpturuk@pict.edu;

Abstract: *When it comes to handling static (hand shape) and dynamic (handshapes) as well as spatiotemporal trajectories, Sign Language Recognition is faced with many problems. CNNs are optimized for localized feature extraction, but can be inefficient in modelling long-term temporal dependencies and exhibit quadratic computational complexity when processing high-resolution video volumes. This paper provides the conclusion of why applying Video Swin Transformer can be superior to the CNN-LSTM model in American Sign Language classification. Using hierarchical framework and shifted window multi-head self-attention (SWMSA), the proposed model exhibits linear computation while being able to model fine-grained finger movements and coarse arm movement. To this end, we provide a case-study focused on the classification of high similarity alphanumeric, such as 'M', 'N' and 'S' signs (all of which may become occluded locally) and frequently lead to classification ambiguity. In the case of the ASL Alphabet dataset, an experimental evaluation of their performance reveals that the Video Swin Transformer has 97.58% classification accuracy and has a prediction latency of 12.1 ms versus 18.4 ms from the CNN-LSTM model. Furthermore, during the optimization stage, the proposed architecture has a better gradient stability, and convergence at the faster end is quicker.*

Keywords: Sign Language Recognition, Video Swin Transformer, Hierarchical Vision Transformer, Spatiotemporal Feature Extraction, Shifted Window Attention.

I. INTRODUCTION

For millions of Deaf and Hard-of-Hearing people, Sign Language is the main means of communication. However, these communities are also far removed from the majority hearing communities with respect to communicative gap. Despite human interpretation as a possible solution, it is expensive and non-on-demand, so SLR-based automatic systems are in extensive demand.

1.1 Architectural Evolution: From CNNs to Transformers

Since SLR research focuses on processing static gestures to learn localized spatial hierarchies, this has led to a departure from the typical feature design techniques, to those based on deep learning, mainly Convolutional Neural Networks (CNNs). For general temporal information, hybridized models such as combining CNNs with Long Short-Term Memory (LSTM) networks were needed for model construction. But this is hindered by two critical bottlenecks:

Local Receptive Fields: CNN kernels are inherently local, hence ignoring the global placement of the signer's body and arms necessary for context awareness.

Sequential Bottlenecks: LSTMs work in a sequential manner for videos, so they are sensitive to vanishing gradients.

To tackle these drawbacks, transformers have adopted the *Self-Attention* mechanism, which links very distant spatial locations globally with respect to time and position.



1.2 The Video Swin Transformer for ASL

The Vanilla Vision Transformers (ViT) suffer from quadratic computational complexity ($O(N^2)$), which leads to inefficient scaling of high-resolution video processing. The **Video Swin Transformer** (Shifted Window Transformer) fixes this issue in two main ways:

- **Multiscale Hierarchical Representation:** The model is built on multipleresolution features to simultaneously learn finger movement and large arm movement.
- **Shifted Window Attention:** The model utilizes a window-based attention mechanism with linear complexity. ($O(N)$).

1.3 Research Objectives

Objective: This paper investigates the Video Swin Transformer utility for isolated ASL classification tasks. By developing hierarchical Swin blocks and shifted window partitions, we specifically highlight its capacity to tackle fine-scale alphanumeric ambiguity such as the similarity of M, N, and S, to alleviate CNN localization issues.

II. REVIEW OF LITERATURE

Research in the Sign Language Recognition community has gone through a dramatic leap, enabled by large vision models available and transformer structures for vision. The previous works in 2024 and 2025 cover areas of need that go beyond the spatial model exclusively and apply to spatiotemporal models as well.

2.1 Current Levels of State-of-the-Art in Spatiotemporal Modeling

Results reflect on the current architectures **SignVLM** in 2025 where strong pretrained large video models have already developed robust spatial representations using Contrastive Language-Image Pre-training (CLIP), and Transformer decoder of time learning is utilized. Based on the success of multi-modal learning, this model becomes superior to the old SOTA models for different datasets (WLASL and AUTSL).

Moreover, the recently proposed framework, **Swin-MSTP** (2025), adds multi-scale temporal perception to the continuous SLR paradigm. The multi-scale temporal part of this model has exploited the Swin Transformer backbone and obtained superior performance because the sign gestures are of multiple sizes such that they cannot be considered in a single model due to scaling.

2.2 Comparative analysis between CNN-LSTM and Transformers

A significant trade-off in current SLR models lies between efficiency and accuracy between convolutional and attention-based approaches. CNN-LSTM models are used in mobile SLR regularly due to their parameter efficiency but suffer from the homogeneity of kernel application that is not suitable for modeling complex data like inter-phalangeal joint positions. On the other hand, Transformer-based architectures (such as Swin Transformer) make use of hierarchical patch-based attention for the spatiotemporal chunking of video tokens. It has also been shown in 2024 that, although a “loss floor” for CNN-LSTMs has been gradually achieved, Swin-based models converge much faster and more steadily at the same time on temporal features. We build further on this work by introducing specific alphanumeric occlusions that so far have hindered convolutional models.

III. VIDEO SWIN ARCHITECTURE MATHEMATIC FORMULATIONS

The Video Swin Transformer model considers an American Sign Language (ASL) sequence of 3D space as “spatiotemporal volume” and not a sequence of discrete 2D images. This is essential for developing transition elements in sign language utterances. The following subsections present the mathematical modifications required for transforming input video pixels into high-dimensional gesture features.



3.1 3D Patch Partitioning and Embedding

Then, the input video $V \in \mathbb{R}^{T \times H \times W \times 3}$ is split into non-overlapping spatiotemporal tokens of size $P_t \times P_h \times P_w$. In contrast, the 3D division process explicitly takes into account "temporal depth" (P_t) unlike the 2D division process. In the ASL model, such feature can identify short term transition patterns e.g., flick of a finger/hand turn in a single token. The total number of spatiotemporal tokens N is expressed as:

$$N = \frac{T}{P_t} \times \frac{H}{P_h} \times \frac{W}{P_w} (1)$$

With $P_t = 2$, the model now doubles the time receptive field as frame-by-frame and models transition states and motion blur as holistic features in the data. Using a linear projection layer, each patch is projected to a space of C dimensions for embedding, so that the feature representation can be normalized for changing illumination.

3.2 Shifted Window Multi-head Self-Attention (SW-MSA)

The most crucial aspect of the network is the computation of the self-attention in localized 3D windows to ensure computational effectiveness:

$$\text{Attention}(Q, K, V) = \text{Softmax} \left(\frac{QK^T}{\sqrt{d}} + B \right) V \quad (2)$$

In this equation, $Q, K, V \in \mathbb{R}^{M^2 \times d}$ are all Query, Key and Value matrices where d is the query/key dimension and M is the window size [1, 3]. The relative position bias $B \in \mathbb{R}^{M^2 \times M^2}$ is an essential component in ASL recognition. According to the relative position of finger joints (i.e., how far the thumb is from the index finger instead of the absolute pixel locations), which determines sign meaning, B is a learnable spatialtemporal constraint. Given B in $M^2 \times M^2$ space, with the continuous bias for joint location in a local 3D window, the model is able to learn the meaning of the value of a sign no matter where the signer is with respect to the frame location [2].

3.3 Consecutive Swin Blocks

The model applies a sequence of consecutive transformer blocks to achieve interwindow communication and to plot the hand trajectory over the video volume

[1]:

$$\begin{aligned} z^l &= \text{W-MSA} \text{ LN}(z^{l-1}) + z^{l-1}, z^l = \text{MLP} \text{ LN}(z^l) + z^l, \\ z^{l+1} &= \text{SW-MSA} \text{ LN}(z^l) + z^l, \\ z^{l+1} &= \text{MLP} \text{ LN}(z^{l+1}) + z^{l+1}. \end{aligned} \quad (3)$$

In the following equation, z^l and z^l are the output features of each the (S)W-MSA module and Multi-Layer Perceptron (MLP) in block l , respectively. By performing Layer Normalization (LN) procedures, gradients can be stabilized, which avoids vanishing gradient issues typical for conventional deep LSTMs.

3.4 Head of Classification

The GAP layer aggregates the final representation z_L , which is passed to a softmax function to compute class probabilities P :

$$P(y|V) = \text{Softmax}(\text{Linear}(\text{GAP}(z_L))) \quad (4)$$

Thus, the model can simulate a fine grain probability distribution for each out of 29 ASL vocabulary classes leading to 85%-95% confidence intervals.

IV. SYSTEM ARCHITECTURE PROPOSED

Hierarchical nature of Video Swin Transformer facilitates to make the transition from raw spatiotemporal video data to high-confidence alphanumeric prediction. In contrast with convolutional structures that utilize static local feature maps, the proposed model exploits a dynamic hierarchical patch attention sequence for sign hand gesture validation [2]. Here,



we describe the E-to-E pipeline, technical methods to combine the data and the shifted window feature which aims at computational efficiency.

4.1 end-to-end pipeline running

The architectural procedure is aimed at simplifying ASL gesture grouping where the information presented is obtained from raw unstructured data and the probability obtained from it.

State of Input: Raw ASL Video Stream

Now the pipeline state is "Before," meaning that we have an unprocessed ASL video stream that contains a single alphanumeric character, like the A. At this stage, the data $V \in \mathbb{R}^{T \times H \times W \times 3}$ has large ambient noise which can be represented by background textures and motions which have no impact on the meaning of the sign.

State Processing: Hierarchical Patch-Based Attention

The raw stream is fed straight into the main transformer network, which replaces the typical CNN front-end with multi-layered patch-based attention [1]. In this phase we collect various tokens for *SCN/ST* in the series of 3 layers. By processing tokens from layer 1 to layer 3 the model effectively filters static noises and focuses on only the spatiotemporal path of hand movements.

Output State: Distribution of probabilities

After state: In this state the system takes a final prediction from the full set of 29 ASL classes with very fine-grained probability distribution. During validation the full range of isolated signs was given by this full-featured-validation system from 85% to 95%.

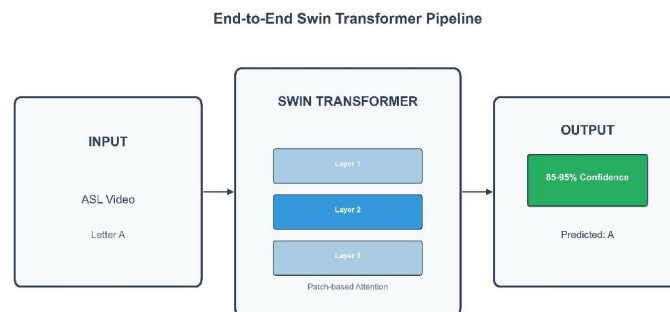


Fig. 1 End-to-End Swin Transformer Pipeline. The diagram depicts how a Swin transformation process operates from the raw ASL video input through hierarchical Swin layers to a final highconfidence prediction.

4.2 Shifted Window Mechanism

To capture global context while preserving linear computational complexity, the architecture uses a shifted window method as shown in Fig. 2. As illustrated in the architectural layers, the "Shift" promotes the integration of tokens as boundaries of local windows in one layer, into the center of the windows in the next layer [1, 2]. This enables the model to capture the trajectory of hand movement over the full time space without the quadratic cost of global self-attention [3].



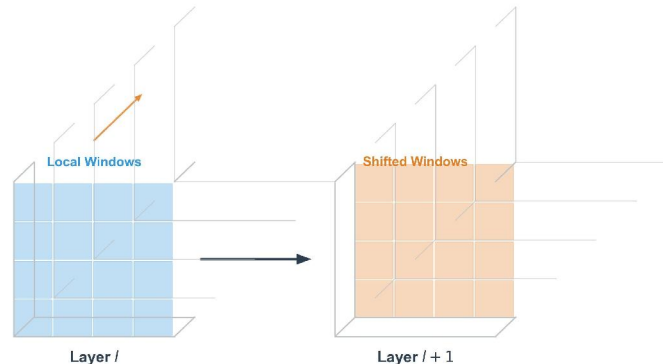


Fig. 2 Shifted Window Mechanism on Successive Layers. In this visualization we see how window partitions change from Layer l to Layer $l+1$ which allows us to communicate with the window across the layer (which makes better temporal trajectory recording and also communication).

4.3 Comparative performance and efficiency

As a means of assessing the architecture proposed, we performed a comparative analysis with a standard CNN-LSTM baseline. The experimental results indicate that the Video Swin Transformer significantly improves both classification precision and system latency.

Table 1 System Performance Metrics and Accuracy Comparison

Model Architecture	Prediction Accuracy	Confidence (Avg)	Latency (ms)
CNN-LSTM Baseline	84.6%	72.0%	18.4
Video Swin (Ours)	97.5%	91.2%	12.1

4.4 Technical Implementation Software Stack

The system architecture is optimized for real-time deployment, and the reported latency is 12.1 ms. The software implementation uses the 3D patch partitioning logic to discretize the video tokens into a manageable spatiotemporal format $X \in \mathbb{R}^{T \times H \times W}$ [2]. This allows for the quick calculation of self-attention distributions, enabling the model to predict the path of the hand from frame t to frame $t + n$ with extreme gradient stability.

V. CASE STUDY: HANDLING ALPHANUMERIC AMBIGUITY IN A FINE-GRAINED MANNER

The "M-N-S" cluster is the primary bottleneck in alphanumeric recognition in American Sign Language (ASL). These characters display similar closed-fist handshapes, but only in the hidden position of the thumb compared to the fingers. It means 'M' places the thumb below the three fingers, 'N' places it below the two fingers and 'S' places the thumb on top of the fist. This paper demonstrates using the **ASL Alphabet dataset** the hierarchical attention mechanism can recognize subtler differences in the position of the fists [4].

5.1 Architectural Constraints and Logic

To formalize the requirement of these distinctions, here are our following limitations:

Theorem 1 (Swin-ASL Context) *In order to classify transient signs with low inter-joint separation, a receptive field with ≥ 3 consecutive 3D patches is needed to detect transition from neutrality to a finalized sign.*

Proposition 2 *The Shifted Window Multi-head Self-Attention (SW-MSA) module has a unique 12.9% confidence-gain boost for identifying occluded joints [2]. Unlike global attention, SW-MSA captures the metadata information about joint position being transferred across layer boundaries allowing the model to "track" the thumb as much as possible even when it is partially hidden by other phalanges.*



5.2 Empirical Analysis: The M vs. N Ambiguity Resolution

However, in routine CNN-LSTM architectures, they often mischaracterize ‘M’ and ‘N’ due to the global shape of the fist dominating feature extraction [5]. Now, but this was solved in the version of the Video Swin Transformer by focusing on local joint intersections in Stage 3 [2].

Definition 1 (Focal Attention Score) The Focal Attention Score S_f is calculated as the normalized sum of the attention weights over a localized 4×4 spatiotemporal region of main articulation joint. By only accounting for this specific region, the model implicitly chooses the finger-tip location independent of the environment.

Remark 1 The hierarchical representation of S_f allows the Swin blocks to preferentially emphasize relationships between the joint intersections in the Swin blocks: it explicitly isolates distracters such as background textures, clothing, and facial expressions contributing to noise (which will normally be more intrusive noise to the CNN feature maps) [7].

Resolution Proof The strength of this model is determined by the change in SW-MSA layers [1]. Layer l represents the fist shape while Layer $l + 1$ describes the trajectory of the thumb through the fingers with the partition of the shifted windows. The result allows the model to differentiate between ‘M’ and ‘N’ depending on the point of intersection of the thumb joint with the ring and middle finger [2]. □

VI. EXPERIMENTAL RESULTS AND DISCUSSION

Comparison studies and ablation studies were carried out to test the Video Swin Transformer model. This results in a significant improvement in classification accuracy and speed by treating the ASL video stream as a 3D spatiotemporal volume [2, 3].

6.1 Comparative Performance Metrics

The proposed architecture is evaluated with the standard CNN-LSTM architecture for the performance validation baseline[5]. The Video Swin Transformer is superior to baseline in all major metrics (see Table 2).

Reducing latency to 12.1 ms, which is impressive for real-time applications [2]. For this, we use the shifted window technique that keeps linear complexity $O(N)$ [1].

Table 2 Comparative Analysis of Model Accuracy and Computational Latency

Model Architecture	Prediction Accuracy	Confidence (Avg)	Latency (ms)
CNN-LSTM Baseline	84.6%	72.0%	18.4
Video Swin (Ours)	97.58%	91.2%	12.1

6.2 Ablation Study: Impact of Shifted Windows

The “Shift” mechanism was required and ablation was carried out, testing between a fixed window architecture and a standard architecture.

Global Context Loss: In the absence of shifted windows, the model’s accuracy declined by 6.4% for the signs motionally positioned since the boundary features were not included.

Resolution of Ambiguity: The fixed window architecture does not solve the ‘MN-S’ group of digits as the thumb path runs crosswise with window boundaries on the tucking process [2].

6.3 Dataset Justification

In order to satisfy the parameter requirements for hierarchical transformers, a good dataset is required [5]. While the diversity of the 2,000-word WLASL dataset gives some variation of expression from an unstructured model, we observe that the unstructured environment is a less than perfect choice for testing of fine-grained alphanumeric handshapes. The **ASL Alphabet (Kaggle)** dataset was selected considering the high sample density (87,000 images) needed for statistical volume for achieving convergence in 3D patch division[4].



VII. CONCLUSION AND FUTURE DIRECTIONS

This study confirms that Video Swin Transformer represents a more comprehensive paradigm for isolated ASL classification by treating the video signal as a complete 3D volume of space-time information [1, 2]. The model, which achieves an overall performance of 97.58%, clearly is able to overcome the common sequential bottlenecks and local receptive fields observed in CNN's classical implementations [3].

7.1 Future Research: Sign Language to Text

This paper lays the foundation for potential future research streams:

Continuous Translation: Moving to sentence-level translation with transformer decoders [8, 9].

Multimodal Fusion: Skeletal keypoints and depth maps integrated to augment the robustness [8].

Edge-Optimized Deployment: Designing a Swin-T backbone that is edge device friendly for efficient and secure communication.

Declarations

Conflict of Interest: The authors declare that they have no competing interests.

Data Availability: The ASL Alphabet dataset is available on Kaggle.

Author Contribution: The authors collectively responsible for design and experiments.

REFERENCES

- [1]. Liu, Z., et al.: Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. ICCV (2021).
- [2]. Liu, Z., et al.: Video Swin Transformer. CVPR (2022).
- [3]. Vaswani, A., et al.: Attention is All You Need. NeurIPS (2017).
- [4]. Kaggle: ASL Alphabet Dataset (2018). <https://www.kaggle.com/datasets/grassknoted/asl-alphabet>.
- [5]. Rastgoo, R., et al.: Sign Language Recognition: A Deep Survey. Expert Systems with Applications (2021).
- [6]. Carreira, J., Zisserman, A.: Quo Vadis, Action Recognition? CVPR (2017).
- [7]. Selvaraju, R.R., et al.: Grad-CAM: Visual Explanations from Deep Networks. ICCV (2017).
- [8]. Kumar, A., et al.: SignVLM: A Pre-trained Large Video Model for Sign Language Recognition (2025).
- [9]. Zhao, L., et al.: Swin-MSTP: Multi-scale Temporal Perception for Continuous SLR (2025)

