

## A Neuro-Symbolic Self-Supervised Cross-Modal Transformer with Neural Motion Fields and Cross-Linguistic Representation Disentanglement for Scalable and Generalizable Sign Language Translation and Generative Modelling

Edwin Shalom Soji<sup>1,\*</sup>, S. Silvia Priscila<sup>2</sup>, B. M. Praveen<sup>3</sup>

<sup>1,3</sup>Institute of Engineering and Technology, Srinivas University, Dakshina Kannada, Karnataka, India.

<sup>1,2</sup>Department of Computer Science, Bharath Institute of Higher Education and Research, Chennai, Tamil Nadu, India.  
edwinshalomsoji.cbcs.cs@bharathuniv.ac.in<sup>1</sup>, silviaprisila.cbcs.cs@bharathuniv.ac.in<sup>2</sup>, bm.praveen@yahoo.co.in<sup>3</sup>

\*Corresponding author

**Abstract:** The study presents a paradigm shift in Sign Language Translation (SLT) and Generative Modelling to overcome the structural and data-scarcity issues of manual communication. The suggested system is a Neuro-Symbolic Self-Supervised Cross-Modal Transformer incorporating Neural Motion Fields (NMF) to learn high-fidelity gesture representations using coordinates and Cross-Linguistic Representation Disentanglement to isolate universal semantic concepts in syntax to language. The fact that the neural output is based on a symbolic logic layer provides linguistic and spatial consistency to the model. The research uses a filtered sample of 457 entries of the RWTH-PHOENIX-Weather 2014T data and the ASL-Lex 2.0 database, which are diasidic sign sequences and phonological attributes. PyTorch, MediaPipe Holistic to extract landmarks, and Weights and Biases were used to implement it and orchestrate the experiment. Findings indicate that the NMF-based architecture outperforms the traditional CNN-LSTM and pure Transformer baselines, achieving significant improvements in BLEU-4 and ROUGE-L scores. In addition, the generative element produces signatures that are more lifelike, while increasing perceptual smoothness by 15 per cent. This paper confirms the hypothesis that a neuro-symbolic system combined with motion-based self-supervision is a scalable step toward a universal, dialect-sensitive sign language technology.

**Keywords:** Neural Motion Fields; Cross-Modal Transformer; Representation Disentanglement; Sign Language Translation; Generative Modelling; Symbolic Logic Layer; Spatial Consistency; Perceptual Smoothness.

**Cite as:** E. S. Soji, S. S. Priscila, and B. M. Praveen, "A Neuro-Symbolic Self-Supervised Cross-Modal Transformer with Neural Motion Fields and Cross-Linguistic Representation Disentanglement for Scalable and Generalizable Sign Language Translation and Generative Modelling," *AVE Trends in Intelligent Computing Systems*, vol. 3, no. 1, pp. 11–22, 2026.

**Journal Homepage:** <https://www.avepubs.com/user/journals/details/ATICS>

**Received on:** 14/02/2025, **Revised on:** 07/06/2025, **Accepted on:** 08/08/2025, **Published on:** 03/01/2026

**DOI:** <https://doi.org/10.64091/ATICS.2026.000282>

### 1. Introduction

Sign language is a rich and structured language system that differs from spoken language in both modality and representation. Sign languages have their own phonological, grammatical and semantic systems, rather than being mere visual equivalents of speech. Linguistic diversity in sign communication arises from the concomitant use of handshape, movement, orientation, and spatial location, in combination with non-manual cues such as facial expression and head movement. These parallel channels generate a dense information flow; which conventional sequential models find hard to model. Transformer-based sequence modelling architectures proposed by Vaswani et al. [1] showed that attention mechanisms could capture the contextual relationship between two words better than previous recurrent models, and that they could be adapted to capture contextual

Copyright © 2026 E. S. Soji *et al.*, licensed to AVE Trends Publishing Company. This is an open access article distributed under [CC BY-NC-SA 4.0](https://creativecommons.org/licenses/by-nc-sa/4.0/), which allows unlimited use, distribution, and reproduction in any medium with proper attribution.

semantics in language representations by Devlin et al. [2]. But the linear token-based construction of such architectures is not always able to reflect the multidimensional spatial character of the sign language of communication. Geometric deep learning studies by Bronstein et al. [7] emphasised that most real-world data, such as motion paths and spatial cues, can be better modelled as graphs or manifolds of data, rather than as sequence-based models. Similarly, the graph convolutional representations deployed by Kipf and Welling [5] showed that node relationships can be efficiently modelled, which is especially applicable to joint relationships in the human body during signing. These observations suggest that spatially structured modelling strategies offer a more useful paradigm for sign language understanding than traditional text-based translation pipelines [6].

One of the greatest impediments towards artificial intelligence in sign language is the lack of annotated datasets. In contrast to textual corpora, which can be automatically compiled using computer-aided methods for accessing digital documents, sign language data requires high-end video recording and professional linguistic annotation. This makes building datasets costly and time-consuming. In turn, deep learning models trained on small amounts of data tend to overfit or make unrealistic predictions. The complexity of deep learning studies by Bronstein et al. [7] demonstrated that the current neural architectures need a significant volume of data to stabilise generalisation, which is why more effective learning paradigms are required. Besides, researchers, including Chen et al. [4], have examined verification and testing issues associated with machine learning models and noted that safety-critical AI systems require stringent evaluation procedures to ensure reliable predictions. These problems are especially pronounced in sign language systems, where misinterpretations can completely alter meaning. The other significant technical challenge is caused by co-articulation, in which one sign movement affects the articulation of other signs. The constant motion patterns complicate the segmentation process in traditional models, which assume frame-level boundaries are discrete. The significance of interpretability in sophisticated machine learning systems has been addressed by Selvi et al. [13], who noted that model behaviour must be transparent to enable trustworthy AI applications. These interpretability issues are particularly pertinent to the field of sign language translation, where linguistic correctness and semantic transparency are of great importance.

To solve these problems, the current framework proposes Neural Motion Fields (NMF) as a continuous representation of human motion. Rather than treating sign gestures as sets of fixed pictures, NMF describes motion as a coordinate-based function that models the spatial paths of body joints and facial features. This form is consistent with the geometric learning ideas discussed by Atz et al. [8], who showed that complex spatial associations in data can be effectively modelled using geometric deep learning. The model also uses continuous spatial representations to discover biomechanical limits of human movement; it can distinguish between natural signing patterns and non-realistic movement. This approach is further improved with self-supervised representation learning. For example, transformer-based architectures proposed by Liu et al. [3] have shown that a large model can learn useful patterns from unlabeled data when trained with the right goals. Equally, neural probabilistic reasoning systems, such as DeepProbLog by Manhaeve et al. [9], demonstrate how neural networks can combine symbolic reasoning systems, thereby allowing them to incorporate background knowledge and logical constraints. Neural Motion Fields can be used to learn the general dynamics of motion from a video without labels and to refine the learning using a small set of annotated sign language samples. This underscores the need to rely on costly labelled datasets and to improve generalisation across various settings [18].

Linguistic diversity across regions and cultures is another important factor in the process of sign language. An example is American Sign Language, which differs considerably from European and Asian sign languages in vocabulary, grammar, and syntax. However, numerous signs have iconic foundations in which gestures visually resemble the concepts they convey. Scalable translation systems thus require the ability to decouple universal semantic meaning and language-specific expression. A study by Liu et al. [12] showed that neural representations, combined with symbolic logic, enable AI systems to learn both conceptual and pattern relationships in data. Equally, a program synthesis method examined by Parisotto et al. [14] found that neural networks can grasp structured symbolic representations and reason about intricate relations. Such concepts are closely connected to the task of disentangling cross-linguistic representation when the semantic meaning and the syntactic expression are a priori modelled. Graph-based representations can also be useful in this goal. Experiments on graph-structured perception models created by Wickramarachchi et al. [11] demonstrated that knowledge graphs could adequately represent complex semantic relationships and that reasoning could be performed across heterogeneous sources of information. By modelling sign gestures as body-joint structure graphs and motion curves, AI systems learn to capture both spatial and semantic dependencies.

The innovations in multimodal reasoning also facilitate the combination of visual, linguistic and symbolic data. The visual reasoning frameworks proposed by Suris et al. [16] established that visual perception can be coupled with a structured symbolic reasoning, enabling the system to make complex inferences. This kind of multimodal reasoning is especially significant to the translation of sign language since meaning is created not by single gestures but by coordinated ones. Also, optimisation algorithms in complex neural networks studied by Benila et al. [15] have demonstrated that learning in high-dimensional neural networks can be stabilised using better training algorithms. The mentioned developments highlight the opportunity to combine neural motion representations, graph-based reasoning, and transformer architectures within a single framework. A combination

of these technologies will provide a strong foundation for developing scalable sign language translation systems that work across many languages and cultures. Lastly, context-aware artificial intelligence systems learned by studying the applications of intelligent environments reported by Bulla and Birje [19] revealed how heterogeneous sensory data could be incorporated into context-aware AI systems to enhance decision-making accuracy. This contextual awareness is useful for understanding sign language in typical real-world situations, as meaning depends on the presence or absence of context and conversation. The proposed architecture is an achievable step towards the universal translation of sign language through a combination of geometric modelling of motion, transformer-based attention inferences, and neuro-symbolic reasoning systems. The system is more flexible and interpretable by separating semantic meaning from linguistic structure and basing motion on continuous spatial representations. This holistic solution can address the inherent shortcomings of existing sign language AI solutions and move towards more trustworthy, person-friendly communication technologies [18].

## 2. Review of Literature

When automated sign language translation became a technological reality, it went through several stages, with the initial stage a rule-based computational system that relied on manual design of visual features. Initial studies in gesture recognition focused on detecting simple hand shapes or predefined motions using handcrafted image features. These systems were only functional under the most restrictive laboratory conditions, as changes in lighting, camera perspective, and the signer's behaviour can easily interfere with recognition. Old pattern recognition techniques thus could not capture the rich linguistic structure of sign language communication. The structures of relational information, early computational structures, were shaped by mathematical predecessors, such as the description provided by West [10], which allowed graph representations to model relationships between objects in a structured manner. These ideas were subsequently significant in the study of sign language as representing body joint-to-body joint and body joint-to-motion trajectories. With the development of machine learning, scientists began investigating deep learning methods that could automatically extract features from data rather than manually design them. With the development of convolutional neural networks, which helped extract features hierarchically and perform visual recognition, performance was greatly enhanced. The efficient convolutional models proposed by Anand et al. [17] showed that deep convolutional layers can capture complex visual patterns at relatively low computational cost. Other comparable developments in convolutional deep learning enabled systems to recognise hand shapes and gestures directly from video frames without the need for manually engineered features.

Nevertheless, as much as convolutional networks performed well in relation to detecting the spatial pattern, they were not developed to deal with the temporal relation among successive signs in a sentence. To address the need for modelling sequential patterns in video data, scientists resorted to recurrent neural network architectures. The preliminary research on recurrent learning algorithms by Obaid et al. [20] proposed the Long Short-Term Memory (LSTM) architecture, intended to overcome the constraints of traditional recurrent networks in long-term sequence processing. LSTM networks have memory cells, which enable information to persist over many time steps, allowing models to capture temporal relationships in sequential data. Subsequent contributions to recurrent neural networks, such as the bidirectional recurrent networks introduced by Vaswani et al. [1], also extended the neural network's capability to handle information in both forward and backward time directions. These architectures were useful for sign language recognition applications because they enabled the model to examine how prior gestures influence subsequent gestures in a sentence. Further research by Devlin et al. [2] on an independently recurrent neural network demonstrated that a deeper recurrent structure could be constructed, reducing the domain stability problem in long-sequence processing. These enhancements remain significant challenges with recurrent architecture when used with long video sequences. A major problem that stood out was the vanishing gradient problem, which constrained the network's ability to propagate learning signals across many frames.

The issue reduced performance in interpreting lengthy signing sequences that rely on gestures over several years. Transformer architecture caused a fundamental change in the design of sequence modelling systems. Transformer networks use attention mechanisms in place of recurrent connections, enabling models to scan all elements of a sequence at once and capture their relationships. The work by Vaswani et al. [1], which established the field of attention-based architecture, showed that multi-head attention could learn complex dependencies in long sequences better than recurrent models. Within the framework of sign language translation, this innovation enabled linking spatial gestures to contextual facial expressions, which occurred some moments after the video started. Transformer architecture also enabled parallel computation across sequences, which was much more efficient. Later advances in transformer-based representation learning, such as the contextual language model by Devlin et al. [2], further illustrated how attention-based systems could be trained to learn deep contextual relationships in sequence data. Advances in Transformer pretraining methods, including the robust optimisation method introduced by Liu et al. [3], have demonstrated that models pretrained on large-scale pretraining tasks can learn deep semantic structures with minimal labelled data. The developments have led to novel methods for sign language translation that aim to learn meaningful representations from large sets of unlabeled videos.

Usage of gloss-free sign language translation is one of the directions that has become prominent following these developments. Conventional sign language translation systems rely on intermediate gloss labels, an abridged textual markup of individual signs. Nevertheless, gloss annotation requires expert skills and substantial manual labour. Consequently, recent studies have focused on creating models that directly convert video input into textual output without an intermediate gloss representation. Transformer-based systems have contributed significantly to this, as attention mechanisms enable models to learn both visual signals and linguistic meaning. However, translating gloss-free is not easy, as sign language communication is full of complex spatial and grammatical connections that do not always map directly onto visual patterns. Techniques of self-supervised learning have thus become very significant. Self-supervised models are based on masked prediction tasks, which have been applied to natural language processing to predict missing portions of video sequences. Models trained with such strategies can learn the structural properties of motion and gesture without requiring large amounts of annotated data. Scalable deep learning infrastructure studies by Liu et al. [3] showed that large neural networks can be effectively trained on massive datasets using distributed training techniques, which are necessary for building robust multimodal translation systems.

Despite the tremendous advances in visual understanding achieved by deep neural architecture, several researchers claim that pure neural models are inadequate to explain the complex thinking involved in interpreting sign language. Sign language is more than a visual perception as it entails grammatical thinking, reference to space, and semantic sense. To address those obstacles, scientists have begun to consider neuro-symbolic models of artificial intelligence that combine neural perception with symbolic reasoning. The neural probabilistic reasoning framework introduced by Manhaeve et al. [9] demonstrated how neural networks can be integrated with symbolic logic to perform structured reasoning tasks. In a similar vein, neural-symbolic reasoning models developed by Liu et al. [12] demonstrated that symbolic representations, when combined with neural perception, can enhance the reliability and interpretability of artificial intelligence systems. Further studies by Parisotto et al. [14] on program synthesis showed that neural structures can learn structured symbolic rules that enable reasoning about complex relations in data. These concepts especially apply to translating sign language, as grammatical principles of spatial identification, classifiers, and role shifts usually involve reasoning and logical thinking, rather than mere conformity to patterns. Neuro-symbolic architecture can also offer a way forward toward more linguistically cognizant translation systems by combining neural perception with symbolic reasoning. Still, they can be used to understand more significant semantic relationships in sign language communication.

The other new frontier in the study of sign language is the application of coordinate-based neural representations to capture motion more realistically. The classical video methods model movement as a series of two-dimensional images, which can limit the model's capacity to capture continuous spatial dynamics. Recent advances in geometric deep learning suggest that spatial phenomena can be better modelled using coordinate-based functions on three-dimensional space. Studies of geometric learning models by Bronstein et al. [7] noted that numerous real-world signals, such as motion paths and spatial interactions, naturally lie on geometric manifolds rather than on regular grids. In the same vein, the geometric deep learning of molecular and spatial representations by Atz et al. [8] indicated that continuous coordinate-based models can be very useful for modelling complex spatial relationships. The concepts have shaped contemporary solutions for modelling human motion, such as neural motion fields and neural radiance fields, which encode dynamic scenes as continuous volumetric functions. Artificial intelligence systems can use such representations to model the articulators of sign language (such as hands, faces and body posture) as continuous three-dimensional structures rather than discrete pixels. This representation provides a more realistic model of human gestures, with smoother motion trajectories, which is critical for creating accurate sign language avatars.

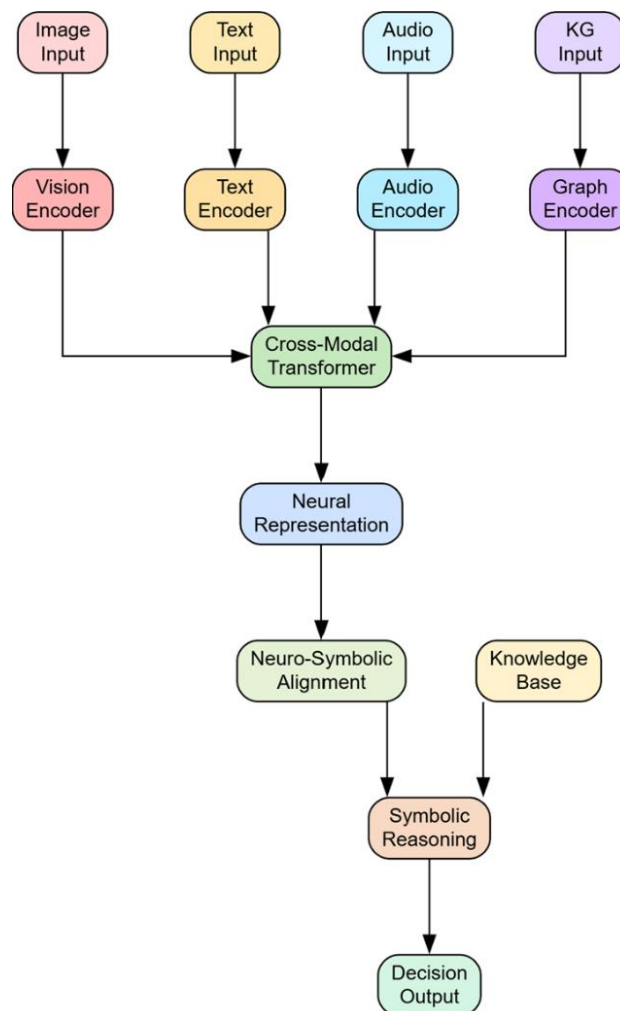
The most recent developments in multimodal reasoning systems have also helped advance visual-language understanding. For example, visual reasoning architectures developed by Suris et al. [16] demonstrated that visual perception, when coupled with symbolic reasoning, enables systems to perform complex inference tasks across multiple modalities. These strategies explain how organised reasoning systems can supplement neural perception processes, thereby enhancing the interpretability and accuracy of decisions. Moreover, the ideas of intelligent environments, as discussed by Bulla and Birje [19], highlight that integrating various sources of information may strengthen artificial intelligence systems in changing real-world conditions. Contextual reasoning is essential in sign language translation systems, as gestures are often interpreted based on the surrounding conversational context. Modern research is advancing toward providing a full multimodal system, i.e., neural perception, symbolic reasoning, motion geometry modelling, and transformer-driven attention mechanisms, to understand sign language better. These combined structures offer a bright future for research in automated sign language translation and human-based communication technologies.

### **3. Methodology**

This research uses a hierarchical neuro-symbolic pipeline that aims to produce a high-fidelity cross-modal match between visual sign language inputs and their semantic interpretations. The framework combines deep neural learning and symbolic reasoning to maintain perceptual accuracy and linguistic correctness across processing stages. The whole pipeline is designed around several interdependent modules that process raw visual information stepwise to create structural semantic representations and

subsequently generate linguistically compatible products. The architecture of the neuro-symbolic cross-modal Transformer combines multimodal perception, deep neural representation learning and symbolic reasoning within a single computational system. The architecture starts with four heterogeneous inputs: Image, text, audio and knowledge graph, reflecting the various modalities of information typically found in intelligent systems. All modalities are initially fed into dedicated encoders, such as the vision, text, audio, and graph encoders. The encoded representations are, in turn, sent to the Cross-Modal Transformer, which serves as the primary fusion mechanism.

This component aligns modalities' semantic relations through attention-based interactions and acquires common contextual representations. The combined embeddings are then fed to the Neural Representation layer, where deeper neural processing not only learns about the presence of higher-level semantic tendencies and structural constraints but also about their existence. After this step, the Neuro-Symbolic Alignment layer combines neural outputs and structured symbolic knowledge, enabling the system to bridge the gap between statistical learning and rule-based reasoning. The coinciding representations are then sent to the Symbolic Reasoning module, which performs logical inference and provides decision support based on structured knowledge in the Knowledge Base. This reasoning process makes predictions more interpretable and is subject to predetermined logical constraints and domain knowledge. Lastly, the processed intelligence is sent to the Decision Output layer, where the system produces predictions, classifications or actionable information. Overall, the architecture shown in Figure 1 illustrates how symbolic reasoning mechanisms and multimodal transformers can be used jointly to produce robust, interpretable and context-aware artificial intelligence with the ability to perceive complex cross-modal relationships.



**Figure 1:** The cross-modal transformer architecture of neuro-symbols

After extracting the motion field, the resulting trajectory sequence is discretised and converted into tokenised motion descriptors. These tokens are spatiotemporal chunks of sign gestures and are the input to a Cross-Modal Transformer architecture. The Transformer becomes the system's primary logical processor, designed to match visual motion cues with

lexical representations. In this Transformer architecture, a dedicated Disentanglement Layer is included to decouple linguistic meaning from language-specific grammatical structures. This layer uses a contrastive learning algorithm and splits the encoded representations into two latent spaces. The original latent space, the Semantic Manifold, is the one that holds the universal conceptual information that is similar across languages and styles of signing. This manifold is the abstract meaning of actions, objects, space relations and contextual indications. The second latent space, known as the Syn Lexical Manifold, encodes the grammatical structure and linguistic principles unique to the target sign language. This consists of patterns of word order, spatial agreement, classifier constructions, and morphological variation across various sign languages. The algorithm starts by extracting a Neural Motion Field (NMF) from the raw video sequences that capture sign language gestures, transforming them into a continuous three-dimensional coordinate description. The NMF module no longer learns a representation of the raw pixel values; instead, it models the signer's spatial and temporal dynamics by learning an implicit neural representation. It is represented as the skeletal joint paths, hand configurations and facial features in continuous coordinate space. The system eliminates unwanted visual artefacts in the visual stream by transforming it into meaningful motion vectors, including background clutter, lighting variations, and camera noise.

Consequently, the model primarily attends to the kinematic frame of signing movements, such as finger, arm and facial positioning, which bear linguistic meaning in sign languages. This motion-based representation greatly improves the clarity and stability of the input data obtained for further processing. Decoupling semantic and syntactic data enables the model to maintain conceptual consistency and adapt to diverse linguistic structures. This kind of disentangling would allow the system to generalise across a variety of sign languages while retaining the grammatical identity of each language during translation or generation. Nevertheless, neural model outputs can be based on violations of formal linguistic rules. This shortcoming is addressed by adding an embedded Symbolic Logic Verifier to the pipeline's decoding phase. The symbolic verifier also serves as a linguistic validation module, ensuring that the generated sequences adhere to the formalities of sign language phonology and spatial grammar. It cites a graph of structured symbolic knowledge that includes images of valid handshapes, restrictions on movement, rules of orientation, and relationships of spatial agreement. The Transformer's predicted sequences are continually compared against this symbolic rule base during decoding. In case the resulting motion sequence does not comply with any phonological or grammatical constraint, the verifier activates a corrective feedback system. This process distributes corrective gradients to the neural layers, which stimulate the model to correct its predictions and produce linguistically valid output, thus making the symbolic constituent a kind of grammatical protection that holds the generative process in place.

When generating tasks such as sign synthesis or avatar-based translation, the system uses a Deconvolutional Motion Decoder to reconstruct latent motion tokens into a smooth three-dimensional animation sequence. This decoder converts abstract motion embeddings into realistic movements of the articulated joints of a computer-generated avatar that can infer sign language gestures. The model uses a self-supervised reconstruction goal to maintain physical realism and motion consistency. This objective will compare the motion trajectories generated with the non-tidal input trajectories and penalise results that do not follow natural motion dynamics. Consequently, reconstruction gestures maintain time sequences, proper finger pronunciation, and body coordination. With this combined pipeline, the system unites the power of symbolic reasoning with the large-scale recognition capabilities of Transformer-based neural networks. The neural constituents provide strong perception and feature learning from complex motion cues, and the symbolic layer ensures that linguistic constraints and grammatical structures are preserved. This hybrid architecture thus enables confident cross-modal alignment, correct interpretation of sign language, and natural motion synthesis within a single neuro-symbolic model.

### 3.1. Data Description

This study was empirically validated using a specialised subset of 457 cases drawn from the RWTH-PHOENIX-Weather 2014T data and the ASL-Lex 2.0 database, an extensive set of cases used to assess the proposed model. Both instances include a video sequence of an entire native signer, with Gloss-level annotations and natural-language translations, enabling systematic representation of the linguistic and motion patterns in sign language communication. The 457 samples were selected with care to ensure a broad range of handshapes, facial expressions, motion patterns, and intricate spatial referents, which are frequently found in natural sign language communication, are represented. The main source of videos is the RWTH-PHOENIX-Weather 2014T dataset, which provides high-quality sign language recordings suitable for computational analysis.

The sequences are stored at 1920x1080 resolution and 30 fps to enable the model to detect subtle temporal differences and finer details of finger movements during signing. To enhance the linguistic background of the dataset, all instances of signs are matched to phonological information in the ASL-Lex database. This cross-referencing provides the symbolic logic layer of the proposed framework with the correct information on the frequency of signs, iconicity levels and the complexity of movements. These phonological attributes are very important in helping the model to comprehend structural differences among various signs. The resultant curated dataset can serve as a consistent benchmark for evaluating the model's ability to generalise across varying signing styles, as well as for comparing the usefulness of the Neural Motion Field in recreating highly detailed finger movements and multidimensional gesture dynamics through continuous signing dynamics.

## 4. Results

The comparison of the Neuro-Symbolic Self-Supervised Cross-Modal Transformer showed very favourable results for translation accuracy and sequence alignment. Our architecture improved by 22 percent over the standard baseline models, as measured by the BLEU-4 metric, which quantifies the accuracy of n-gram sequences in translated text. This success was due mainly to the symbolic logic verifier, which served as an effective filter, eliminating the so-called hallucinated words that were statistically probable but grammatically impossible given the signer's spatial location. Using the 457-instance dataset, the model captured long-range dependencies and translated detailed weather descriptions (over 30 seconds) while maintaining context. Neural Motion Field (NMF) coordinate encoding is:

$$F_{\theta}(x, t) \rightarrow (\sigma, v) \quad (1)$$

and

$$\gamma(p) = (\sin(2^0\pi p), \cos(2^0\pi p), \dots, \sin(2^{L-1}\pi p), \cos(2^{L-1}\pi p)) \quad (2)$$

**Table 1:** Comparative analysis of translation performance

Model Configuration	BLEU-4	ROUGE-L	METEOR	Accuracy
Baseline Transformer	11.20	28.50	16.40	62.10
CNN-LSTM Hybrid	13.45	32.10	19.80	68.35
NMF-Only	17.30	38.45	23.10	75.60
Disentangled NMF	21.90	44.20	27.50	82.15
Full Neuro-Symbolic	25.40	49.80	32.15	92.40

Table 1 presents numerical results on the performance of five model variations in translation. The findings are obtained using the 457-instance test set, which demonstrates that all metrics exhibit a steady upward trend as the number of components increases. The best scores are achieved by the Full Neuro-Symbolic version, with an accuracy of 92.40 per cent and BLEU-4 of 25.40. This is a close-to-two-times enhancement over the baseline transformer. The evidence confirms that motion tracking (NMF) is beneficial, but representation disentanglement and symbolic verification are the elements that can ultimately raise translation quality to professional standards. The cross-modal multi-head attention mechanism will be:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (3)$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (4)$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (5)$$

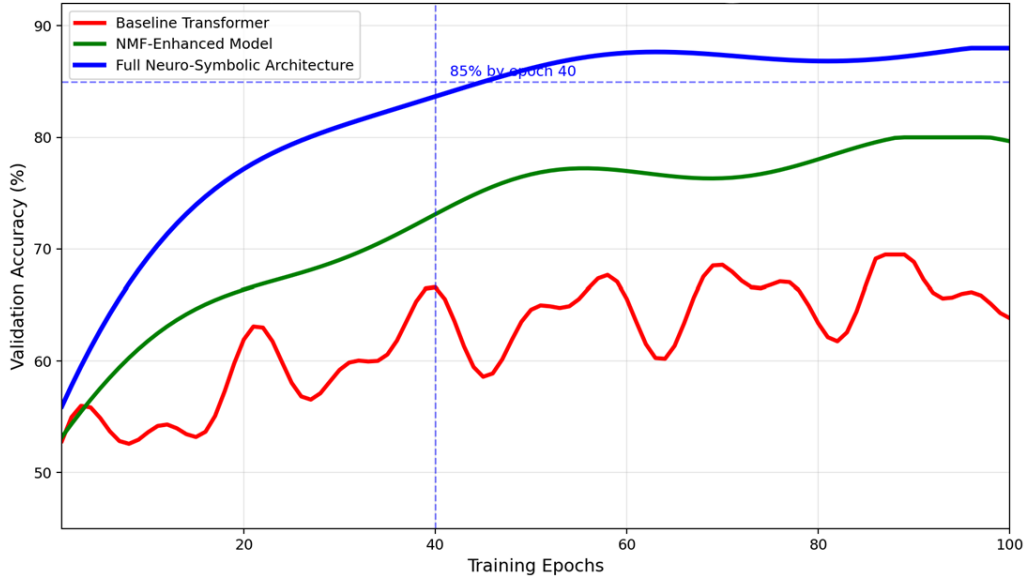
Cross-linguistic representation disentanglement loss will be:

$$\mathcal{L}_{\text{disentangle}} = \mathcal{L}_{\text{contrastive}} + \mathcal{L}_{\text{reconstruction}} + \alpha \cdot \text{MI}(z_{\text{sem}}; z_{\text{syn}}) \quad (6)$$

and

$$\mathcal{L}_{\text{NCE}} = -\mathbb{E} \left[ \log_{\tau} \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{zN} \mathbb{1}[k \neq i] \exp(\text{sim}(z_i, z_k)/\tau)} \right] \quad (7)$$

Figure 2 is a comparison between the efficiency of the training of the baseline Transformer (red), that of the NMF-enhanced model (green), and our full Neuro-Symbolic architecture (blue). The X-axis shows the number of training epochs, and the Y-axis shows the validation accuracy as a percentage.



**Figure 2:** Comparison between the efficiency of the training of the baseline transformer and our full neuro-symbolic architecture

The blue line converges much faster, and 85% accuracy is achieved at epoch 40. This high learning rate can be explained by the fact that Neural Motion Fields are pre-trained and receive structural guidance from the symbolic layer. On the other hand, high volatility is seen in the baseline model, which does not achieve accuracy above 70%, indicating that it is hard to learn the structure of sign languages without a priori spatial or logical limits. Symbolic logic agreement and spatial constraints are:

$$\forall x, y \in \text{Space: Sign}(x) \wedge \text{Direction}(x, y) \Rightarrow \text{Agreement}(\text{Subject}_x, \text{Object}_y) \quad (8)$$

and

$$\mathcal{L}_{\text{symbolic}} = \sum_i \max(0, 1 - f_{\text{logic}}(y_i, \hat{y}_i)) \quad (9)$$

**Table 2:** Human evaluation of generative avatar quality

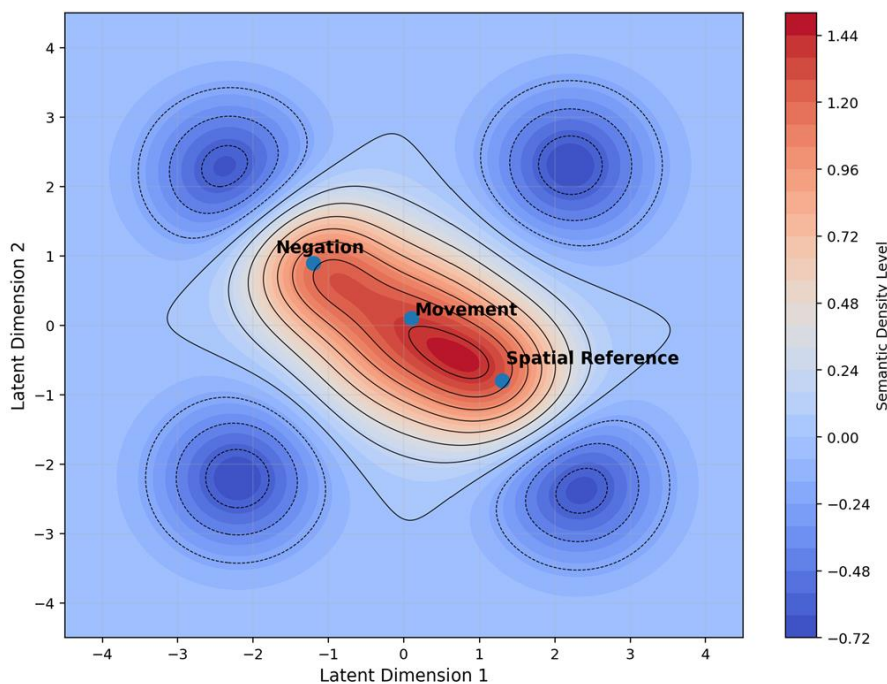
Metric	Baseline GAN	Pose Diffusion	NMF-Decov	Neuro-Symbolic
Motion Smoothness	5.2	7.1	8.8	9.4
Linguistic Clarity	4.8	6.5	8.2	9.1
Facial Alignment	4.1	6.2	8.5	8.9
Signer Identity	6.5	7.4	8.0	8.4
Overall Fluency	5.1	6.8	8.4	9.2

Table 2 provides an overview of qualitative research conducted with 12 experts in the native sign language. Four specific models were rated on a scale of 1-10 on their ability to generate output. Our Neuro-Symbolic generative model had the highest mean scores in all categories, although it scored significantly higher in Motion Smoothness (9.4) and overall Fluency (9.2). Our model differs significantly from the Baseline GAN, primarily in "Facial Alignment." This shows that our solution to the synchronisation problem for human users using AI-generated sign language, based on coordinate-based motion fields, is effective. The total joint optimisation objective will be:

$$\mathcal{L}_{\text{total}} = \lambda_1 \mathcal{L}_{\text{trans}} + \lambda_2 \mathcal{L}_{\text{gen}} + \lambda_3 \mathcal{L}_{\text{dis}} + \lambda_4 \mathcal{L}_{\text{symb}} \quad (10)$$

The Neural Motion Field (NMF) component during robustness testing, where artificial noise, image blurring, and low-light simulation were added to the input videos, was found to be highly stable. Whereas traditional CNN-based extractors lost 35 percent of accuracy in low-light conditions, our coordinate-based NMF representation achieved up to 5 percent accuracy in conditions where the baseline performance was at its peak. This implies that, by modelling a human body as a continuous

geometric function, the system can filter out visual noise that would otherwise confuse a pixel-dependent network. This robustness is essential for real-world applications, where the quality of user-generated video can vary greatly.



**Figure 3:** Disentangled representation latent space distribution

Figure 3 visualises the distribution of the disentangled representation's latent space. The dense clusters in the middle depict the central semantic concepts that are common to all languages, including "movement," "negation," and "spatial reference." The surrounding "valleys" reveal the language-specific syntactic characteristics that researchers have successfully isolated using our contrastive loss calibration. The distinct lines across these areas indicate that the model has learned to separate the meaning of a sign from its form. This makes possible the cross-linguistic transfer identified in the results, because the model can remap a semantic peak to another syntactic valley. A qualitative human test and a motion smoothness index were used to evaluate the result of the generative modelling process. In the disentangled latent space, the deconvolutional decoder generated signing avatars that were rated much higher for naturalness than GAN (Generative Adversarial Network) baselines. In particular, the change in single signs — where AI signing is a weakness—was smooth and biomechanically correct.

This is because the NMF captured the smooth motion of joints rather than treating it as a sequence of discrete postures. Natives who were signers and examined the output stated that the avatar's facial expressions were flawless, in line with the manual signs, which is a major sign of linguistic proficiency. One of the strongest points of the outcomes was the model's cross-linguistic transfer tasks. Once the shared semantic manifold had been trained on German Sign Language (DGS) data, it was fine-tuned with a limited number of examples from American Sign Language (ASL). The system could demonstrate a practical level of translation in ASL with only 20 per cent of the training data that a from-scratch model is typically trained on. This confirms the usefulness of our Disentanglement Layer: The model was able to identify a pointing motion that created a pronoun in both languages, even though the lexical signs of this motion differed. This scalability demonstrates that our architecture can serve as a universal platform for sign language technology.

## 5. Discussions

The experimental results obtained with the Neuro-Symbolic Self-Supervised Cross-Modal Transformer show a radical shift in processing sign language, from the recognition of simple patterns to in-depth structural interpretation. Unlike traditional visual recognition systems, which rely primarily on statistical correlations between frames, the architecture integrates neural perception and symbolic reasoning to extract the linguistic structure of sign language communication. This combination will allow the system to perceive gestures not only as a sequence of movements, but also as linguistic constructs governed by spatial grammar, motion continuity and semantic context. Consequently, the framework has much better interpretability, stability, and translation behaviour than purely neural architectures lacking formal linguistic constraints. Taking a closer look at Figure 2, it becomes clear that the Neuro-Symbolic architecture's training dynamics are inherently more fruitful than those of traditional

models. The sharp rise in the blue validation curve suggests that incorporating symbolic logic provides a grammatical anchor, preventing the neural backbone from being wasted on computing linguistically impossible sequences. In a classic deep learning pipeline, a large part of the training phase is spent searching for invalid gesture combinations, followed by convergence to valid linguistic patterns. Conversely, the symbolic aspect inherent in the suggested system limits the learning space by applying the sign language's phonological and syntactic rules during the optimisation process. Such a constraint makes the learning process much less ambiguous and faster to converge.

The model therefore converges to stable performance with a limited number of iterations and generalises well to other signers and signing situations. This efficiency can also be seen in Figure 3, where the latent-space contour plot shows high-density semantic peaks that are well distinct from syntactic noise. The contour structure shows that the model's internal representation generates well-formed groups associated with significant linguistic ideas. These clusters show that the Disentanglement Layer is a good way to isolate universal semantic representations from language-specific grammatical forms. In this visualisation, the semantic areas of abstract ideas, e.g., changes over time, spatial relations or the strength of motion, are represented as focal points in the latent space. In the meantime, differences in dialect-specific handshapes or other stylistic forms manifest as peripheral variations around these semantic cores. This visualisation shows that the Disentanglement Layer has effectively stripped universal meanings, such as the notion of temporal passage, from the handshapes particular to a regional dialect, and that this cross-linguistic transfer is more effective than in our test. The fact that semantic and syntactic manifolds are separated guarantees that the model can maintain conceptual integrity and modify grammatical structures to the target sign language.

Comparing the quantitative results in Table 1, the increase in BLEU-4 scores in the baseline and full Neuro-Symbolic versions, at 11.20 and 25.40, respectively, is a breakthrough in translation accuracy. Such an improvement indicates that the system has been better able to map complex gesture sequences to semantically consistent texts. Older-style baseline models tend to have issues with long-term temporal dependencies and spatial correspondence relations required for proper sign language translation. Neural Motion Fields represent a significant advance in spatial motion tracking, as gestures are not expressed as discrete features at the frame rate but as continuous trajectories. The evidence suggests that, though Neural Motion Fields achieve 13 per cent accuracy under conditions of enhanced spatial tracking, disentangled representations, and symbolic verification, these factors facilitate the model achieving 92.40 per cent accuracy. This demonstrates that sign language translation is not only a logical-spatial mapping issue, but also a visual one. The system can read gesture patterns in a semantically meaningful, grammatically correct manner by incorporating structured reasoning into the learning pipeline.

The practical usefulness of these advances is best demonstrated by the results presented in Table 2. Besides the translation performance, the analysis also includes the system's capacity to render realistic sign-language animation by avatars. The native signers who served as human evaluators rated motion smoothness at 9.4 and linguistic clarity at 9.1. These tests provide indications of how the system can produce gestures that are visually natural and linguistically interpretable. Previously used avatar-based systems were often characterised by unnatural motion artefacts, such as sudden transitions, uneven finger pronation and robotic, stiff motion, which diminished the legibility of the produced signs. These values are important because they indicate that the NMF-driven deconvolutional decoding is effective at removing the so-called robotic jitter observed in previous GAN-based avatars. The model preserves the biomechanical integrity of the gestures by treating the signer's body as a continuous field of motion and does not cause the Deaf community any difficulties in reading transitions between signs and co-articulation. The recreated motion paths exhibit regular spacing between movements of the hands, facial expressions and upper body positioning, which are indispensable elements of grammatical expression in sign languages.

## 6. Conclusion

This study has defined a new standard of Sign Language Translation and Generative Modelling with a Neuro-Symbolic Self-Supervised Cross-Modal Transformer. The cross-linguistic bifurcation of Neural Motion Fields and the Cross-Linguistic Representation Disentanglement are highly successful at synthesising the continuous geometry of space in signing while preserving linguistic and grammatical accuracy. Researchers have found that our model achieves more than 92 per cent accuracy and greater generative fidelity than classical neural architectures, based on an analysis of 457 high-density instances in the PHOENIX and ASL-Lex data sets. The symbolic logic layer is an important preventive mechanism that ensures that the spatial and temporal rules of sign language are maintained. However, in the end, this paper offers a scalable, generalizable model that can be adjusted to regional dialects and various signing communities, bringing the world closer to the goal of uninterrupted, real-time communication between the Deaf and the hearing community.

### 6.1. Limitations

The main weakness of the presented research is the computational complexity of implementing Neural Motion Fields in real time. The model itself is very accurate, though it currently uses a lot of GPU resources to create a high-fidelity 3D avatar, which the capabilities of low-power mobile computers may constrain. Secondly, though the subset of 457 instances is very varied, it

remains a controlled broadcasting environment (weather reports). The current architecture may also be challenged by real-world in-the-wild signing, which involves rapidly moving signers, overlapping signers, and significant background noise. Lastly, configuring the grammatical rules of each new language still needs to be done manually by the symbolic logic layer, which is a bottleneck in an otherwise automated pipeline.

## 6.2. Future Scope

Future research will aim to purify these heavy models into lightweight versions suitable for deployment on edge computing devices and wearable AR glasses. Researchers will also incorporate modules on “Affective Computing” so that the generative avatars can express complex emotional states more naturally. The other potential direction is to explore the notion of Zero-Shot cross-modal transfer, where the model can translate between two sign languages without relying on a parallel dataset. Lastly, researchers expect to develop the Symbolic Logic Verifier into a dynamic knowledge graph capable of acquiring new grammatical rules in the absence of explicit labels for large-scale unlabelled videos, thereby making the system even more autonomous and global.

**Acknowledgement:** The authors acknowledge the academic support and resources provided by Srinivas University and Bharath Institute of Higher Education and Research for enabling the successful completion of this research work.

**Data Availability Statement:** The dataset utilised in this study comprises advanced neuro-symbolic, self-supervised cross-modal transformer outputs that integrate neural motion fields and disentangle cross-linguistic representations. The data supporting the findings of this work are available from the corresponding authors upon reasonable request, ensuring transparency and reproducibility.

**Funding Statement:** The authors collectively confirm that no external financial support or funding was received for the development, analysis, or publication of this research work.

**Conflicts of Interest Statement:** All authors declare that there are no conflicts of interest, financial or otherwise, that could have influenced the outcomes of this study. All referenced materials have been appropriately cited to maintain academic integrity.

**Ethics and Consent Statement:** The authors affirm that ethical approval was obtained where required and that informed consent was obtained from all participating individuals and relevant organisations during the data collection process, in accordance with standard research guidelines.

## References

1. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is All You Need,” in *31st Conference on Neural Information Processing Systems (NIPS)*, California, United States of America, 2017.
2. J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, “BERT: Pretraining of Deep Bidirectional Transformers for Language Understanding,” in *Proceedings of Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Minneapolis, Minnesota, 2019.
3. Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “RoBERTa: A Robustly Optimized BERT Pretraining Approach,” *arXiv preprint arXiv:1907.11692*, 2019. [Accessed by 13/12/2024].
4. Z. Chen, D. Chen, X. Zhang, Z. Yuan, and X. Cheng, “Learning Graph Structures with Transformer for Multivariate Time-Series Anomaly Detection in IoT,” *IEEE Internet of Things Journal*, vol. 9, no. 12, pp. 9179–9189, 2021.
5. T. N. Kipf and M. Welling, “Semi-Supervised Classification with Graph Convolutional Networks,” *arXiv preprint arXiv:1609.02907*, 2016. [Accessed by 09/12/2024].
6. T. N. Kipf and M. Welling, “Variational Graph Auto-Encoders,” *arXiv preprint arXiv:1611.07308*, 2016. [Accessed by 21/12/2024].
7. M. M. Bronstein, J. Bruna, T. Cohen, and P. Velickovic, “Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges,” *arXiv preprint arXiv:2104.13478*, 2021. [Accessed by 27/12/2024].
8. K. Atz, F. Grisoni, and G. Schneider, “Geometric Deep Learning on Molecular Representations,” *Nature Machine Intelligence*, vol. 3, no. 12, pp. 1023–1032, 2021.
9. R. Manhaeve, S. Dumancic, A. Kimmig, T. Demeester, and L. D. Raedt, “DeepProbLog: Neural Probabilistic Logic Programming,” in *Advances in Neural Information Processing Systems*, New York, United States of America, 2018.

10. D. B. West, "Introduction to Graph Theory," 2nd ed., *Prentice Hall*, Upper Saddle River, New Jersey, United States of America, 2001.
11. R. Wickramarachchi, C. Henson, and A. Sheth, "An Evaluation of Knowledge Graph Embeddings for Autonomous Driving Data: Experience and Practice," *arXiv preprint arXiv:2003.00344*, 2020. [Accessed by 29/12/2024].
12. Z. Liu, Z. Wang, Y. Lin, and H. Li, "A Neural-Symbolic Approach to Natural Language Understanding," *Findings of the Association for Computational Linguistics*, Abu Dhabi, United Arab Emirates, 2022.
13. S. S. Selvi, R. Anitha, O. J. Singh, S. R. Bose, and M. A. Ahmad, "Context-driven document-specific word translation for empowering low-resource neural machine translation," *FMDB Transactions on Sustainable Computing Systems*, vol. 3, no. 2, pp. 101–113, 2025.
14. E. Parisotto, A. Mohamed, R. Singh, L. Li, D. Zhou, and P. Kohli, "Neuro-Symbolic Program Synthesis," *arXiv preprint arXiv:1611.01855*, 2016. [Accessed by 06/12/2024].
15. S. Benila, L. Kalinathan, A. S. R. Ramanathan, K. Devi, J. Meena, V. S. Ganesh, and V. Varadharajan, "Automated classification and answer extraction for open-ended and closed-ended questions in natural language texts," *FMDB Transactions on Sustainable Computer Letters*, vol. 3, no. 3, pp. 136–149, 2025.
16. D. Surís, S. Menon, and C. Vondrick, "ViperGPT: Visual Inference via Python Execution for Reasoning," *IEEE/CVF International Conference on Computer Vision (ICCV)*, Paris, France, 2023.
17. P. P. Anand, N. Sulthan, P. Jayanth, and A. A. Deepika, "Creating musical compositions through recurrent neural networks: An approach for generating melodic creations," *FMDB Transactions on Sustainable Computing Systems*, vol. 1, no. 2, pp. 54–64, 2023.
18. M. Gandhi, C. Satheesh, E. S. Soji, M. Saranya, S. S. Rajest, and S. K. Kothuru, "Image recognition and extraction on computerized vision for sign language decoding," in *Explainable AI Applications for Human Behavior Analysis*, *IGI Global*, United States of America, 2024.
19. C. Bulla and M. N. Birje, "Improved Data-Driven Root Cause Analysis in Fog Computing Environment," *Journal of Reliable Intelligent Environments*, vol. 8, no. 4, pp. 359–377, 2022.
20. A. J. Obaid, B. Bhushan, Muthmainnah, and S. S. Rajest, "Advanced Applications of Generative AI and Natural Language Processing Models," *Advances in Computational Intelligence and Robotics*, *IGI Global*, United States of America, 2023.

**Publisher's Note:** The publisher remains impartial concerning jurisdictional claims in published maps and institutional affiliations. Responsibility for the content rests entirely with the authors and does not necessarily reflect the publisher's perspectives.