# Sign Language Recognition via Skeleton-Aware Multi-Model Ensemble

Songyao Jiang, Bin Sun, Lichen Wang, Yue Bai, Kunpeng Li and Yun Fu
Northeastern University, Boston MA, USA

*Abstract*—Sign language is commonly used by deaf or mute people to communicate but requires extensive effort to master. It is usually performed with the fast yet delicate movement of hand gestures, body posture, and even facial expressions. Current Sign Language Recognition (SLR) methods usually extract features via deep neural networks and suffer overfitting due to limited and noisy data. Recently, skeleton-based action recognition has attracted increasing attention due to its subject-invariant and background-invariant nature, whereas skeleton-based SLR is still under exploration due to the lack of hand annotations. Some researchers have tried to use off-line hand pose trackers to obtain hand keypoints and aid in recognizing sign language via recurrent neural networks. Nevertheless, none of them outperforms RGB-based approaches yet. To this end, we propose a novel Skeleton Aware Multi-modal Framework with a Global Ensemble Model (GEM) for isolated SLR (SAM-SLR-v2) to learn and fuse multi-modal feature representations towards a higher recognition rate. Specifically, we propose a Sign Language Graph Convolution Network (SL-GCN) to model the embedded dynamics of skeleton keypoints and a Separable Spatial-Temporal Convolution Network (SSTCN) to exploit skeleton features. The skeleton-based predictions are fused with other RGB and depth based modalities by the proposed late-fusion GEM to provide global information and make a faithful SLR prediction. Experiments on three isolated SLR datasets demonstrate that our proposed SAM-SLR-v2 framework is exceedingly effective and achieves state-of-the-art performance with significant margins. Our code will be available at https://github.com/jackyjsy/SAM-SLR-v2
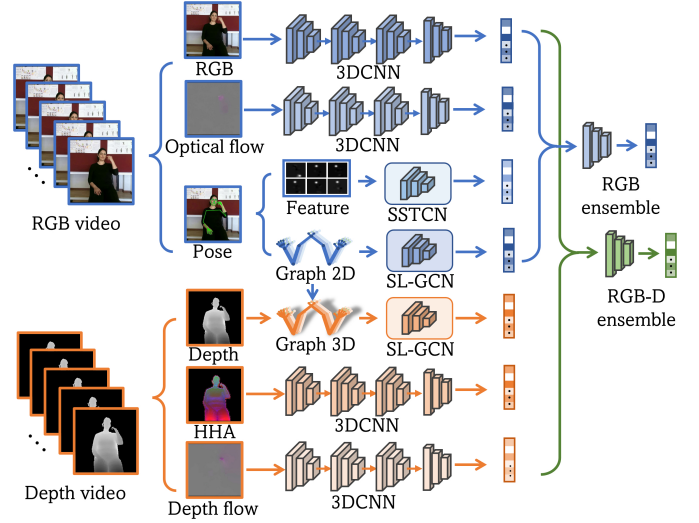
Fig. 1. Concept of the Skeleton Aware Multi-modal Sign Language Recognition Framework with Global Ensemble Model (SAM-SLR-v2). All local and global motion information is extracted and fused to make final predictions.

## I. INTRODUCTION

Sign language [1] is a visual language performed with the dynamic movement of hand gestures, body posture, and facial expressions. It is a widely-used alternative approach for deaf and mute people to communicate effectively. Understanding and performing sign language requires a substantial time of learning which is beyond feasible for the public, which leads to a barrier between deaf-mute people and others. Moreover, sign language is dependent on language [2], [3], [4] (*e.g.*, English and Chinese) and culture [5] (*e.g.*, American and British) that further limits its popularity. Sign Language Recognition (SLR) aims to help deaf-mute people communicate smoothly with others in their daily life by automatically interpreting sign language. As machine learning and computer vision make great progress in the past decade, SLR has drawn much research attention.

SLR contains two tasks, isolated SLR and continuous SLR. The isolated setting is a fine-grained and fully supervised classification at word (gloss) level, while the continuous setting maps whole videos into sentences (*i.e.*, sequences of glosses) in a weakly supervised manner. SLR is a more challenging problem compared with conventional action recognition for

the following reasons. First, sign language requires both delicate hand gestures and global body motion to distinctly and accurately express its meaning. Facial expressions may also be used to express emotions and emphasis. Similar gestures may sometimes express different and even opposite meanings. Second, different signers may perform sign language differently (*e.g.*, left-hander, right-hander, different speed, body shape, and localism), which makes SLR more challenging. Collecting more samples from as many signers as possible is desired yet expensive. To this end, it requires a finer-grained model to capture the delicate dynamics of the whole human body. Besides, we expect the model to effectively merge the information from different cues when making the final recognition.

Traditional SLR methods mainly deploy handcrafted features (*e.g.*, HOG [6] and SIFT [7]) with conventional classifiers (*e.g.*, kNN and SVM [4], [8], [9]). As deep learning achieves significant progress in video representation learning, general temporal dynamics learning methods (*e.g.*, RNN, LSTM, and 3DCNNs) are first explored for SLR in [10], [11], [12], [13]. To more effectively capture the local motion and further improve the accuracy, attention mechanisms are introduced in [14], [15]. Multi-cue approaches are proposed in [16], [17], [18] as well. However, due to the extensive vocabulary sizes and the limitation of annotated data, the

current methods are still not effective enough for robust SLR.

Recently, in human action recognition, skeleton-based approaches have become more popular and drawn increasing attention due to their strong adaptability to dynamic circumstances and complicated backgrounds [19], [20], [21], [22], [23]. Besides, the skeleton-based methods provide complementary information to RGB-based modalities, hence their ensemble results further improve the overall recognition rates. However, there exist some barriers that hinder the extension of skeleton-based methods to the SLR task. The skeleton-based methods for action recognition rely on ground-truth annotations of human bodies provided by motion capture systems, which restrict themselves to limited datasets captured in lab-controlled environments. Apart from that, most motion capture systems only focus on the fewer body keypoints excluding hand keypoints. The provided keypoints are insufficient to recognize sign language performed with delicate hand gestures and motions. Some researchers attempt to obtain hand poses using separate hand detectors with keypoint estimators and propose to use an RNN-based model to recognize the sign language [24]. Unfortunately, their estimated hand keypoints are unreliable, and the RNN-based model cannot learn the dynamics of the human skeleton properly.

In this work, we focus on the isolated SLR task. We propose a Skeleton Aware Multi-modal SLR framework with a Global Ensemble Model (SAM-SLR-v2) to explore the potential of skeleton-based SLR and fuse with other modalities in both RGB and RGB-D scenarios to further improve the recognition rate. Specifically, we design a new spatio-temporal skeleton graph for SLR using whole-body keypoints extracted by a pretrained whole-body pose estimator. Then we propose a multi-stream Sign Language Graph Convolution Network (SL-GCN) to model the embedded dynamics. To fully exploit the information in whole-body keypoints, we propose a novel Separable Spatial-Temporal Convolution Network (SSTCN) to learn from the whole-body skeleton features. Moreover, studies on action recognition reveal that data from different modalities complement each other, provide knowledge of latent correlation, and further improve the final performance. Although we can simply add predictions from all modalities together to achieve higher accuracy, we desire a method that tunes the best weight for each modality in a data-driven way. Hence, we propose a Global Ensemble Model (GEM) to automatically learn the multi-modal ensemble and improve the overall recognition rate. Our main contributions can be summarized as follows:

- We construct novel 2D and 3D skeleton graphs designed for SLR using a pretrained whole-body pose estimator and graph reduction, which requires no extra annotation effort.
- We propose a novel SL-GCN to model dynamics in the skeleton graphs. To our best knowledge, this is the first successful attempt to tackle the SLR task using 2D/3D whole-body skeleton graphs that surpasses RGB-based methods.
- We propose a novel SSTCN to further exploit whole-body skeleton features. It can significantly improve the accuracy compared with the traditional 3D convolution.

- We propose an ensemble model GEM for both RGB and RGB-D based SLR, which learns weights from seven modalities and achieves state-of-the-art performance on three isolated SLR datasets with significant margins.

A preliminary version of this work [25] has been reported in the corresponding workshop of the SLR challenge [26], during which we won the championships in both RGB[1] and RGB-D tracks[2]. Compared with our workshop version, we have made the following improvements: (1) We introduce a new modality Keypoint3D, which considers 3D coordinates in space and deal with occlusions. It improves the overall recognition rate of RGB-D ensemble. (2) We propose a new learning-based late-fusion ensemble method named GEM for multi-modal ensemble, which achieves higher recognition rates than our previous version and saves efforts on weights tuning. (3) The test labels of the AUTSL dataset has been released, so we update our performance from the validation set to the test set. (4) Beyond the challenge dataset (AUTSL), we report our performance on two additional large-scale datasets for isolated (*i.e.*, SLR500 and WLASL2000) compared with the recent state-of-the-art methods. Our approach significantly surpasses their performance with notable margins. (5) We update our figures, provide more details of models, analyze ensemble sensitivity, and discuss challenging cases, which may inspire future research on SLR.

## II. RELATED WORK

**Sign Language Recognition (SLR)** achieves significant progress in obtaining high recognition accuracy in recent years due to the development of practical deep learning architectures and the surge of computational power [15], [10], [13], [27]. Researchers have been modeling the spatio-temporal information in the videos using 2D CNN with RNN or 3D CNN. A 3D-convolutional SLR network associated with attention modules is proposed in [15] to learn the spatio-temporal features from raw videos. CNN, Feature Pooling Module, and LSTM Networks associated with adaptive weights are utilized in [10] to obtain distinctive representations. Besides, researchers have been extending isolated SLR to weakly supervised continuous SLR and Sign Language Translation (SLT) by incorporating sequence learning methods [28], [29], [14], [30], [31]. Recently, multi-cue methods using upper-body poses, hand keypoints, and mouth features have been developed to improve the recognition rate. For example, hand pose priors are introduced in hand-aware frameworks for further performance improvement in [17], [16]. A spatial-temporal multi-cue network for continuous SLR and SLT is proposed in [18]. However, these methods are still not effective enough to capture the complete motion information for robust sign language recognition.

**Skeleton-based Action Recognition** mainly focuses on learning dynamic patterns from the motion of the human skeleton and effectively recognizing human action and activities [32], [33], [34], [35], [36], [37]. In the meanwhile, aiming for a higher recognition accuracy, it can collaborate with other

---

[1]RGB: https://chalearnlap.cvc.uab.es/challenge/43/track/41/result/

[2]RGB-D: https://chalearnlap.cvc.uab.es/challenge/43/track/42/result/

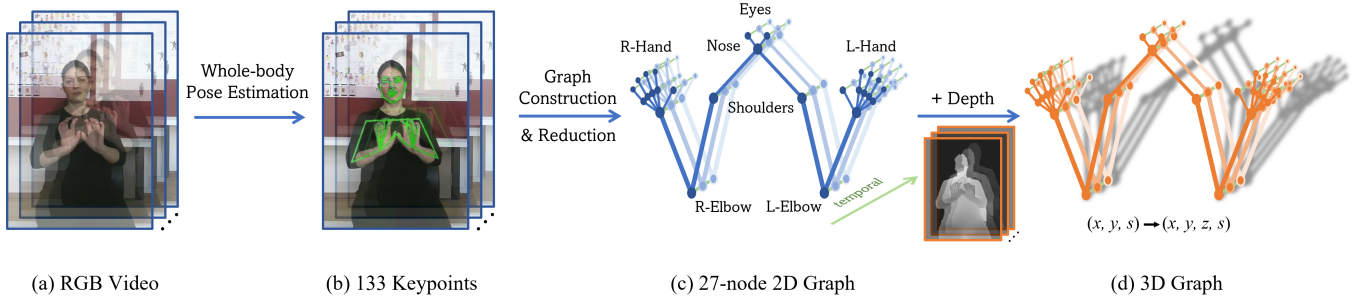(a) RGB Video     (b) 133 Keypoints     (c) 27-node 2D Graph     (d) 3D Graph

Fig. 2. Construction of 2D and 3D graphs. A pretrained whole-body pose estimator is applied on an input RGB video (a) to obtain 133-point whole-body keypoints (b). We then construct a spatio-temporal 2D graph following the natural connections of the human body. The graph is reduced to 27-node (c) using graph reduction to mitigate noises. Each node in the 2D graph is represented by $(x, y, s)$ where $x$-$y$ are 2D coordinates and $s$ is the confidence score. The 3D graph is obtained by overlaying the 2D coordinates on the depth video, reading the depth encoding $z$ at the corresponding location, and treating it as an additional dimension as $(x, y, z, s)$.

modalities (*e.g.*, RGB, depth, and EMG) and benefit from multi-modal learning [38], [39], [40], [41], [42]. Recurrent neural networks (*e.g.*, RNN and LSTM) are once popular in modeling the temporal information of skeleton data [33], [37], [36], [43]. Recently, a graph-based approach that implements a Graph Convolutional Network (GCN) to model the dynamic patterns in skeleton data is first explored by ST-GCN [19]. It draws much research attention that a few improved models are developed to further improve the performance [35], [44], [45], [46], [21], [22], [23]. Specifically, AS-GCN [35] digs the latent joint connections to boost the recognition performance. A two-stream approach is presented in [45] and further extended to four streams in [21]. DecoupleGCN [22] increases the GCN capacity while introducing no extra computational cost. Inspired by ResNet [47], ResGCN [23] incorporates a bottleneck architecture to boost model capacity, reduce parameters, and avoid gradient vanishing. However, skeleton-based SLR methods are still under exploration. Simply applying the ST-GCN to SLR has been unsuccessful, which only reaches around 60% top-1 accuracy on 20 classes (much worse than RGB-based approaches) [48].

**Multi-modal Approach** aims to explore data captured from different resources, by different devices, and from distinctive views to improve the overall performance. Its motivation lies in an assumption that different modalities contain unique and view-specific information that complements each other. The assumption suggests that multi-modal information can be fused together to further boost performance. Frameworks that learn robust and view-invariant feature representation for downstream tasks are proposed in [49], [50]. A weight-sharing model is developed in [51] to obtain modality hallucination for multi-modal image classification. DA-Net [52] adopts a view-independent module collaborated with a view-specific module to capture the multi-modal information and effectively merges the prediction scores. PM-GANs [53] utilizes a cross-view generative strategy associated with a novel fusion to learn the prediction correlations effectively. A feature factorization network is proposed in [54] which learns the specific view-shared information for RGB-D action recognition. A cascaded residual autoencoder is designed to handle incomplete view classification settings [55]. Researchers also propose to use a super vector to fuse the multi-modal representations [56]. En-

couraged by the success of those multi-modal approaches, we aim to incorporate more modalities (*e.g.*, visual, depth, body gesture, and hand gesture) to capture both view-dependent and view-specific information from all aspects. We aim to learn via a universal framework to achieve higher performance.

## III. METHODOLOGY

This section first gives an overview of the proposed SAM-SLR-v2 framework. Then we introduce the SL-GCN for keypoint graphs and the SSTCN for skeleton features, respectively. After that, we present an effective 3DCNN baseline model for the other modalities. Last, we describe the late-fusion GEM, for the multi-modal ensemble.

### A. SAM-SLR-v2 Framework Overview

An overview of the proposed SAM-SLR-v2 is shown in Figure 1. Seven modalities processed from RGB and depth videos are considered in the proposed framework. Three different architectures (*i.e.*, SL-GCN, SSTCN, and 3DCNN) are used to extract features from the seven modalities and make predictions of sign language glosses independently. The late-fusion ensemble model (*i.e.*, GEM) takes predictions from all the modalities and outputs final predictions in the RGB and the RGB-D scenarios.

### B. SL-GCN for Skeleton Keypoints

**Graph Construction and Reduction**. For skeleton-based action recognition, researchers rely on ground-truth skeleton keypoints annotated by motion capture systems such as Kinect v2 [57], which unfortunately is not capable of providing hand and finger annotations. Since hand gestures play a crucial role in performing sign language, we use a pose estimator pretrained with whole-body annotations to predict whole-body keypoints, which include 133 landmarks of face, body, hands, and feet. A spatial 2D graph can then be constructed by connecting every pair of adjacent keypoints according to the human-body natural connections. This graph is further extended to a spatio-temporal graph by connecting all the nodes to themselves in the temporal dimension. Mathematically, the
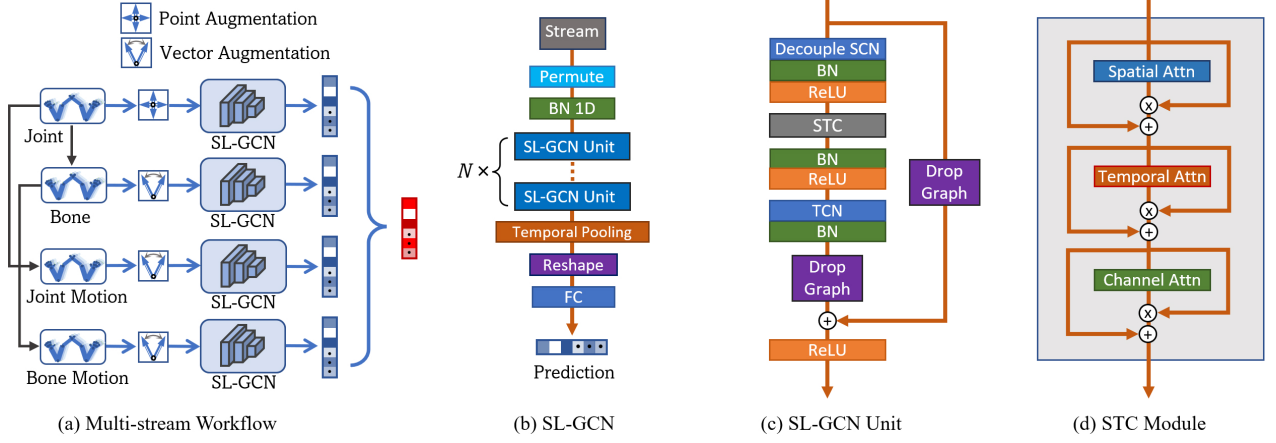
Fig. 3. Illustration of the SL-GCN pipeline: (a) The multi-stream workflow includes streams of joint, bone, joint motion, and bone motion; (b) Illustration of the SL-GCN architecture; (c) Network details of the SL-GCN Unit; (d) The STC self-attention module used in the SL-GCN unit consists of a spatial attention module, a temporal attention module, and a channel attention module that connect in a cascaded way.

node set $V = \{v_{i,t} | i = 1, ..., N, t = 1, ..., T\}$ consists of all 133 whole-body nodes. Their adjacent matrix $\mathbf{A}$ is defined as

$$\mathbf{A}_{i,j} = \begin{cases} 1 & \text{if } d(v_i, v_j) = 1 \\ 0 & \text{else} \end{cases} \quad (1)$$

where $d(v_i, v_j)$ calculate the minimum distance (*i.e.*, the minimum number of nodes) between Node $v_i$ and $v_j$ in the shortest path.

However, contrasting with the graph used in skeleton-based action recognition which contains around 17 nodes, the whole-body skeleton graph contains too many nodes and edges which introduce high-level unexpected noise. Besides, if the distance between two nodes is too far (*i.e.*, have many nodes in between), it is inaccurate to explore their interactions. Our experiment shows that simply using the whole-body skeleton graph results in lower accuracy. Therefore, we operate a graph reduction on the whole-body skeleton graph which trims down the 133 nodes to 27 nodes based on our observations of the videos and visualizations of the GCN activation. The resulted graph consists of seven nodes for the upper body (nose, eyes, shoulders, and elbows) and ten nodes for each hand, as illustrated in Figure 2(c). Graph reduction leads to faster model convergence and significantly higher recognition rates. Each node in the 2D graph is represented by $(x, y, s)$ where $x$-$y$ are 2D coordinates and $s$ is the confidence score. When depth information is available, we construct a 3D graph by reading the corresponding depth $z$ at keypoint locations $x$-$y$ and treating it as an additional dimension as $(x, y, z, s)$, as illustrated in Figure 2(d).

**Graph Convolution**. We adopt the spatio-temporal graph convolution with the spatial partitioning strategy [19] to capture the embedded dynamics in the whole-body skeleton graph. We implement the spatial graph convolution as

$$\mathbf{x}_{\text{out}} = \mathbf{D}^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I})\mathbf{D}^{-\frac{1}{2}}\mathbf{x}_{\text{in}}\mathbf{W}, \quad (2)$$

where $\mathbf{A}$ represents an adjacent matrix of intra-body connections, $\mathbf{I}$ represents an identity matrix of self-connections, $\mathbf{D}$ stands for the diagonal degree of $(\mathbf{A} + \mathbf{I})$, and $\mathbf{W}$ is the trainable weights of the convolution. Practically, Equation 2

is implemented by performing standard 2D convolution and multiplying the results by $\mathbf{D}^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I})\mathbf{D}^{-\frac{1}{2}}$. To perform temporal graph convolution, we implement a standard 2D convolution on the temporal dimension with a kernel size $k_t \times 1$ as the reception field. We adopt an extended variation of the spatial GCN called decoupling GCN [22]. In a decoupling GCN layer, to further boost the model capacity, the extracted features are grouped into $G$ groups that each group has its trainable adjacent matrix $\mathbf{A}$. We then concatenate the outputs of all groups back together as the output features.

**Multi-stream SL-GCN**. Inspired by [21] which adopts a multi-stream workflow for action recognition, we find that it is also worth investigating 1st-order coordinates (joints), 2nd-order vector (bone vector), and their motion vectors for sign language recognition, as shown in Figure 3(a). Following the natural connections of the human body, we generate the bone vectors pointing from the starting joints to their ending joints. Thus the bone stream is represented by a tree graph where the nose acts as the root node. Mathematically, the starting and ending joints can be represented as

$$v_{i,t}^{\text{J}} = (x_{i,t}, y_{i,t}, [z_{i,t}], s_{i,t}), \quad (3)$$

$$v_{j,t}^{\text{J}} = (x_{j,t}, y_{j,t}, [z_{j,t}], s_{j,t}), \quad (4)$$

where $(x, y, [z], s)$ represents 2D coordinates, optional depth, and the confidence score. The bone vectors $v^{\text{B}}$ are calculated by subtracting Equation 3 from Equation 4 as

$$v_{j,t}^{\text{B}} = (x_{j,t} - x_{i,t}, y_{j,t} - y_{i,t}, z_{j,t} - z_{i,t}, s_{j,t}), \forall (i, j) \in \mathbb{H} \quad (5)$$

where the set $\mathbb{H}$ contains all natural connections of the human body. Motion streams are obtained by subtracting the difference between adjacent frames. The joint motion vectors $v^{\text{JM}}$ and the bone motion vectors $v^{\text{JM}}$ are represented as

$$v_{i,t}^{\text{JM}} = v_{i,t+1}^{\text{J}} - v_{i,t}^{\text{J}}, \quad (6)$$

$$v_{i,t}^{\text{BM}} = v_{i,t+1}^{\text{B}} - v_{i,t}^{\text{B}}. \quad (7)$$

We train every stream separately, multiply their predictions by assigned weights, and summing up the results as the final prediction.
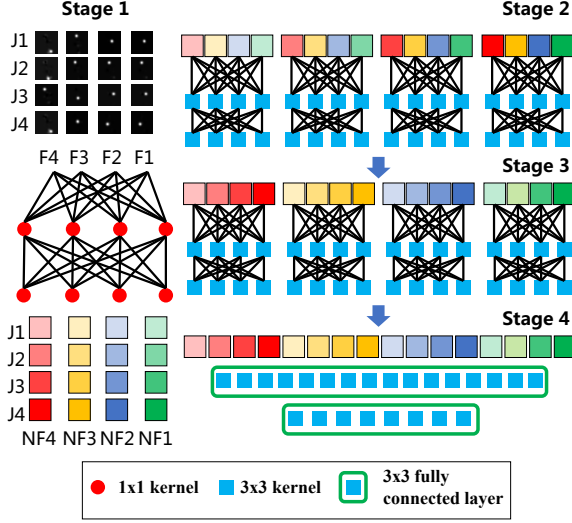
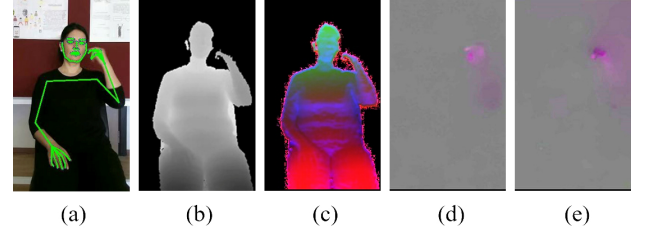Fig. 4. The architecture of SSTCN for skeleton features. Abbrevs: J=Joints; F=Frames; NF=New Features.



Fig. 5. Visualization of modalities: (a) RGB with whole-body keypoints overlay; (b) Depth; (c) Masked HHA; (d) Optical flow; (e) Depth flow. (better viewed in color)

**SL-GCN Structure**. The structure of our proposed SL-GCN is presented in Figure 3(b). The input stream is permuted and normalized before feeding into $N$ instances of SL-GCN units for spatio-temporal graph convolution. An average pooling is then applied to the temporal dimension of the resulted features. The result is reshaped and fed into a fully connected layer for classification. Our proposed SL-GCN unit is illustrated in Figure 3(c). We find that deep graph models are easier to overfit on the video classification task and the ordinary dropout layer works poorly in GCNs. We propose to construct the basic SL-GCN Unit with a decoupled spatial convolutional layer (DecoupleSCN) [22] to mitigate overfitting. We also introduce an STC (spatial, temporal, and channel-wise) self-attention mechanism inspired by [21]. As illustrated in Figure 3(d), the STC module consists of modules of spatial attention, temporal attention, and channel attention connected in a cascaded configuration. In our experiment, we use $N = 10$ such SL-GCN units in the proposed SL-GCN.

### C. SSTCN for Skeleton Features

We propose an SSTCN model to exploit the whole-body features in addition to the keypoint coordinates. We can learn from ResNet2+1D [58] that the performance of an action recognition model can be further improved via factorizing the network into a temporal part and a spatial part. Therefore, our SSTCN model is separated into four stages to handle the features from different dimensions. The pipeline is shown in Figure 4. We save the features of 33 keypoints including 1 nose keypoint, 4 mouth keypoints, 6 upper-body keypoints, and 22 hands keypoints before the argmax operation in the pose estimator. We then uniformly sample 60 frames in each video for SSTCN for recognition. The saved features are resized to $24 \times 24$ using maximum pooling. We process the input features with a 2D separable convolution layer that reduces the parameters and converges easily. At Stage one, features are reshaped from $60 \times 33 \times 24 \times 24$ to $60 \times 792 \times 24$, and fed to $1 \times 1$

convolution layers for temporal convolution. At Stage two, the extracted features are then shuffled and grouped into 60 groups for $3 \times 3$ grouped convolution to extract spatial features of each frame. At the next stage, we start processing the features in the feature dimension. Specifically, the features are shuffled again and grouped into 33 groups. We use grouped $3 \times 3$ convolution on spatial and temporal dimensions for each keypoint. At the last stage, a few fully connected layers are used to make final predictions. We adopt a residual path for the first three stages to avoid gradient vanishing. Besides, random dropouts are deployed in every module to avoid overfitting [59].

Researchers have found that one-hot labels and cross-entropy loss may be easy to overfit with limited data [60]. So we use the technique of label smoothing to avoid overfitting. Mathematically, label smoothing is defined as

$$q'(k|x) = (1 - \epsilon)\delta_{k,y} + \epsilon u(k), \tag{8}$$

where $k$ is the number of classes, $q'(k|x)$ is a label-smoothed predicted distribution, $\epsilon$ is a hyper-parameter in $(0, 1)$, and $u(\cdot)$ represents the uniform distribution. The cross-entropy loss is then modified as

$$H(q', p) = -\sum_{k=1}^{K} \log p(k)q'(k) = (1 - \epsilon)H(q, p) + \epsilon H(u, p), \tag{9}$$

where $H(\cdot)$ is the cross-entropy function and $q$ is the real data distribution, and $p$ stands for the predicted distribution. The modified cross-entropy can be explained as the penalty to the difference between the predicted distribution and the real distribution combined with the difference between the predicted distribution and a prior distribution (*e.g.*, uniform distribution). To further improve the performance, we replace all activations to the Swish [61] activation function as

$$f(x) = x \cdot \text{Sigmoid}(x). \tag{10}$$

### D. 3DCNN Baselines for the Other Modalities

As mentioned in Section II, studies reveal that ensembles from multiple cues and modalities could further improve the overall performance. To benefit from the other modalities (*i.e.*, RGB frames, optical flows, HHA, and depth flow), we hence build a simple yet effective baseline using 3D CNNs. Most 3D CNN architectures are easy to overfit, especially on smaller datasets. In our experiment, The ResNet2+1D [58] architecture that separates the temporal and spatial convolution of 3D
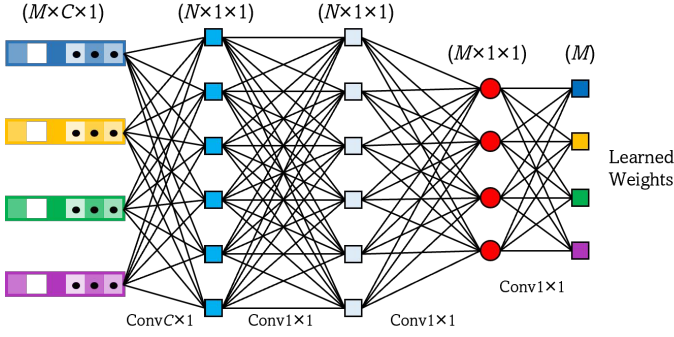
Fig. 6. Illustration of the global ensemble model. The predictions of $M$ modalities are concatenated together, then fed into a convolutional layer with kernel size $C \times 1$, where $C$ is the number of classes. After a few $1 \times 1$ fully connected layers with filter size $N$ and layer normalization, we obtain $M$ weights to fuse the modalities together as the final prediction.

CNNs, provides the best result compared with other popular 3D CNN networks. We also notice that deeper model depth does not always lead to better performance while making the network to overfit easier. So we choose the ResNet2+1D-18 variation pretrained on the Kinectics dataset [62] as our backbone. Apart from that, to improve the accuracy further, we pretrain the model on the largest available SLR dataset SLR500 [63] for RGB frames. Pretraining increases the final accuracy by around 1% and improves the model convergence. Similar to the SSTCN presented in Section III-C, we adopt the Swish activation described in Equation 10 instead of using ReLU. Besides, to mitigate overfitting, we implement the label smoothing technique in Equation 8 with the corresponding cross-entropy loss in Equation 9.

### E. Multi-modal Late-fusion Ensemble

**Model-free Late Fusion**. In our previous version, we use a simple late-fusion approach to fuse predictions from all modalities. To be specific, for all modalities, we save the predictions from the last fully connected layers before the softmax. We assign weights to all modalities manually and add them up as the final prediction

$$q_{\text{RGB}} = \alpha_1 q_{\text{skel2D}} + \alpha_2 q_{\text{RGB}} + \alpha_3 q_{\text{flow}} + \alpha_4 q_{\text{feat}}, \quad (11)$$

$$q_{\text{RGB-D}} = \alpha_1 q_{\text{skel3D}} + \alpha_2 q_{\text{RGB}} + \alpha_3 q_{\text{flow}} + \alpha_4 q_{\text{feat}} + \alpha_5 q_{\text{HHA}} + \alpha_6 q_{\text{depthflow}}, \quad (12)$$

where $q$ represents the predictions before softmax, $\alpha_{1,2,3,4,5,6}$ are hyper-parameters to be tuned based on validation accuracy. In our previous version, we use $\boldsymbol{\alpha} = \{1, 0.9, 0.4, 0.4\}$ for RGB track and $\boldsymbol{\alpha} = \{1.0, 0.9, 0.4, 0.4, 0.4, 0.1\}$ for RGB-D track. For the other datasets without validation and test splits, we keep those weights unchanged without further hard-tuning.
**Global Ensemble Model**. Since finding the best weights for fusion is time-consuming, we propose a learning-based global ensemble model (GEM) to fuse all the modality automatically,

TABLE I
A STATISTICAL SUMMARY OF SLR DATASETS.
* NO LONGER PUBLICLY AVAILABLE
† NOT RELEASED AT THE MOMENT.

| Datasets | #Signs | #Signers | #Samples | Languages |
|---|---|---|---|---|
| AUTSL [64] | 226 | 43 | 21,083 | Turkish |
| SLR500 [63] | 500 | 50 | 125,000 | Chinese |
| WLASL2000 [11] | 2,000 | 119 | 21,083 | American |
| MS-ASL [65] * | 1,000 | 222 | 25,513 | American |
| BSL-1K [66] † | 1,064 | 40 | 273,000 | British |

as illustrated in Figure 6. The whole process can be described as

$$\{\alpha_i\}_{i=1}^4 = G_{\text{RGB}}([q_{\text{Key2D}}, q_{\text{RGB}}, q_{\text{flow}}, q_{\text{feat}}]),$$
$$\{\alpha_i\}_{i=1}^6 = G_{\text{RGB-D}}([q_{\text{Key3D}}, q_{\text{RGB}}, q_{\text{flow}}, q_{\text{feat}}, q_{\text{HHA}}, q_{\text{depthflow}}]), \quad (13)$$

where $q$ represents the results of each modality, $G$ represents the prediction procedure of GEM, and $[\cdot]$ stands for the concatenation operation. The first layer of the model is used to down-sample the modalities effectively with a global convolution. Therefore, the kernel of the first layer is $C \times 1$. Then we use $1 \times 1$ fully connected layers to predict the final weights of all modalities. The final weights are multiplied with the inputs respectively and the weighted modalities are summed up as the final prediction.

## IV. EXPERIMENTS

In this section, we report the evaluation results of our proposed SAM-SLR-v2 on three major sign language recognition datasets of different sign languages. We show the single-modality performance of our proposed multi-stream SL-GCN, SSTCN, and 3DCNN models, as well as different combinations of their ensembles. We demonstrate the effectiveness of proposed approaches via ablation studies on the components of SL-GCN, SSTCN and 3DCNN networks. We also study the proposed GEM compared with simple model-free late fusion. Last, we show some challenging cases and discuss the limitations of SAM-SLR-v2.

### A. Datasets

**AUTSL Dataset** [64] is a Turkish SLR dataset collected using Kinect V2 sensor [57], [67]. Statistically, 226 different sign glosses are performed by 43 signers with 20 backgrounds. The dataset contains 38,336 videos that split into training, validation, and testing subsets. We use their currently released balanced test set in our experiments and report both top-1 and top-5 recognition rates.
**SLR500 Dataset** [63] is a balanced Chinese sign language dataset for isolated sign language recognition (sometimes referred to as CSL isol.) that contains 500 words performed by 50 signers. Each word is performed by all 50 signers 5 times, so there are 125,000 videos in total. This dataset is captured in a controlled lab environment with a solid-color background. The former 36 signers are used for training and the later 14 signers are used for testing.
**WLASL2000 Dataset** [11] is a American Sign Language with a vocabulary size of 2000 words. It is a challenging dataset

TABLE II
PERFORMANCE OF MULTI-STREAM SL-GCN.

| Streams | AUTSL | | SLR500 | | WLASL2000 | |
|---|---|---|---|---|---|---|
| | Top-1 | Top-5 | Top-1 | Top-5 | Top-1 | Top-5 |
| Joint | 95.35 | 99.49 | 97.90 | 99.92 | 45.61 | 77.79 |
| Bone | 95.69 | 99.55 | 97.93 | 99.88 | 43.27 | 75.58 |
| Joint Motion | 93.21 | 99.12 | 97.04 | 99.80 | 27.23 | 56.73 |
| Bone Motion | 93.29 | 99.31 | 97.24 | 99.81 | 31.26 | 60.35 |
| Multi-stream | **96.47** | **99.76** | **98.16** | **99.95** | **51.50** | **84.94** |

TABLE III
ABLATION STUDIES ON SL-GCN ON AUTSL VALIDATION SET.

| Variations | Top-1 |
|---|---|
| SL-GCN (Joint) | **95.02** |
| w/o Graph Reduction | 63.69 |
| w/o Decouple GCN | 94.66 |
| w/o Drop Graph | 94.81 |
| w/o Keypoints Augmentation | 90.16 |
| w/o STC Attention | 93.53 |

TABLE IV
RESULTS OF SINGLE MODALITIES ON AUTSL TEST SET.

| Modality | Top-1 | Top-5 |
|---|---|---|
| Keypoints 2D | 96.47 | 99.76 |
| Keypoints 3D | **96.53** | **99.81** |
| Features | 93.37 | 99.20 |
| RGB Frames | 95.00 | 99.47 |
| RGB Flow | 90.41 | 98.74 |
| Depth HHA | 93.75 | 99.28 |
| Depth Flow | 90.78 | 98.50 |

collected from web videos performed by 119 signers. It contains 21,083 samples with unconstrained recording conditions. The dataset is imbalanced and the average samples per video are much lower than the above two datasets.

We follow the signer-independent settings for all the above datasets. Besides, we are also aware of other large-scale isolated sign language datasets. Microsoft American Sign Language dataset (MS-ASL) [65] is an American sign language dataset that is no longer publicly available due to the deletion of the YouTube-hosted videos. A new large-scale British Sign Language dataset (BSL-1K) [66] has not been released to the public yet at the moment of writing. A statistical summary of the mentioned datasets is provided in Table I.

### B. Multi-modal Data Preparation

**Whole-body Pose Keypoints and Features**. MMPose [68] provides a whole-body pose estimator pretrained with whole-body keypoint annotations and HRNet [69] as its backbone. We use the pretrained model to obtain 133 keypoints from RGB videos, based on which we construct the 27-node 2D graph in Section III-B and process the 3D graphs using depth videos. Keypoint coordinates are normalized to [-1,1]. Random sampling, mirroring, rotating, scaling, jittering, and shifting are applied as data augmentations. We use a sample duration of 150 frames. We repeat videos with lesser than 150 frames to 150 frames. To obtain skeleton features for each video, we downsample the estimated heatmaps, choose 33 joint channels, and uniformly sample 60 frames.

**RGB Frames and Optical Flow**. We extract all frames from RGB videos to load and process them faster in parallel. Besides, we use the Denseflow API implemented with OpenCV and CUDA that provided by OpenMMLab [70] to obtain the optical flows using the TVL1 algorithm [71]. The output x and y flow maps are concatenated in the channel dimension. We cropped all RGB and optical flow frames to the bounding boxes from the pose estimator, and then resize them to $256 \times 256$. When we train the model, we sample 32 consecutive frames randomly from the input video. When testing, we sample five such clips uniformly from the input videos and average on their predicted scores.

**Depth HHA and Depth Flow**. HHA stands for the horizontal disparity, height above the ground, and angle normal that encode a depth map into a three-channel RGB image. Because HHA features enable better scene understanding, we extract HHA features from depth videos as another modality instead of inputting gray-scale depth maps directly. The depth videos of the AUTSL dataset come with masks, so we mask out those regions and fill them with zeros when generating HHA. An

HHA example is shown in Figure 5(c), where black regions are masked out pixels. Cropping and resizing operations are performed the same way as RGB frames. We also apply the same data augmentations to HHA features as the RGB modality. Apart from that, we use Denseflow to extract optical flow from the depth modality as well. A sample of depth flow is shown in Figure 5(e). Compared with optical flows (e.g., Figure 5(d)), depth flows contain less noise and capture distinctive motion information.

### C. Performance of Multi-stream SL-GCN

Both top-1 and top-5 recognition rates of the proposed multi-stream SL-GCN are reported in Table II. Among the four streams, the joint stream provides the best accuracy. The multi-stream approach further improves the overall accuracy, which demonstrates that our proposed whole-body skeleton graph and multi-stream SL-GCN are very effective. Table IV shows that SL-GCN (Keypoints 2D and 3D) perform the best in all single-modality methods. Since the skeleton graphs are less complicated, the graph-based method is lighter-weight and much faster to inference compared with large-capacity 3DCNN-based models, which is another advantage of the proposed SL-GCN.

Table III presents ablation studies on the proposed SL-GCN. We find that the graph reduction technique contributes the most to the recognition rate. Without that, it can barely learn from the noisy dynamics in the 133-node skeleton graph. Due to the limitation of annotated data, the model is easy to overfit, thus the data augmentation techniques are also important in learning the embedded dynamics. Besides, we find that the DropGraph module, the decoupling GCN module, and the STC attention mechanism all contribute to the final performance.

### D. Multi-modal Performance on AUTSL Dataset

The results of all single-modal methods on the AUTSL balanced test set are reported in Table IV. The Keypoints 2D and Keypoints 3D method represents our proposed multi-stream SL-GCN using 2D and 3D skeleton graphs. They

TABLE V
MULTI-MODAL ENSEMBLE RESULTS EVALUATED ON AUTSL TEST SET.
ABBREVS: K2=KEYPOINTS 2D; K3=KEYPOINTS 3D; F=FEATURES;
R=RGB; O=OPTICAL FLOW; H=HHA; D=DEPTH FLOW.

| Ensemble | K2 | K3 | F | R | O | H | D | Top-1 | Top-5 |
|---|---|---|---|---|---|---|---|---|---|
| Skel 2D | ✓ | | ✓ | | | | | 96.90 | 99.84 |
| RGB+Flow | | | ✓ | ✓ | ✓ | | | 96.10 | 99.73 |
| RGB All | ✓ | | ✓ | ✓ | ✓ | | | **97.62** | **100** |
| Skel 3D | | ✓ | ✓ | | | | | 97.01 | 99.87 |
| Depth | | | | | | ✓ | ✓ | 95.38 | 99.57 |
| RGB&D | | | | ✓ | ✓ | ✓ | ✓ | 96.73 | 99.73 |
| RGBD All | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | **98.02** | **100** |

TABLE VI
PERFORMANCE OUR ENSEMBLE RESULTS (WITH AND WITHOUT
FINE-TUNING USING VALIDATION SET) EVALUATED ON AUTSL TEST SET.

| Model | Fine-tune | RGB | | RGB-D | |
|---|---|---|---|---|---|
| | | Top-1 | Top-5 | Top-1 | Top-5 |
| Baseline [64] | - | 49.22 | 75.78 | 62.02 | 83.45 |
| VTN-PF [72] | w/ val | 92.92 | - | 93.32 | - |
| wenbinwuee | w/ val | 96.55 | - | 96.55 | - |
| USTC-SLR | w/ val | 97.62 | - | 97.65 | - |
| SAM-SLR [25] | No | 97.51 | **100** | 97.68 | **100** |
| SAM-SLR-v2 | No | **98.00** | **100** | **98.10** | **100** |
| w/o Keypoint 3D | No | - | - | 98.02 | 99.95 |
| + Extra data | w/ val | **98.42** | **100** | **98.53** | **100** |

TABLE VII
MULTI-MODAL PERFORMANCE ON SLR500 DATASET.

| Modality | Top-1 | Top-5 |
|---|---|---|
| GLE-Net [73] | 96.80 | - |
| Hand + RGB [16] | 98.30 | - |
| Keypoints | 98.16 | **99.95** |
| Features | 97.34 | 99.80 |
| RGB Frames | **98.26** | 99.84 |
| RGB Flow | 95.94 | 99.63 |
| Key + Feat | 98.56 | 99.96 |
| Key + RGB | **98.98** | **99.97** |
| RGB + Flow | 98.45 | 99.88 |
| SAM-SLR [25] | 98.98 | 99.98 |
| SAM-SLR-v2 | **99.00** | 99.98 |

TABLE VIII
MULTI-MODAL PERFORMANCE ON WLASL2000 DATASET.
* PRETRAINED WITH EXTRA BSL-1K DATASET WHICH IS NOT PUBLICLY
AVAILABLE AT THE MOMENT.

| Modality | per-instance | | per-class | |
|---|---|---|---|---|
| | Top-1 | Top-5 | Top-1 | Top-5 |
| Baseline [11] | 32.48 | 57.31 | - | - |
| Fusion-3 [17] | 38.84 | 67.58 | - | - |
| I3D [66] * | 46.82 | 79.36 | 44.72 | 78.47 |
| Hand + RGB [16] | 51.39 | 86.34 | 48.75 | 85.74 |
| Keypoints | **51.50** | **84.94** | **48.87** | **84.02** |
| Features | 46.84 | 79.63 | 44.41 | 78.35 |
| RGB Frames | 47.51 | 80.31 | 44.53 | 78.93 |
| RGB Flow | 40.46 | 73.23 | 37.88 | 71.86 |
| Key + Feat | 55.03 | 89.90 | 52.48 | 87.03 |
| Key + RGB | **57.55** | **90.34** | **54.83** | **89.75** |
| RGB + Flow | 53.53 | 86.50 | 50.63 | 85.77 |
| SAM-SLR [25] | 58.73 | 91.46 | 55.93 | **90.94** |
| SAM-SLR-v2 | **59.39** | **91.48** | **56.63** | 90.89 |

perform the best compared with the other methods. We also find that the depth flow achieves a little better recognition rate compared with the RGB flow resulted from lesser noisy data. The fused ensembles in both RGB and RGB-D scenarios using different choices of modalities are summarized in Table V as three groups. The skeleton-based ensemble (Skel 2D and Skel 3D) achieves better accuracy than "RGB + Flow" and Depth ensemble (HHA + Depth Flow), which demonstrates that the proposed skeleton-based methods are very effective in SLR. The RGB and RGB-D ensemble results that use all available modalities show that the skeleton-based methods also complement RGB and depth modalities. Their collaboration further improves the overall performance. Comparisons with other players in the challenge are reported in Table VI. Our improved SAM-SLR-v2 achieves better recognition rates compared with its previous version SAM-SLR that won the championships of the SLR challenge. During the challenge, we are allowed to use the validation labels to fine-tune our models. Our results with fine-tuning on the validation set are shown at the bottom. Our results without fine-tuning still surpass the other methods fine-tuned with validation labels. Our method achieves the state-of-the-art with significant margins in both RGB and RGB-D scenarios.

*E. Multi-modal Performance on SLR500 Dataset*

Performance on the isolated Chinese sign language dataset (SLR500) is reported in Table VII. Since this dataset is collected in a lab-controlled environment with a large number of repetitions and invariant background, all modalities obtain higher recognition rates compared with the other two datasets. For the same reason, compared with skeleton-based methods, RGB-based 3DCNN achieves a higher Top-1 recognition rate but a lower Top-5 recognition rate. Since the SLR500 dataset

provides only a validation set instead of train-val-test splits, we use the weights learned from the AUTSL dataset to obtain the final ensemble results. Compared with the recent state-of-the-art method GLE-Net [73] and Hand+RGB [16], our two-modality methods and all-modality SAM-SLR-v2 all achieve much better recognition rates.

*F. Multi-modal Performance on WLASL2000*

The WLASL2000 dataset is the most challenging variation of the WLASL dataset with unconstrained backgrounds and camera conditions. Since the dataset is not balanced, we report both per-instance accuracy as well as per-class accuracy following [65], [16]. As reported in Table VIII, the keypoints modality performs the best (51.50% per-instance) among all the single-modal methods due to its independence on backgrounds and resistance to noises. The keypoints and RGB modalities complement each other and boost the two-modal performance. Using all the available modalities, the top-1 performance is improved to 57.55% per-instance. Our proposed late ensemble model is capable of raising the recognition rates by around 0.7%. Compared with other state-of-the-art methods, our proposed SAM-SLR-v2 achieves the best recognition rates in both per-instance and per-class metrics with large leading gaps of around 8%.

*G. Ablation Study on SSTCN and 3DCNN*

An ablation study on SSTCN is provided in Table IX. Compared with ResNet3D [74], [75] and ResNet2+1D [58] on
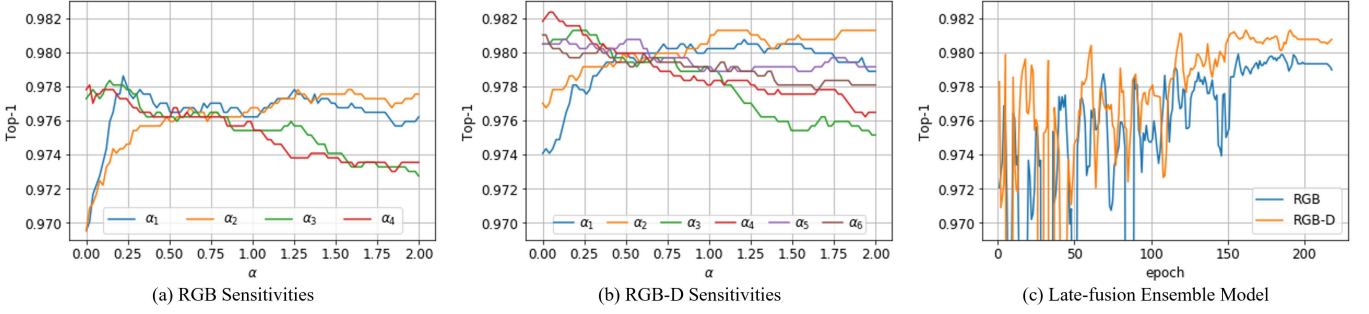
Fig. 7. Sensitivity analysis of ensemble models evaluated using AUTSL test set. (a) Sensitivities of four RGB ensemble parameters ($\alpha_1 \ldots \alpha_4$); (b) Sensitivities of on six RGB-D ensemble parameters ($\alpha_1 \ldots \alpha_6$); (c) Learning the ensemble weights using the proposed late-fusion GEM (both RGB and RGB-D).

TABLE IX
COMPARING SSTCN WITH RESNET3D AND RESNET2+1D AND ABLATION STUDY ON FEATURE SIZES EVALUATED ON AUTSL VALIDATION SET.

| Methods | Feature size | Top-1 |
|---|---|---|
| ResNet3D | $12 \times 12$ | 92.82 |
| ResNet2+1D | $12 \times 12$ | 93.03 |
| SSTCN | $12 \times 12$ | 93.60 |
| SSTCN | $24 \times 24$ | **94.32** |

TABLE X
AN ABLATION STUDY OF 3D CNN ARCHITECTURE USING RGB FRAMES EVALUATED ON AUTSL VALIDATION SET.

| 3D CNN Variations | Top-1 |
|---|---|
| Ours (RGB Frame) | **94.77** |
| w/o Label Smoothing | 93.75 |
| w/o Swish Activation | 92.88 |
| w/o Pretraining on CSL | 93.41 |
| w/ ResNet3D-18 Backbone | 93.10 |

the same feature size, SSTCN is more effective on recognizing sign glosses using skeleton features. Our study also reveals that the proposed SSTCN can achieve even higher accuracy with larger input feature sizes.

Table X shows an ablation study of 3D CNN architecture using the RGB frames. It shows that both swish activation and label smoothing are effective, where they improve the top-1 accuracy by 2% and 1%, respectively. Pretraining on the CSL dataset [63] improves the overall accuracy by 1.4%. The ResNet2+1D-18 backbone provides 1.7% better performance than ResNet3D-18.

### H. Comparison between Multi-modal Ensembles

The model-free simple ensemble method requires a lot of effort in tuning the weights of different modalities. With our proposed late-fusion GEM, those weights are automatically learned from the data. We study the effectiveness of GEM compared with the simple ensemble method. In Table VI, VII, and VIII, performance of model-free simple ensemble is shown as SAM-SLR. Compared with the simple ensemble, our proposed late-fusion GEM can improve the overall performance as well as avoid hard tuning on modality weights.

### I. Sensitivity Analysis of Ensemble Methods

A sensitivity study on ensemble methods is conducted on the AUTSL test dataset. In the experiment, we use the manually tuned parameters in Section III-E as our basis and change one parameter at a time (varying from 0.0 to 2.0), while keeping the other parameters fixed. We plot the resulted top-1 accuracy to analyze the sensitivity of those ensemble parameters. In the RGB scenario, we have four modalities and we analyze the sensitivity of $\alpha_1$ to $\alpha_4$, as shown in Figure 7(a). In the RGB-D scenario, we analyze the sensitivity of $\alpha_1$ to $\alpha_6$, as shown in Figure 7(b).

In this experiment, the ensemble parameters tuned on the validation set do not provide the highest performance on the test set. The ensemble performance is sensitive to the values of the parameters that small changes may lead to large variations in the ensemble accuracy. Moreover, when there are more modalities (*e.g.*, in the RGB-D scenario), it is hard to optimize all those parameters manually. Our proposed late-fusion GEM solves this problem by automatically learning these parameters from data. Figure 7(c) shows the top-1 accuracy during training the proposed model. Initially, the ensemble model gives sensitive results due to its large search space. The model converges after 100 epochs and delivers steady and accurate ensemble results. In summary, the proposed late-fusion GEM is very effective in learning the multi-modal ensemble, gives robust higher performance than simple ensemble, and saves a lot of effort to tune the best parameters.

### J. Challenging Cases and Model Limitations

Figure 8 shows some challenging cases of sign language recognition from the AUTSL dataset. The offline full-body pose estimator may fail due to off-screen or occlusion, especially for fingers, which is a common issue for pose estimation methods. However, fingers play a critical role in expressing signs. Some of those failures can be corrected by RGB-based features. That is the reason why multi-modal ensemble significantly boosts the recognition rate. Moreover, a same sign may be performed very differently by different signers (*e.g.*, signers may use their left hands, right hands, or both hands to perform the same gloss). Hence, mirror augmentation is very important in model training, but more data collected from more signers are desired. Last, opposite signs could be visually similar in a single frame (*e.g.*, heavy vs light-weight). Distinguishing between those signs requires delicate modeling of spatio-temporal dynamics. In our experiment, we find that

(a) Keypoints errors      (b) Single-hand vs Dual-hand

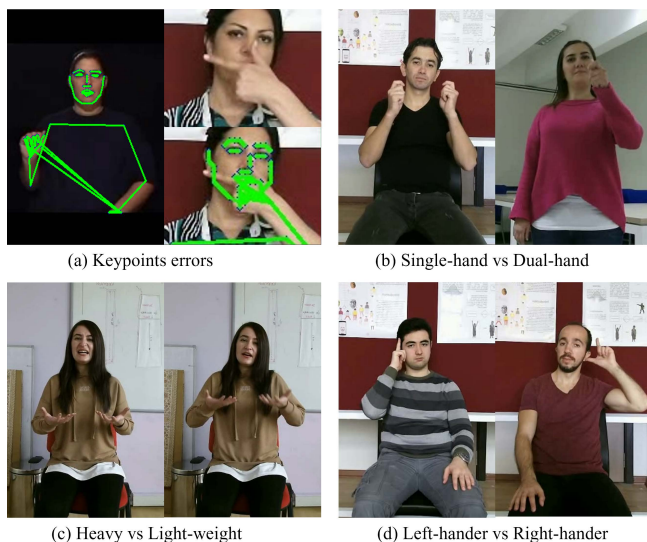(c) Heavy vs Light-weight      (d) Left-hander vs Right-hander

Fig. 8. Examples of challenging cases: (a) Errors in whole-body keypoint estimation; (b) A same sign performed differently by different signers; (c) Two opposite signs can be visually similar; (d) Left-handers and right-handers perform mirror-symmetrically.

skeleton-based methods are smarter choices over RGB-based methods.

## V. CONCLUSION

In conclusion, we propose a novel SAM-SLR-v2 framework to learn multi-modal feature representations from RGB-D videos towards more effective and robust isolated SLR. Among those modalities, our proposed skeleton-based methods are the most effective in modeling motion dynamics due to their signer-independent and background-independent characteristics. Specifically, we construct novel 2D and 3D spatio-temporal skeleton graphs using pretrained whole-body keypoint estimators and propose a multi-stream SL-GCN to model the embedded motion dynamics. Our method does not require any additional effort on annotating hands and provides more reliable hand keypoint estimation than off-line hand detectors. Besides modeling motion dynamics of keypoint coordinates via graphs, we propose SSTCN to predict using skeleton features. Furthermore, we study the multi-modal fusion problem based on the other modalities (*i.e.*, RGB frames, optical flow, HHA, and depth flow) via a learning-based late-fusion ensemble model named GEM. Experimentally, we show that our proposed SAM-SLR-v2 framework achieves the state-of-the-art performance on three challenging datasets for isolated SLR (*i.e.*, AUTSL, SLR500, and WLASL2000) as well as won the championships in both RGB and RGB-D tracks during the CVPR 2021 challenge on isolated SLR. We hope our work could facilitate and inspire future research on SLR.

## REFERENCES

[1] K. Emmorey, *Language, cognition, and the brain: Insights from sign language research*. Psychology Press, 2001. 1

[2] C. Valli and C. Lucas, *Linguistics of American sign language: An introduction*. Gallaudet University Press, 2000. 1

[3] T. Johnston and A. Schembri, *Australian Sign Language (Auslan): An introduction to sign language linguistics*. Cambridge University Press, 2007. 1

[4] Q. Yang, "Chinese sign language recognition based on video sequence appearance modeling," in *Proceedings of IEEE Conference on Industrial Electronics and Applications*, 2010, pp. 1537–1542. 1

[5] A. Mindess, *Reading between the signs: Intercultural communication for sign language interpreters*. Nicholas Brealey, 2014. 1

[6] Q. Zhu, M.-C. Yeh, K.-T. Cheng, and S. Avidan, "Fast human detection using a cascade of histograms of oriented gradients," in *Proceedings of IEEE Computer Vision and Pattern Recognition*, vol. 2, 2006, pp. 1491–1498. 1

[7] D. G. Lowe, "Object recognition from local scale-invariant features," in *Proceedings of IEEE International Conference on Computer Vision*, vol. 2, 1999, pp. 1150–1157. 1

[8] N. H. Dardas and N. D. Georganas, "Real-time hand gesture detection and recognition using bag-of-features and support vector machine techniques," *IEEE Transactions on Instrumentation and measurement*, vol. 60, no. 11, pp. 3592–3607, 2011. 1

[9] A. Memiş and S. Albayrak, "A Kinect based sign language recognition system using spatio-temporal features," in *Proceedings of International Conference on Machine Vision*, vol. 9067. International Society for Optics and Photonics, 2013, p. 90670X. 1

[10] O. M. Sincan, A. O. Tur, and H. Y. Keles, "Isolated sign language recognition with multi-scale features using LSTM," in *2019 27th Signal Processing and Communications Applications Conference (SIU)*. IEEE, 2019, pp. 1–4. 1, 2

[11] D. Li, C. Rodriguez, X. Yu, and H. Li, "Word-level deep sign language recognition from video: A new large-scale dataset and methods comparison," in *Proceedings of IEEE Winter Conference on Applications of Computer Vision*, 2020, pp. 1459–1469. 1, 6, 8

[12] L. Pigou, A. Van Den Oord, S. Dieleman, M. Van Herreweghe, and J. Dambre, "Beyond temporal pooling: Recurrence and temporal convolutions for gesture recognition in video," *International Journal of Computer Vision*, vol. 126, no. 2, pp. 430–439, 2018. 1

[13] A. O. Tur and H. Y. Keles, "Isolated sign recognition with a siamese neural network of RGB and depth streams," in *IEEE International Conference on Smart Technologies*, 2019, pp. 1–6. 1, 2

[14] J. Huang, W. Zhou, Q. Zhang, H. Li, and W. Li, "Video-based sign language recognition without temporal segmentation," in *Proceedings of AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018. 1, 2

[15] J. Huang, W. Zhou, H. Li, and W. Li, "Attention-based 3D-CNNs for large-vocabulary sign language recognition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 9, pp. 2822–2832, 2018. 1, 2

[16] H. Hu, W. Zhou, and H. Li, "Hand-model-aware sign language recognition," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, 2021, pp. 1558–1566. 1, 2, 8

[17] A. A. Hosain, P. S. Santhalingam, P. Pathak, H. Rangwala, and J. Kosecka, "Hand pose guided 3d pooling for word-level sign language recognition," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 3429–3439. 1, 2, 8

[18] H. Zhou, W. Zhou, Y. Zhou, and H. Li, "Spatial-temporal multi-cue network for sign language recognition and translation," *IEEE Transactions on Multimedia*, 2021. 1, 2

[19] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *Proceedings of AAAI conference on Artificial Intelligence*, vol. 32, no. 1, 2018. 2, 3, 4

[20] H. Wang and L. Wang, "Modeling temporal dynamics and spatial configurations of actions using two-stream recurrent neural networks," in *Proceedings of IEEE Computer Vision and Pattern Recognition*, 2017. 2

[21] L. Shi, Y. Zhang, J. Cheng, and H. Lu, "Skeleton-based action recognition with multi-stream adaptive graph convolutional networks," *IEEE Transactions on Image Processing*, vol. 29, pp. 9532–9545, 2020. 2, 3, 4, 5

[22] K. Cheng, Y. Zhang, C. Cao, L. Shi, J. Cheng, and H. Lu, "Decoupling GCN with DropGraph module for skeleton-based action recognition," in *Proceedings of European Conference on Computer Vision*, 2020. 2, 3, 4, 5

[23] Y.-F. Song, Z. Zhang, C. Shan, and L. Wang, "Stronger, faster and more explainable: A graph convolutional baseline for skeleton-based action recognition," in *Proceedings of ACM International Conference on Multimedia*, 2020, pp. 1625–1633. 2, 3

[24] Q. Xiao, M. Qin, and Y. Yin, "Skeleton-based chinese sign language recognition and generation for bidirectional communication between

deaf and hearing people," *Neural Networks*, vol. 125, pp. 41–55, 2020. 2

[25] S. Jiang, B. Sun, L. Wang, Y. Bai, K. Li, and Y. Fu, "Skeleton aware multi-modal sign language recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3413–3423. 2, 8

[26] O. M. Sincan, J. C. S. Jacques Junior, S. Escalera, and H. Y. Keles, "Chalearn LAP large scale signer independent isolated sign language recognition challenge: Design, results and future research," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2021. 2

[27] K. M. Lim, A. W. C. Tan, C. P. Lee, and S. C. Tan, "Isolated sign language recognition using convolutional neural network hand modelling and hand energy image," *Multimedia Tools and Applications*, vol. 78, no. 14, pp. 19 917–19 944, 2019. 2

[28] O. Koller, S. Zargaran, H. Ney, and R. Bowden, "Deep sign: Enabling robust statistical continuous sign language recognition via hybrid CNN-HMMs," *International Journal of Computer Vision*, vol. 126, no. 12, pp. 1311–1325, 2018. 2

[29] D. Guo, W. Zhou, H. Li, and M. Wang, "Hierarchical LSTM for sign language translation," in *Proceedings of AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018. 2

[30] J. Pu, W. Zhou, and H. Li, "Iterative alignment network for continuous sign language recognition," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4165–4174. 2

[31] R. Cui, H. Liu, and C. Zhang, "A deep neural framework for continuous sign language recognition by iterative training," *IEEE Transactions on Multimedia*, vol. 21, no. 7, pp. 1880–1891, 2019. 2

[32] J. Cai, N. Jiang, X. Han, K. Jia, and J. Lu, "Jolo-gcn: Mining joint-centered light-weight information for skeleton-based action recognition," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 2735–2744. 2

[33] Y. Du, W. Wang, and L. Wang, "Hierarchical recurrent neural network for skeleton based action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1110–1118. 2, 3

[34] Z. Huang, C. Wan, T. Probst, and L. Van Gool, "Deep learning on lie groups for skeleton-based action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6099–6108. 2

[35] M. Li, S. Chen, X. Chen, Y. Zhang, Y. Wang, and Q. Tian, "Actional-structural graph convolutional networks for skeleton-based action recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3595–3603. 2, 3

[36] S. Li, W. Li, C. Cook, C. Zhu, and Y. Gao, "Independently recurrent neural network (indrnn): Building a longer and deeper rnn," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5457–5466. 2, 3

[37] J. Liu, A. Shahroudy, D. Xu, and G. Wang, "Spatio-temporal lstm with trust gates for 3d human action recognition," in *European conference on computer vision*. Springer, 2016, pp. 816–833. 2, 3

[38] F. Baradel, C. Wolf, and J. Mille, "Human action recognition: Pose-based attention draws focus to hands," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2017, pp. 604–613. 3

[39] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308. 3

[40] V. Choutas, P. Weinzaepfel, J. Revaud, and C. Schmid, "Potion: Pose motion representation for action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7024–7033. 3

[41] J.-F. Hu, W.-S. Zheng, J. Pan, J. Lai, and J. Zhang, "Deep bilinear learning for rgb-d action recognition," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 335–351. 3

[42] M. Zolfaghari, G. L. Oliveira, N. Sedaghat, and T. Brox, "Chained multi-stream networks exploiting pose, motion, and appearance for action classification and detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2904–2913. 3

[43] C. Si, Y. Jing, W. Wang, L. Wang, and T. Tan, "Skeleton-based action recognition with spatial reasoning and temporal stack learning," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 103–118. 3

[44] L. Shi, Y. Zhang, J. Cheng, and H. Lu, "Skeleton-based action recognition with directed graph neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, 7912–7921. 3

[45] ——, "Two-stream adaptive graph convolutional networks for skeleton-based action recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 026–12 035. 3

[46] C. Si, W. Chen, W. Wang, L. Wang, and T. Tan, "An attention enhanced graph convolutional lstm network for skeleton-based action recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1227–1236. 3

[47] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778. 3

[48] C. C. de Amorim, D. Macêdo, and C. Zanchettin, "Spatial-temporal graph convolutional networks for sign language recognition," in *International Conference on Artificial Neural Networks*. Springer, 2019, pp. 646–657. 3

[49] J. Zheng, Z. Jiang, and R. Chellappa, "Cross-view action recognition via transferable dictionary learning," *IEEE Transactions on Image Processing*, vol. 25, no. 6, pp. 2542–2556, 2016. 3

[50] H. Wang, A. Kläser, C. Schmid, and C.-L. Liu, "Dense trajectories and motion boundary descriptors for action recognition," *International Journal of Computer Vision*, vol. 103, no. 1, pp. 60–79, 2013. 3

[51] J. Hoffman, S. Gupta, and T. Darrell, "Learning with side information through modality hallucination," in *Proceedings of IEEE Computer Vision and Pattern Recognition*, 2016, pp. 826–834. 3

[52] D. Wang, W. Ouyang, W. Li, and D. Xu, "Dividing and aggregating network for multi-view action recognition," in *Proceedings of European Conference on Computer Vision*, September 2018. 3

[53] L. Wang, C. Gao, L. Yang, Y. Zhao, W. Zuo, and D. Meng, "PM-GANs: Discriminative representation learning for action recognition using partial-modalities," in *Proceedings of European Conference on Computer Vision*, 2018, pp. 384–401. 3

[54] A. Shahroudy, T.-T. Ng, Y. Gong, and G. Wang, "Deep multimodal feature analysis for action recognition in RGB+D videos," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 5, pp. 1045–1058, 2018. 3

[55] L. Tran, X. Liu, J. Zhou, and R. Jin, "Missing modalities imputation via cascaded residual autoencoder," in *Proceedings of IEEE Computer Vision and Pattern Recognition*, 2017, pp. 1405–1414. 3

[56] Z. Cai, L. Wang, X. Peng, and Y. Qiao, "Multi-view super vector for action recognition," in *Proceedings of IEEE Computer Vision and Pattern Recognition*, 2014, pp. 596–603. 3

[57] D. Pagliari and L. Pinto, "Calibration of kinect for Xbox One and comparison between the two generations of microsoft sensors," vol. 15, pp. 27 569–27 589, 10 2015. 3, 6

[58] D. Tran, H. Wang, L. Torresani, J. Ray, Y. LeCun, and M. Paluri, "A closer look at spatiotemporal convolutions for action recognition," in *Proceedings of IEEE Computer Vision and Pattern Recognition*, 2018, pp. 6450–6459. 5, 8

[59] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014. 5

[60] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *arXiv preprint arXiv:1506.01497*, 2015. 5

[61] P. Ramachandran, B. Zoph, and Q. V. Le, "Searching for activation functions," *arXiv preprint arXiv:1710.05941*, 2017. 5

[62] J. Carreira, E. Noland, A. Banki-Horvath, C. Hillier, and A. Zisserman, "A short note about kinetics-600," *arXiv preprint arXiv:1808.01340*, 2018. 6

[63] J. Zhang, W. Zhou, C. Xie, J. Pu, and H. Li, "Chinese sign language recognition with adaptive HMM," in *Proceedings of IEEE International Conference on Multimedia and Expo*, 2016, pp. 1–6. 6, 9

[64] O. M. Sincan and H. Y. Keles, "AUTSL: A large scale multi-modal turkish sign language dataset and baseline methods," *IEEE Access*, vol. 8, pp. 181 340–181 355, 2020. 6, 8

[65] H. Vaezi Joze and O. Koller, "Ms-asl: A large-scale data set and benchmark for understanding american sign language," in *The British Machine Vision Conference*, 2019. 6, 7, 8

[66] S. Albanie, G. Varol, L. Momeni, T. Afouras, J. S. Chung, N. Fox, and A. Zisserman, "Bsl-1k: Scaling up co-articulated sign language recognition using mouthing cues," in *European Conference on Computer Vision*, 2020, pp. 35–53. 6, 7, 8

[67] C. Amon, F. Fuhrmann, and F. Graf, "Evaluation of the spatial resolution accuracy of the face tracking system for Kinect for windows V1 and V2," in *Proceedings of AAAI Conference on Artificial Intelligence*, 2014, pp. 16–17. 6

[68] M. Contributors, "OpenMMLab pose estimation toolbox and benchmark," https://github.com/open-mmlab/mmpose, 2020. 7

[69] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *Proceedings of IEEE Computer Vision and Pattern Recognition*, 2019, pp. 5693–5703. 7

[70] S. Wang, Z. Li, Y. Zhao, Y. Xiong, L. Wang, and D. Lin, "denseflow," https://github.com/open-mmlab/denseflow, 2020. 7

[71] C. Zach, T. Pock, and H. Bischof, "A duality based approach for real-time TV-L1 optical flow," in *Proceedings of Joint Pattern Recognition Symposium.* Springer, 2007, pp. 214–223. 7

[72] M. De Coster, M. Van Herreweghe, and J. Dambre, "Isolated sign recognition from rgb video using pose flow and self-attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3441–3450. 8

[73] H. Hu, W. gang Zhou, J. Pu, and H. Li, "Global-local enhancement network for nmf-aware sign language recognition," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 17, pp. 1 – 19, 2020. 8

[74] K. Hara, H. Kataoka, and Y. Satoh, "Can spatiotemporal 3D CNNs retrace the history of 2D CNNs and ImageNet?" in *Proceedings of IEEE Computer Vision and Pattern Recognition*, 2018, pp. 6546–6555. 8

[75] H. Kataoka, T. Wakamiya, K. Hara, and Y. Satoh, "Would mega-scale datasets further enhance spatiotemporal 3D CNNs?" *arXiv preprint arXiv:2004.04968*, 2020. 8