

FineHand: Learning Hand Shapes for American Sign Language Recognition

Al Amin Hosain, Panneer Selvam Santhalingam, Parth Pathak, Huzefa Rangwala and Jana Košecká
Department of Computer Science, George Mason University, Fairfax, Virginia, USA

Abstract—American Sign Language recognition is a difficult gesture recognition problem, characterized by fast, highly articulate gestures. These are comprised of arm movements with different hand shapes, facial expression and head movements. Among these components, hand shape is the vital, often the most discriminative part of a gesture. In this work, we present an approach for effective learning of hand shape embeddings, which are discriminative for ASL gestures. For hand shape recognition our method uses a mix of manually labelled hand shapes and high confidence predictions to train deep convolutional neural network (CNN). The sequential gesture component is captured by recursive neural network (RNN) trained on the embeddings learned in the first stage. We will demonstrate that higher quality hand shape models can significantly improve the accuracy of final video gesture classification in challenging conditions with variety of speakers, different illumination and significant motion blur. We compare our model to alternative approaches exploiting different modalities and representations of the data and show improved video gesture recognition accuracy on GMU-ASL51 benchmark dataset.

I. INTRODUCTION

Despite numerous efforts in the past addressing different components of American Sign Language (ASL) gesture recognition, automated sign language parsing in the wild remains challenging. The presented work is motivated by the need to design interfaces that can enable interactions between a Deaf and Hard-of-Hearing (DHH) user and a digital assistant (e.g. Amazon Echo, Google Now). Intelligent virtual assistant devices are becoming ubiquitous and the types of services they offer continues to expand. They can help users in answering questions, managing schedules, describing weather and many more. However, most of these devices are voice controlled. Hence, DHH people are deprived of the benefits from using these assistants.

An ASL sign is performed by a combination of hand gestures, facial expressions and postures of the body. Further, the sequential motion of specific body locations (such as hand-tip, neck and arm) provide informative cues about a sign.

In this work we consider the problem of classification of individual ASL signs captured by short video snippets in unrestricted settings, by multiple signers. We focus on learning effective hand shape representations robust to changes of style, motion blur and illumination (see Figure 1). To localize the hands, we exploit recent advances in human pose estimation methods from video, which are effective in estimating 2D hand joint locations [2], [34], [35]. In order

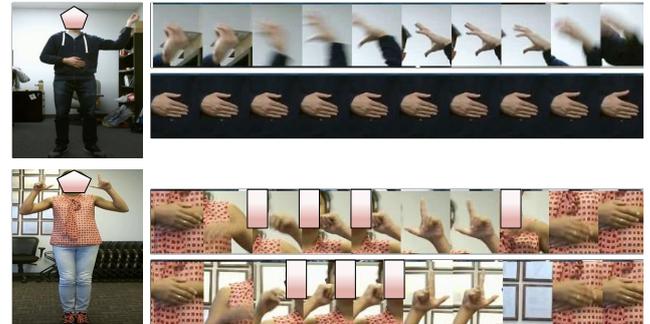


Fig. 1. Subjects performing different ASL signs and corresponding hand-shape pattern. Faces are masked for privacy concern (Note: All faces are masked in this paper for the same reason).

train discriminative hand shape model, frame level annotation of hand shapes classes is required. Such labelling tedious and time consuming process. For this reason, we first manually select a small set of training hand shape examples from each video gesture to train the initial model and keep increasing this collection using the predictive power of deep convolutional neural networks (CNNs). We will show that, this type of supervision can benefit gesture recognition accuracy while keeping the manual annotation effort at minimum. By training a network for hand shape classification, we obtain discriminative hand shape embeddings as penultimate layer of the network. These are subsequently used to learn sequential dynamics of video gestures video using recursive neural network (RNN) [16]. In summary, the contributions of this work are:

- We propose an iterative learning mechanism to train a deep CNN to learn robust hand shape embeddings.
- We implement sequential RNN models for video gesture recognition using the learned hand shape embeddings.
- We show evaluation of our method varying different factors such as fraction of hand-shape supervision and single hand vs both hands and compare it with other multi-modal methods for ASL gesture recognition.
- We create and release per frame hand-shape annotations for GMU-ASL51 dataset. This will pave the way for rigorous sequential machine learning on the dataset.

We will demonstrate superior performance of the proposed model on the GMU-ASL51 dataset [8] and compare it quantitatively with several baseline approaches which use different representations and different sensing modalities.

II. RELATED WORK

Previous work on sign language recognition focused either on video gestures, finger spelling or sentence parsing and deployed variety of sensing modalities. With the exception of few many benchmarks were acquired in controlled laboratory settings, with small variation in speakers, clothing and lighting or used speaker augmentation to make the hand detection, tracking and recognition more robust.

In [37] authors developed HMM based framework for ASL phrase verification and subjects used colored gloves during data collection. The follow up work used Kinect skeletal data and accelerometers worn on hands [36]. Starner [29] demonstrated two HMM based approaches to recognize sentence level ASL using a single camera to track the user's unadorned hands, but used carefully chosen clothing and backgrounds for hand detection and tracking purposes.

More recently, approaches based on deep neural network have attracted much attention in modeling sign language and enabled to bypass the cumbersome feature engineering stage. With the advent of deep learning techniques, several approaches for ASL gesture recognition have been developed for learning discriminative spatio-temporal features from video [9], [12], [30] and applied to finger spelling recognition. Huang [9] showed the effectiveness of using Convolutional neural network (CNN) with RGB video data for sign language recognition. Three dimensional CNN has been used to extract spatio-temporal features from the video in [12]. Similar architecture was implemented for Italian gestures [25]. Sun *et al.* [30] hypothesized that not all RGB frames in a video are equally important and assigned a binary latent variable to each frame in training videos for indicating the importance of a frame within a latent support vector machine model. Zaki *et al.* [38] proposed two new features with existing hand crafted features and developed the system using HMM based approach. Some have used appearance-based features and divided the approach into sub-units of RGB and tracking data, with a HMM model for recognition [3]. These methods either tackled few number of sign class variation with uniform background, trained models using a fraction of test subject's data or did not explore the hand shapes rigorously.

Compared to RGB methods, skeletal data has received little attention in ASL recognition. However, in a closely related human action recognition task, a significant amount of work has been done using body joint information. Shahroudy [27] released the largest dataset for human activity recognition. They proposed an extension of long short term memory (LSTM) model which leverages group motion of several body joints to recognize human activity from skeletal data. A different adaptation of the LSTM model was proposed by Liu [18] where spatial interaction among joints was considered in addition to the temporal dynamics. Some researches focused on capturing salient motion pattern of body joints [31] and some leveraged hierarchical properties of joint configuration [4]. Several attention based model were proposed for human activity analysis [19], [28]. Few prior

works converted skeleton sequences of body joints or RGB videos into an image representation and then applied state-of-art image recognition models to achieve good results [13], [17]. In generic activities the whole body moves, which is not the case for ASL sign gestures where primarily the hands move. Estimated body joint poses (2D/3D) can be used to a certain level, because pose data only gives a high level motion pattern of a sign gesture. Hence, sign gestures recognition demands careful hand shape modeling.

Hand segmentation or recognition is also a well studied problem in computer vision [1], [15], [21]. Some of these methods concentrate on hand detection rather than modeling hand shapes [21], some has different viewpoints such as egocentric views [1]. Koller *et al.* [15] trained a CNN for hand shape modeling in an semi supervised manner. This method is closer to a part of our work, however, it lacks fine grained hand shape modeling. On the other hand, our goal is to learn robust hand shape representation using careful supervision.

We use only RGB modality in sign language classification. This makes our system independent of depth sensor. However, unlike traditional RGB based methods which use feature based HMM, we base our model on deep learning methods considering the size and variation in the dataset. Besides, our method focuses on learning fine grained hand-shape features before the sign classification phase. It uses fine supervision from a small fraction of data to learn hand shape and further use sequential modeling using learned representation for final classification task. We show by leveraging careful supervision on hand shapes our methods can achieve significant performance boost.

III. GMU-ASL51 DATASET

All of our hand shape learning as well as ASL classification tasks were evaluated using GMU-ASL51 benchmark [8]. We picked this dataset because it is the only publicly available dataset of this type (isolated word level ASL gestures) with large number of sign variation. GMU-ASL51 has 51 word level ASL signs performed by 12 subjects of different ages, gender and builds. The dataset was collected using depth sensor and has two modalities: RGB videos and 3D skeletal body parts. More detail can be found in the paper. In this dataset, only video level class label is available. One of our main contributions of this work is to systematically annotate per frame hand shape from each video.

IV. OUR APPROACH

We refer to our proposed sign gesture classification pipeline as FineHand. It has two parts. In this section, we start with describing the first part which is a CNN model trained to learn hand shape representation. Hereafter, we refer this model as hand shape model (or embedder). Then, we present the second part which is a sequential RNN classification model learn to recognize different sign gestures using hand shape embeddings.

A. Hand Shape Embeddings

The goal of this part of our pipeline is to learn high-dimensional representation of hand shape which is discriminative for ASL gesture recognition. In this section we are first going to present how we crop hand patches using off the shelf pose estimation method, followed by iterative hand shape learning mechanism. Finally, we present some qualitative results from our learned hand shape model.

a) *Pose Estimation*: Pose estimation is the process of estimating 2D or 3D body joint locations (e.g. wrist, elbow) in single image. Typically it is done by first detecting human subjects in an image frame and then parsing body joint location [5], [11], [23]. or the inferring the body parts without first detecting the person as a whole [2], [10], [26]. For this

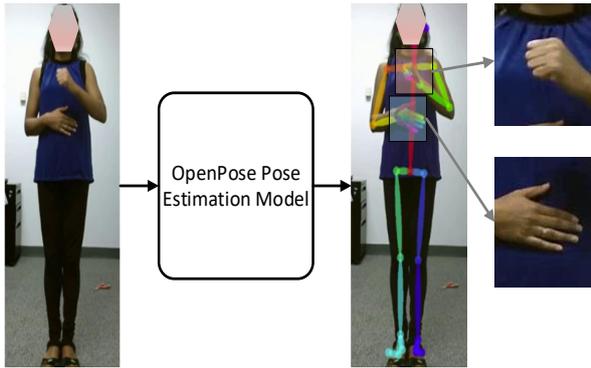


Fig. 2. Pose estimation and hand cropping process.

work we have chosen the state-of-the-art 2D human body pose estimation approach OpenPose [2]. It is to be noted that we only use hand poses to crop a hand patch from each image. Figure 2 shows the whole process of estimating body poses from an RGB frame and cropping hand patches.

TABLE I

ITERATIVE HAND-SHAPE LEARNING PROCESS. IN THE HEADER ROW P, C AND T SYMBOLIZES PREDICTION, CORRECT AND TOTAL COUNT RESPECTIVELY. ITER 1 IS THE MANUAL ANNOTATION, HENCE ALL LABELS ARE CORRECT. T COLUMN OF FINAL PASS DENOTES THE CUMULATIVE COUNT OF HAND-SHAPE SAMPLES FOR THE CLASS REPRESENTED BY ROWS. ITERATION IS ABBREVIATED AS ITER.

Class	Iter 1	Iter 2			Iter 3		
		P	C	T	P	C	T
C1	402	598	534	936	524	511	1447
C2	217	277	277	494	281	277	771
C3	69	88	73	142	102	96	238
C4	328	554	408	735	435	396	1131
C5	163	236	196	359	219	198	557

1) *Iterative Hand Shape Learning*: GMU-ASL51 dataset has 12 subjects and 51 word level sign classes with only sign video level labels. In this phase, we learn per frame hand representation by training hand shape embedder. Our approach depends on some early manual hand shape annotation. To be more specific, We take one gesture sample per subject for each sign class which gives us a total of 612 (12×51) sign videos. We extract hand-shape patches using

pose data and manually group them based on visual similarity into 41 classes. These examples are used to fine-tune ResNet-50 CNN architecture [6]. In the second iteration, we use the high-confidence predictions of our initial model to predict hand shape patches from another 612 sign videos (different from the 612 set of first iteration). At this point we will have some incorrect predictions because the model is trained on small amount of data. We manually correct the incorrect predictions. Although this fix is manual but we can see that significantly less amount of labor is needed in the second iteration than the first one. After this phase, we have more annotated hand-shape patches and we retrain our model. With additional iterations our model becomes more robust. Similarly we can start a third iteration and so on. We perform



Fig. 3. Sample hand-shapes. Each row shows 12 randomly picked samples from the created dataset in the iterative hand shape learning phase.

3 such iterations and end up with 41 classes of hand shapes which are distributed among 51 sign gesture classes of GMU-ASL51 dataset. Table I shows the count of annotated hand-shape samples for five classes and three iterations. Details for all 41 classes will be provided in the supplemental materials. It should be noted here, in three iterations the fraction of sign gesture data we use to train the model were 4.16%, 8.32% and 12.5% respectively. It should be also mentioned that, we use per frame hand-shape annotation on this fraction of data which means that a video gesture sample could possibly generate several hand-shape training examples. Effects of this incremental learning on sign classification is discussed in more detail in the result section. Among 41 hand-shapes two are unusual: garbage and rest-position. Keeping those two classes is significant because, in a sign video most of the frames hands are blurred or are in a resting position. If we have a way to learn these uninformative hand shapes then we can exclude them during sign language modeling and hence have robust feature representation. Figure 3 shows examples of several hand shape classes picked from 41 class hand shape dataset. Each row represents one class. Twelve

samples in a row are picked randomly from our created hand shape dataset. We observe significant intra-class variation. Bottom two rows show the sample hand shapes from the class `garbage` and `rest-position` respectively. This ResNet based hand shape model is a module of the proposed FineHand sign gesture classification pipeline. After being trained on hand shape data, parameters of this model can be frozen during sequential learning of signs.

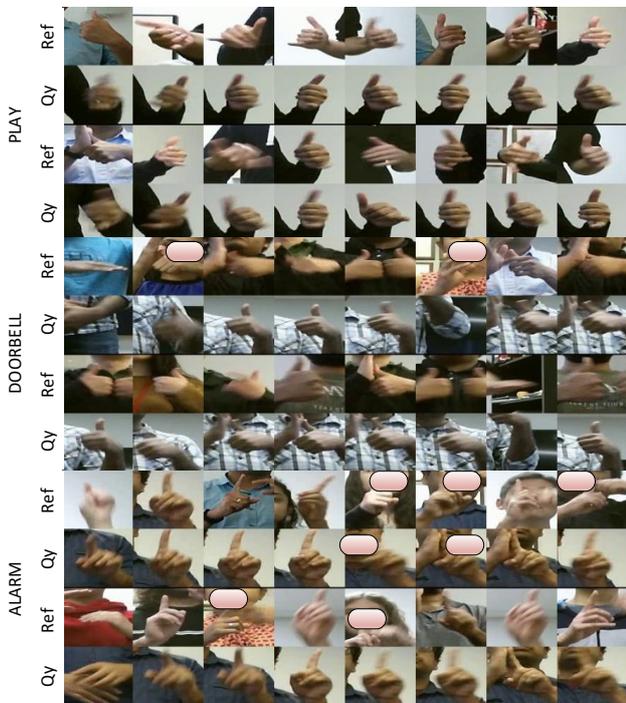


Fig. 4. Predicted hand shape classes from a trained ResNet-50. For each of three sign classes two samples are shown where *Qy* means query hand shape sequence and *Ref* means a reference sample of predicted label for each corresponding query hand patch from the training samples.

a) Qualitative Results: The hand shape model network is trained in an incremental fashion. Manual annotation in the first iteration and all the other generated hand shape labels are hand independent. This means if we look at all the samples from a class, we will find examples from both hands. Of course, a shape will be horizontally flipped or rotated as we look at the left versus right hand. First, second and fourth rows in Figure 3 show such examples. Figure 4 presents predicted and reference examples patches using a trained hand shape model, four top rows show two samples (divided by thin black line) from the sign `play`. For each samples the second row shows the query (*Qy*) hand patch and the corresponding patch in the row above (*Ref*) represents a reference training sample from predicted class. This classification model is trained on cross-subject manner which means none of the hand patches from the test subject (query) is used during training procedure. We observe, in most of the cases, hand shape model learns useful feature representation which is invariant to rotation, scaling and background. We will then keep the penultimate layer

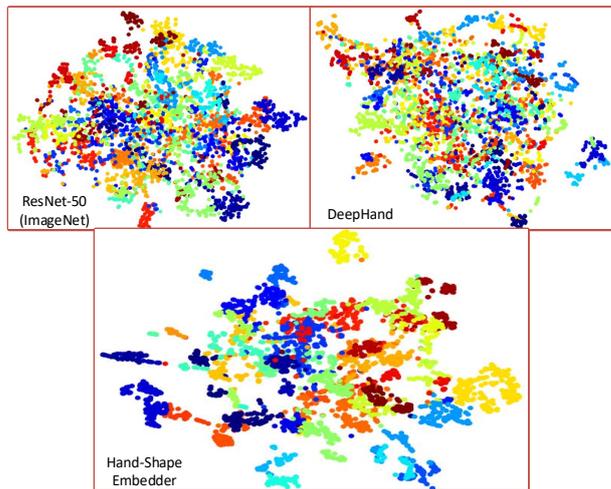


Fig. 5. T-SNE visualization of hand shape embeddings from different model sources for 4033 samples. Top figures show representation obtained from ResNet-50 (ImageNet pre-trained) and DeepHand respectively, from left to right. Bottom figure shows embedding obtained from a trained (cross subject) ResNet model using hand shape labels created by us.

of the model to be high-dimensional embedding each hand patch.

Figure 5 shows the comparison among T-SNE representation of embeddings obtained by our model. It shows that hand shape representation learned by our CNN embedder cluster better than embeddings produced by ResNet (ImageNet trained) [6] or DeepHand [15] models. This gives us the motivation behind using this effective representation in sign video classification task.

B. Sequential Sign Gesture Learning

As pointed out and shown in the previous section, learned hand shape representation could be useful in classifying ASL gesture videos. However, sequential dynamics in video data still needs to be modeled carefully for better classification accuracy. With a trained hand shape embedder, each video can be converted into a sequence of embeddings. We base our sequential modeling using recurrent neural network (RNN).

a) Recurrent Neural Network: For modeling sequential data, recurrent neural network (RNN) have yielded impressive results on a variety of sequence prediction tasks [16]. RNN models can capture temporal dynamics by maintaining an internal state. However, the basic RNN has problems dealing with long term dependencies in data due to the vanishing gradient problem. Some solutions to the vanishing gradient problem involve careful initialization of network parameters or early stopping [24]. But the most effective solution is to modify the RNN architecture in such a way that it has a memory state (cell state) at every time step that can identify what to remember and what to forget. This architecture is referred as long short term memory (LSTM) network. While the basic RNN is a direct transformation of the previous state and the current input, the LSTM maintains

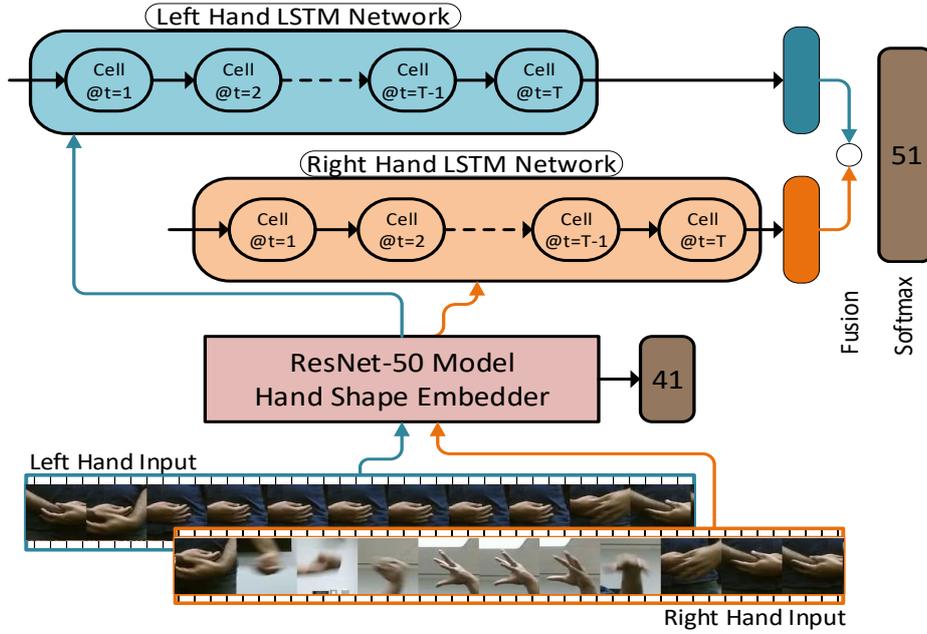


Fig. 6. FineHand RNN model. The ResNet-50 model is trained separately first on 41 hand-shape classes. After training, it provides representation for each hand-patch video which is then used in sequential LSTM classifier for 51 sign classes in the dataset.

an internal memory and has a mechanism to update and use that memory. This is achieved by deploying four separate neural networks also called gates. More detailed description of LSTM model can be found in [7].

1) *FineHand Architecture*: We propose to use an LSTM recurrent neural network for this task. For a sign video, input to this model is a sequence of embeddings obtained from our ResNet based hand shape embedder model. Assume we have a sequence cropped hands images $R^{F \times H \times W}$ where F, H, W are number of frames, height and width of cropped hand patches respectively. Feeding the images to the hand shape model we obtain $D = 2048$ dimensional embedding for each frame and output of the form $R^{F \times D}$. We use this data to train an LSTM model where F is the number of temporal steps. For simplicity to deal with different number of frames, we sample T predetermined number of frames uniformly. Hence, input to the LSTM network is $R^{T \times D}$. Finally we pick the hidden state at the end of last layer of LSTM network and take that as a encoded representation for each sequence. This final representation captures rich temporal as well as spatial hand shape features and is fed to a fully connected neural network layer to produce prediction probability distribution for sign gesture video classification. We train two different networks for left and right hands. We fuse the output of two networks at the end to produce final classification scores. Figure 6 shows the details of our proposed architecture. It shows that the left-hand and right-hand patches are input to the trained hand-shape embedder model and the generated representation is being used as input for the recurrent LSTM networks.

2) *Training Details*: Average cross entropy loss is used to update the network parameters for a batch. Given a true one

hot encoded class label of $y_{i,j}$ and corresponding predicted score of $\hat{y}_{i,j}$ this loss is calculated as Equation 1.

$$L = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C y_{i,j} \log(\hat{y}_{i,j}) \quad (1)$$

Here, N is the size of the minibatch and C is the number of classes. For our experiments, these values are 64 and 51 respectively. For selecting the number of layer of this network, hidden state size and time steps of input, we performed grid search. We found best validation accuracy is obtained with 2 layer networks, 512 as hidden state size and 20 (T) as the input sequence length. Size of the input dimension at each time steps is the size (2048) of the representation produced by hand shape embedder. We used Adam Optimizer for training our networks [14] with learning rate set to 0.0001.

V. EXPERIMENTS

In this section, we are going to describe different methods we used in comparison with our proposed architecture.

A. Comparative Methods

a) *3D Convolution*: Convolutional Neural Network (CNN) with 3D convolutional kernel (3D CNN) has shown promising performance in classifying human activities in video [12]. That's why we chose 3D CNN on RGB hand patch videos for one of our baselines. It consists of four 3D convolutional layers and two fully connected layers at the end. There are two separate networks for left and right hands' patches. Final embedding of these two networks are concatenated before producing softmax score.

TABLE II

TEST ACCURACIES ACROSS 12 SUBJECTS. IN HEADER ROW EACH SUBJECT IS REPRESENTED BY S APPENDED WITH SUBJECT NUMBER. BOTTOM ROW SHOWS THE RESULT OF OUR PROPOSED FINEHAND ARCHITECTURE. OTHER ROWS ARE DIFFERENT COMPARATIVE METHODS. 3D LABELED THREE ROWS SHOW THE RESULTS USING DEEPHAND EMBEDDING WITH KINECT 3D POSE. 2D LABELED ROWS SHOW THE SIMILAR EXPERIMENTS WITH OPENPOSE POSES.

		S01	S02	S03	S04	S05	S06	S07	S08	S09	S10	S11	S12	Average
3D	3D CNN	0.62	0.63	0.60	0.60	0.56	0.56	0.62	0.45	0.39	0.51	0.46	0.17	0.51
	FoaNet [22] Style	0.18	0.22	0.10	0.19	0.11	0.28	0.15	0.17	0.09	0.12	0.11	0.05	0.15
	PoseLSTM	0.81	0.88	0.68	0.84	0.88	0.85	0.85	0.78	0.81	0.83	0.76	0.83	0.82
	RgbLSTM	0.81	0.82	0.85	0.89	0.81	0.89	0.93	0.69	0.83	0.72	0.81	0.37	0.78
2D	FusionLSTM	0.90	0.93	0.88	0.93	0.90	0.95	0.96	0.84	0.91	0.89	0.89	0.67	0.88
	PoseLSTM	0.83	0.90	0.89	0.92	0.90	0.92	0.95	0.88	0.94	0.93	0.94	0.63	0.89
	RgbLSTM	0.80	0.82	0.88	0.93	0.82	0.92	0.95	0.79	0.90	0.75	0.92	0.27	0.81
	FusionLSTM	0.85	0.85	0.93	0.95	0.87	0.92	0.98	0.83	0.93	0.85	0.96	0.40	0.86
FineHand (ours)		0.93	0.98	0.97	0.91	0.91	0.93	0.99	0.88	0.91	0.94	0.96	0.83	0.93

b) Deep Hand Model: The authors in [15] trained a 22 layer deep convolutional neural network (CNN) with more than 1 million images, from videos of Danish and New Zealand sign language. The data is weakly labeled with only video level annotation. The CNN model for estimating the likelihood of hand shapes, is trained using EM algorithm, jointly with Hidden Markov Model (HMM) for parsing sign gestures. The network is trained to recognize 60 hand shape classes plus one garbage class which determines the start or end of a video. We forgo the softmax classification layer of this network and use the embeddings computed by this pre-trained model as representations of crops of hand patches from ASL data. The final layer embedding (1024 dimensional) as a feature vector. We use representation generated by this model in our comparative experiments described next.

c) Kinect 3D Pose: For this experiment we used RGB video and 3D skeletal pose data as computed by Kinect sensor [40]. We do three types of experiments by using these two types of data. Two of those experiments use each modality separately. The other one uses fusion strategy to get maximum out of both modalities. These models are referred as 3D version of PoseLSTM, RgbLSTM and FusionLSTM in result section. We use embedding from Deep Hand model in this comparative method.

d) OpenPose: This experiment is similar to the process described in last paragraph except only uses estimated poses from RGB videos instead of using 3D skeletal poses generated by Kinect sensor. Process of estimating poses from video is described in section IV-A. Rationale behind doing this experiment is to exclude the dependency on depth sensor. Since, this poses are estimated from RGB video, this experiment depends solely on RGB modality. These models are referred as 2D version of PoseLSTM, RgbLSTM and FusionLSTM in result section.

e) Comparison with Similar Work: It was difficult to find a similar work which can be directly compared to our method. This is due to the nature of our set up which is isolated ASL word level sign recognition. We could not find any standard public dataset other than recently released GMU-ASL51 for this purpose. However, there is a public dataset on generic gestures [32]. Various deep learning based methods

have been proposed to model generic isolated gestures [20], [22], [33], [39]. Narayana *et al.* proposed FOANet which uses fusion of different channels on cropped hand patches and full body [22]. Using different modalities (RGB, depth and flow) with those channels this work sparsely fuses 12 channels of inputs to model gestures. Details can be found in the paper. We tried to reproduce this work on GMU-ASL51 as closely as possible. However exact recreation was not possible due to several factors such as number of data channels used, number of location features of hand patches and mechanism of cropping hand patches. Details of these differences will be provided in supplemental materials.

VI. RESULTS

Table II shows our experimental results. All of our experiment shown are cross subject in manner. For a particular test subject, we trained our model using data from all other subjects in the dataset. This cross subject evaluation criteria supports practical usability of our system to subjects unknown to the trained model. Each column in Table II shows test accuracy of one subject in GMU-ASL51. First two rows shows the baseline results: 3D CNN and FoaNet style implementation. Next two blocks of three rows show experiments with 3D and 2D poses respectively. Finally, bottom row shows results of our proposed method (FineHand).

Result shows that methods using 2D poses achieve almost similar performance to 3D Kinect poses (88% vs 86%) which is interesting because 2D pose methods depend only on RGB data. It should be noted that, unlike Kinect sensor, OpenPose provides finger joints. We presume, even though these poses lack depth information, finger joints help to achieve comparable performance with 3D Kinect poses.

From Table II we observe that, our proposed method, FineHand outperforms top models using 3D and 2D poses by 5% and 7% respectively. It should be mentioned that, pose based models use pose data while FineHand model only uses RGB hand patches. Taking this into consideration, it is fair to compare FineHand with RGB only versions (RgbLSTM) of 3D and 2D implementation. In that case, FineHand outperforms those implementations by 15% and 11% respectively. This significant boost in classification accuracy can be justified by the way FineHand learns hand

shapes. Representation used for RgbLSTMs is taken from DeepHand, a pre-trained model on huge amount of hand shapes data from a different class distribution as des V-A. On the other hand, FineHand embedder learns representation from of fine grained hand shapes which has proven to be crucial for this kind of significant performance gain in classification. Our best method outperforms the work came with the dataset [8] by 12%.

The FOANet style implementation on GMU-ASL51 has really bad performance (second row in Table II) even though it is one of the top performing models for generic gestures. One possible reason is the number of data channels used. While original work uses 12 channels, in our implementation we use only 2 channels to make it comparable with our proposed work. Another reason is the training procedure. FOANet architecture proposed to capture sequential dynamics in a video gesture by using sliding window based approach where classification scores for a video gesture was computed by taking averages over all sliding window scores. Our method however, pre samples some fixed number of frames from a video and produces one set of prediction score per video gesture.

TABLE III
AVERAGE CROSS SUBJECT SIGN RECOGNITION ACCURACY ON DIFFERENT ITERATIONS OF HAND SHAPE LEARNING

Iterations (% Train Data)	Accuracy
Iteration 0 (0.00%)	0.65
Iteration 1 (4.17%)	0.89
Iteration 2 (8.32%)	0.91
Iteration 3 (12.5%)	0.93

a) Effect of Hand-shape Learning Iterations: In section IV-A.1, we briefly described how hand-shape learning CNN is trained in successive iterations. We hypothesize, increasing these iterations will boost up the sign classification accuracy. Table VI shows the average recognition accuracy on using different fractions of data for training the embedder CNN. Here, ‘Iteration 0’ represents no hand-shape learning meaning we use embedding representation from ImangeNet pre-trained CNN model as input to LSTM network for sign classification.

TABLE IV
AVERAGE CROSS SUBJECT SIGN RECOGNITION ACCURACY FOR DIFFERENT SCENARIOS OF HAND USAGE.

Input Types	Accuracy
Left Hand	0.65
Right Hand	0.90
Both Hands (max score)	0.86
Both Hands (catenation)	0.92
Both Hands (mean score)	0.93

b) Both Hands vs Single Hand: We are also interested to compare results obtaining from either using left or right hand and using them together. Usually some signs are dominated by single hand while others are double handed. It is obvious that, using both hands’ information will increase the

accuracy. However, we want to see how much improvement is possible using both hands. In case of using both hands, we also show if there is any best fusion mechanism. Table VI-.0.a shows this results. We observe that, good accuracy is achieved using only right hand input. This is because, we notice that, all of the subjects use right hand as dominant hand in GMU-ASL51 dataset. However, using both hands we have 3% improved accuracy which suggests that, in some cases right hand is also important. Among the fusion strategies we found, averaging logits works best.

TABLE V
AVERAGE CROSS SUBJECT SIGN RECOGNITION ACCURACY FOR DIFFERENT LEARNING MECHANISMS (JOINT VS SEPARATE).

Train Type	Accuracy
Separate Learning	0.93
Joint Learning	0.89

c) Effect of Joint Learning: This section shows comparison between learning the whole network jointly vs separately. Our default set up is separate learning where first we train the hand shape CNN model using annotated hand patch data, then freeze it during sequential learning of embedding produced by it on training video hand patches. In case of joint learning, we don’t freeze the hand shape CNN network during sequential sign learning. We notice that, joint learning worsen the performance of the whole network. We presume, since the CNN is first trained on hand-shape supervision, it might get confused when sign level gradient updates are done on its parameters. Also, during back propagation gradient has to travel backwards through the LSTM network before hitting hand-shape CNN model. This causes derogatory updates on CNN parameters. Hence, produced embedding representation differs in each iteration, which impact LSTM learning negatively.

VII. CONCLUSION AND FUTURE WORK

We have demonstrated effectiveness of learning hand-shape representation for ASL sign gesture classification from video. We showed qualitative results of better representation produced by our proposed hand-shape learning mechanism. We also verified that, this representation can achieve superior sign classification accuracy than other sources of embedding learning. Our proposed method is RGB only but outperforms multi-modal (RGB and pose) approaches of sign language recognition. Given the hand shape annotation, exploring sentence level sign language modeling is an interesting direction. We believe, per frame hand shape annotation will help to localize individual sign gestures in a sign video sentence.

REFERENCES

- [1] S. Bambach, S. Lee, D. Crandall, and C. Yu. Lending a hand: Detecting hands and recognizing activities in complex egocentric interactions. In *IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [2] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh. OpenPose: realtime multi-person 2D pose estimation using Part Affinity Fields. In *arXiv preprint arXiv:1812.08008*, 2018.

- [3] H. Cooper, E.-J. Ong, N. Pugeault, and R. Bowden. Sign language recognition using sub-units. *J. Mach. Learn. Res.*, 13(1):2205–2231, July 2012.
- [4] Y. Du, W. Wang, and L. Wang. Hierarchical recurrent neural network for skeleton based action recognition. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1110–1118, June 2015.
- [5] G. Gkioxari, B. Hariharan, R. Girshick, and J. Malik. Using k-poselets for detecting people and localizing their keypoints. In *Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR '14*, pages 3582–3589, Washington, DC, USA, 2014. IEEE Computer Society.
- [6] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, June 2016.
- [7] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, Nov. 1997.
- [8] A. A. Hosain, P. S. Santhalingam, P. Pathak, J. Kosecka, and H. Rangwala. Sign language recognition analysis using multimodal data, 2019.
- [9] J. Huang, W. Zhou, H. Li, and W. Li. Sign language recognition using 3d convolutional neural networks. In *2015 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6, June 2015.
- [10] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele. Deepcut: A deeper, stronger, and faster multi-person pose estimation model, 2016.
- [11] U. Iqbal and J. Gall. Multi-person pose estimation with local joint-to-person associations, 2016.
- [12] S. Ji, W. Xu, M. Yang, and K. Yu. 3d convolutional neural networks for human action recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(1):221–231, Jan. 2013.
- [13] Q. Ke, M. Bennamoun, S. An, F. A. Sohel, and F. Boussaïd. A new representation of skeleton sequences for 3d action recognition. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 4570–4579, 2017.
- [14] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- [15] O. Koller, H. Ney, and R. Bowden. Deep hand: How to train a cnn on 1 million hand images when your data is continuous and weakly labelled. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [16] Z. C. Lipton. A critical review of recurrent neural networks for sequence learning. *CoRR*, abs/1506.00019, 2015.
- [17] J. Liu, N. Akhtar, and A. Mian. Skepxels: Spatio-temporal image representation of human skeleton joints for action recognition. *CoRR*, abs/1711.05941, 2017.
- [18] J. Liu, A. Shahroudy, D. Xu, A. K. Chichung, and G. Wang. Skeleton-based action recognition using spatio-temporal lstm network with trust gates. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2017.
- [19] J. Liu, G. Wang, L. Y. Duan, K. Abdiyeva, and A. C. Kot. Skeleton-based human action recognition with global context-aware attention lstm networks. *IEEE Transactions on Image Processing*, 27(4):1586–1599, April 2018.
- [20] Q. Miao, Y. Li, W. Ouyang, Z. Ma, X. Xu, W. Shi, and X. Cao. Multimodal gesture recognition based on the resc3d network. In *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pages 3047–3055, Oct 2017.
- [21] A. Mittal, A. Zisserman, and P. H. S. Torr. Hand detection using multiple proposals. In *British Machine Vision Conference*, 2011.
- [22] P. Narayana, J. R. Beveridge, and B. A. Draper. Gesture recognition: Focus on the hands. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5235–5244, June 2018.
- [23] G. Papandreou, T. Zhu, N. Kanazawa, A. Toshev, J. Tompson, C. Bregler, and K. Murphy. Towards accurate multi-person pose estimation in the wild, 2017.
- [24] R. Pascanu, T. Mikolov, and Y. Bengio. Understanding the exploding gradient problem. *CoRR*, abs/1211.5063, 2012.
- [25] L. Pigou, S. Dieleman, P.-J. Kindermans, and B. Schrauwen. Sign language recognition using convolutional neural networks. In *Computer Vision - ECCV 2014 Workshops*. Springer International Publishing, 2015.
- [26] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. V. Gehler, and B. Schiele. Deepcut: Joint subset partition and labeling for multi person pose estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [27] A. Shahroudy, J. Liu, T. T. Ng, and G. Wang. Ntu rgb+d: A large scale dataset for 3d human activity analysis. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1010–1019, June 2016.
- [28] S. Song, C. Lan, J. Xing, W. Zeng, and J. Liu. An end-to-end spatio-temporal attention model for human action recognition from skeleton data. In *AAAI Conference on Artificial Intelligence*, pages 4263–4270, 2017.
- [29] T. Starner, A. Pentland, and J. Weaver. Real-time american sign language recognition using desk and wearable computer based video. *IEEE Trans. Pattern Anal. Mach. Intell.*, 20(12):1371–1375, Dec. 1998.
- [30] C. Sun, T. Zhang, and C. Xu. Latent support vector machine modeling for sign language recognition with kinect. *ACM Trans. Intell. Syst. Technol.*, 6(2):20:1–20:20, Mar. 2015.
- [31] V. Veeriah, N. Zhuang, and G. J. Qi. Differential recurrent neural networks for action recognition. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 4041–4049, Dec 2015.
- [32] J. Wan, S. Z. Li, Y. Zhao, S. Zhou, I. Guyon, and S. Escalera. Chalearn looking at people rgb-d isolated and continuous datasets for gesture recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 761–769, June 2016.
- [33] H. Wang, P. Wang, Z. Song, and W. Li. Large-scale multimodal gesture recognition using heterogeneous networks. In *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pages 3129–3137, Oct 2017.
- [34] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh. Convolutional pose machines. In *CVPR*, 2016.
- [35] Y. Xiu, J. Li, H. Wang, Y. Fang, and C. Lu. Pose Flow: Efficient online pose tracking. In *BMVC*, 2018.
- [36] Z. Zafrulla, H. Brashear, T. Starner, H. Hamilton, and P. Presti. American sign language recognition with the kinect. In *Proceedings of the 13th International Conference on Multimodal Interfaces, ICMI '11*, pages 279–286, New York, NY, USA, 2011. ACM.
- [37] Z. Zafrulla, H. Brashear, P. Yin, P. Presti, T. Starner, and H. Hamilton. American sign language phrase verification in an educational game for deaf children, 08 2010.
- [38] M. M. Zaki and S. I. Shaheen. Sign language recognition using a combination of new vision based features. *Pattern Recognition Letters*, 32(4):572 – 577, 2011.
- [39] L. Zhang, G. Zhu, P. Shen, and J. Song. Learning spatiotemporal features using 3dcnn and convolutional lstm for gesture recognition. In *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pages 3120–3128, Oct 2017.
- [40] Z. Zhang. Microsoft kinect sensor and its effect. *IEEE MultiMedia*, 19(2):4–10, Apr. 2012.