

# L<sup>A</sup>T<sub>E</sub>X Author Guidelines for 3DV Proceedings

Anonymous 3DV submission

Paper ID \*\*\*\*

## Abstract

*The ABSTRACT is to be in fully-justified italicized text, at the top of the left-hand column, below the author and affiliation information. Use the word “Abstract” as the title, in 12-point Times, boldface type, centered relative to the column, initially capitalized. The abstract is to be in 10-point, single-spaced type. Leave two blank lines after the Abstract, then begin the main text. Look at previous 3DV abstracts to get a feel for style and length. The ABSTRACT is to be in fully-justified italicized text, at the top of the left-hand column, below the author and affiliation information. Use the word “Abstract” as the title, in 12-point Times, boldface type, centered relative to the column, initially capitalized. The abstract is to be in 10-point, single-spaced type. Leave two blank lines after the Abstract, then begin the main text. Look at previous 3DV abstracts to get a feel for style and length.*

## 1. Introduction

According to National Institute on Deafness and Other Communication Disorders, one in 1000 infants is born totally deaf, while an additional one to six per thousand are born with hearing loss of different levels ???. Sign language is used by Deaf or Hard-of-Hearing (DHH) people to communicate among themselves. In addition to hearing difficulties, DHH people can have hardness in understanding and speaking natural languages spoken by normal hearing people. Because of that, they can not use spoken or written form of natural language to communicate with others. Sign language is the natural language to DHH people. Primary use of sign language is to communicate with other DHH people. It also allows them to develop intellectual abilities and learn social traits which are fundamental attributes any normal human being should possess as a part of society.

American sign language (ASL) is the 3<sup>rd</sup> mostly used language in USA and is used by around half a million of people ???. Sign language can also be used by the DHH people to communicate with the normal hearing people given that both parties know it. But not many normal hearing

people know sign language. That’s why daily life communication with people around becomes challenging for them. For example, requesting someone around to open the door, turning on the thermostat or asking someone about weather condition etc. On the other hand, personal digital assistants (PDA) such as Siri, Cortana, Google Home etc. are becoming popular in making life easier for people. These PDAs can be used to make better quality of life for DHH people. To achieve this goal, we must need an ASL recognizer embedded with PDAs.

The purpose of this work is to build an ASL recognition system. Most of the current system dealing with ASL recognition use RGB video data. ASL sign is performed by moving the hands combined with facial expressions and postures of the body. While video data is good for having a total view of whole body movement, we have noticed that, the motion of specific position of different locations of the hands and head areas are important for ASL recognition. In other word, ASL sign can be treated as a sequence of motion of some specific body parts. Using, video data it is not practical to single out each different body location and motion sequences associated with them from sequences of raw RGB value. Microsoft kinect is a 3D camera sensor which can use the 3D information of a person to provide specific 3D coordinate of a body location ?? which is called skeletal data. To the best of our knowledge, there is no publicly available skeletal dataset in literature. We used this sensor to collect data for ASL recognition problem.

With skeleton data, ASL can be seen as a sequence 3D coordinates or a 3D time series ??. Accurate recognition of such data depends on how well we can model inherent sequential pattern in data. Recurrent neural network (RNN) has shown very good performance in type kind of sequential modeling ??. In this work, we are going to implement two different architectures which are variants of RNN with skeletal ASL data. Success of a practical machine learning system depends on it’s performance after it is being deployed. In other words, the system must do well with new subjects other than those were used in training the system. To work in this direction, we proposed a calibration mechanism with our implemented system which essentially finds

what fraction of data from a new subject is needed to give a performance boost to our system. Overall we have the following contributions in this work.

- We created a new dataset for ASL and made it public
- Implementation of different RNN variants to solve ASL recognition problem with skeletal data.
- Lastly, we propose a calibration mechanism which will tell us what percentage of data from a test subject is needed to give a performance boost to the system.

## 2. Related work

Most of the sign language recognition system in the literature use RGB video data and used Hidden Markov Models (HMM) to develop the recognition system. For example Zafrulla *et al.* used worn colored gloves in hands during data collection and developed an HMM based framework for ASL phrase verification [5]. In their another work, they used hand crafted features from Kinect skeleton data [4] and accelerometer worn in hand. Primarily they used directional unit vectors from one joint to another and angles between different combinations of three joints. Having all features crafted, they used an HMM based framework to build the system. Huang *et al.* showed effectiveness of using Convolutional neural network (CNN) with RGB video data [1]. They used 3D CNN which was proposed in [2] to extract spatio-temporal features from video and built a system which does not depend on hand crafted feature for ASL recognition. Similar type of architecture was implemented by Pigou *et al.* in case of Italian gestures [3].

## 3. Dataset

Most of the dataset in literature used RGB video in ASL recognition. Using RGB data there is no way to track a specific body location, such as hand tip, wrist, neck etc. But, to correctly recognize gestures, motion trajectory of each body part could be useful. With the advent of Microsoft Kinect sensor, it is possible to track several specific body locations. Specifically, Kinect version 2 can track 25 body joints. Kinect does this by exploiting depth information in video and using a machine learning model embedded inside it. This process is called skeletal tracking. Figure 1 how Kinect sees a person as a configuration of skeletal joints ???. Since no skeletal dataset was publicly available for ASL, we decided to build our own dataset to initiate research in this direction.

### 3.1. Data collection

We collected 480 ASL samples for 4 signers and 6 daily activities related signs. Since our main goal of this work is to help deaf and hearing impaired people, we chose six activities which are important for them in daily life. Those

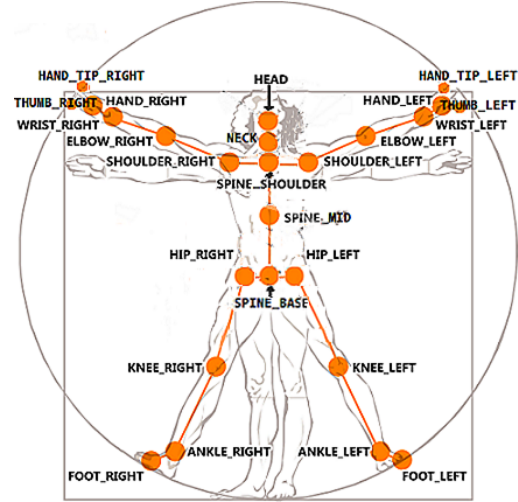


Figure 1: Kinect skeletal body configuration

are listed in Table 1. For each person, we collected 20 samples for each sign. In this way, we ended up with total  $(480)20 \times 6 \times 4$  data samples. We used Kinect version 2 sensor for data collection. For each performer we put the sensor in different distances from the person to make the dataset more challenging and realistic. Figure 2 shows an example of data. It shows six frames from the gesture *air conditioning* at the top and skeletal joints drawn to them at the bottom.

Wake Up	Turn on Thermostat
Rain	Air Condition
Open Door	Close Door

Table 1: Selected Commands for data collection and experiment

### 3.2. Preprocessing

Kinect skeleton sequences are captured in camera coordinate system which means camera position is the origin. To make sequences invariant to different camera positions and body-shapes we had to convert sequences into body coordinate system. To do this, we picked a joint location and made it origin, then we subtracted  $(x, y, z)$  locations of all other joints from it. Finally, different collected sample has different frame length. Also video file with 32 frame rate has negligible difference between two consecutive frames. That's why we sampled a fix number of frames from each sample. We used this fixed number  $T$  to be 20.



Figure 2: Example of a sign (Air Condition) performed by a person. Top row shows six frames from RGB videos increasing in time from left to right. Bottom row shows exactly same frames but with skeleton data visualization on it.

#### 4. Sequential Models for ASL

In this section, we are going to describe several sequential deep learning models and show how we can use those machine learning models for ASL recognition. Like most of the machine learning models, artificial neural networks (ANN) solves problems by translating data from an input domain into an output domain. Usually, inputs to an ANN is a 1D or multidimensional vector and output is a 1D vector. The input data goes through a non-linear transformation along the network. Using the output vector, one can do some prediction or regression tasks. The parameters of the network are learned by feeding training data to it continuously. On each iteration, first the network predicts labels for training data, then calculates loss using actual labels and finally updates the parameters based on the calculated loss. This iteration process is repeated for a predefined number of times or until a desirable accuracy is acquired on test data. While updating parameters, it starts from the output layer and move backward to the input layer. As it progresses backward, it propagates loss and calculates gradients of errors with respect to the parameters of each layer in a systematic way. Formally this is known as back propagation algorithm ?? and it is the heart of any neural network based machine learning model. That is how ANN works in short. Speaking of ANN, we usually imagine a feed forward neural network where it has an input layer, an output layer and one or more hidden layers. Parameters in different layers are different from each other. This type of network has shown tremendous performance accuracy in capturing

non linear distinction among different classes in a dataset. Specifically Convolutional neural network (CNN) which is one variant of feed forward neural network has shown human level accuracy in image classification tasks ???. Images can be thought of as 2D matrices. Patterns a machine learning model needs to learn from an image extends through this height and width dimensions. But some practical problems have also dependency over time *i.e.* information at a certain point depends on information from a previous point or a future point or both. One very common example of such data is natural language. Say for example, we want to predict the next word in the blank position of the sentence “Her mother is a school teacher. \_\_\_\_\_ teaches in ABC high school”. We can easily see that the blank position is a pronoun which refers to the word *mother* in the previous sentence. But if the word were *father* instead of the word *mother* then the pronoun at the “\_\_\_\_\_” position would be changed too. This example shows an temporal dependency in data. Although feed forward neural network has shown tremendous accuracy in finding spatial pattern in data, it can not efficiently capture such kind of temporal dependencies. One of main reasons behind this weakness is that it sees the whole input data at once and then it keeps transforming the data through its hidden layers until reaches the final layer. In other words it has to capture the patterns as a whole from the input data. We can solve this problem if there is a mechanism which at first allows us to input a portion of data in an order, then model it. Then it takes input the second portion of the ordered data, fuses it with modeled output from the first step, then model the fused data. Keep doing this until



the final portion of our ordered data will allow us to capture temporal dependencies in our data. To serve this purpose Recurrent neural network (RNN) was introduced ?? It is one of the variations of ANN which is suitable for capturing temporal dynamics of data. Figure 3 shows a high level view of such network.

$$\begin{aligned} g(x_t, h_{t-1}) &= Ux_t + Vh_{t-1} \\ f &= \phi(g) \end{aligned} \quad (1)$$

Eq 1 shows equations of an usual RNN. Here  $x_t$  is the current input and  $h_{t-1}$  is the previous RNN state.  $U$  and  $V$  are the network parameters associated with  $x_t$  and  $h_{t-1}$  respectively.  $\phi$  is a nonlinear function which we can choose from several alternatives. Some common choices are *sigmoid*, *ReLU*, *tanh* etc. Two of the most important differences between RNN and feed forward neural networks are time ordered input and network parameter sharing. We can see from Figure 3 that input is fed to the network in different time steps, not at once like feed forward case. However the parameters the network learns throughout different time steps are same. At any time step the network process current input along with whatever previously seen and processed data. This helps to capture the temporal dynamics in data. But the basic RNN has problem dealing with long term dependencies in data ?? This happens because the basic RNN has no control over what it should learn and what should forget. It just keeps updating its memory as it goes which is not desirable in most practical cases. Sometimes a memory from the long past could be useful for current prediction than a memory from the recent past. Another problem which makes training an RNN problematic is known as vanishing gradient problem ?? Since network parameters of an RNN are shared over time, at any time step error derivative not only depends on the current input of network but also on the previous state of the network. So there is a multiplicative term of error derivatives back to time step 0 to calculate error derivative of current time. This is a modification of back propagation algorithm and known as back propagation through time (BPTT) ?? If individual gradients are close to zero this multiplicative term would become very close to zero and eventually gradients disappear. The opposite can also happen, which is called exploding gradient problem, however this is more obvious and easy to detect. The vanishing gradient problem is not easy to detect and serves as a silent killer to the network. Several approaches have been adopted to deal with vanishing gradient problem. Some of them are to careful initialization of network parameters, early stopping etc ?? But most effective solution was to modify the RNN architecture so that it has its own memory and using that memory it can control what to remember and what to forget. This architecture is called long short term memory (LSTM) network ?? Another similar architecture called gated recurrent unit (GRU)

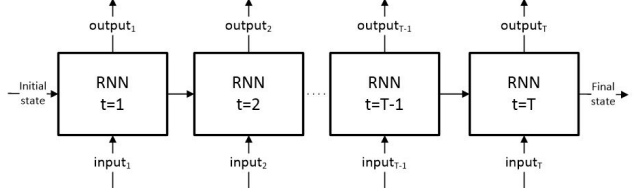


Figure 3: An RNN networks unrolled over T time steps. At each time step the network has one raw input and one processed input from previous time step. Combining these two it produces embedding for current time step

has been also proposed and has shown robustness against vanishing/exploding gradient problem.

#### 4.1. Long short term memory

As we see from eq 1, at each time step, traditional RNN can be seen as a nonlinear function of the current input and the previous step. While vanilla RNN is a direct transformation from the previous state and the current input, LSTM takes a different approach. It maintains an internal memory and it has mechanism to update and use that memory. The mechanism consists of four separate neural networks also called gates. Figure 4 shows a cell of an LSTM network. Forget, input, update and output gate are four different gates

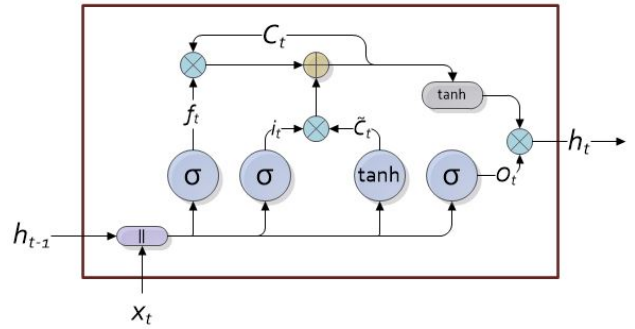


Figure 4: An LSTM cell

represented by four circles and symbolized as  $f_t$ ,  $i_t$ ,  $\tilde{C}_t$  and  $o_t$  respectively.  $\oplus$  and  $\otimes$  represents element wise addition and multiplication respectively. Two vertical bars inside the rounded rectangle means concatenation operation on its inputs.

$$\begin{aligned} h_t &= o_t \otimes \tanh(C_t) \\ o_t &= \sigma(W_o * \text{concat}(h_{t-1}, x_t)) \\ C_t &= (f_t \otimes C_{t-1}) \oplus (i_t \otimes \tilde{C}_t) \\ i_t &= \sigma(W_i * \text{concat}(h_{t-1}, x_t)) \\ \tilde{C}_t &= \tanh(W_{\tilde{C}} * \text{concat}(h_{t-1}, x_t)) \end{aligned} \quad (2)$$

Figure 4 shows input at the current time step  $x_t$  and the previous state  $h_{t-1}$  enter into the cell and get concatenated. Then the forget gate processes it to remove unnecessary information and outputs  $f_t$  which gets multiplied with the previously stored memory  $C_t$  and produces a refined memory for the current time. Meanwhile, the input gate and the update gate process the concatenated input and convert it into a candidate memory for the current time step. The refined memory from the previous step and proposed candidate memory of the current step get added to produce the final memory for the current step. This addition could render the output out of scale. Because of that, a squashing function is followed which is a *hyperbolic tangent* in our case. Its job is to scale the elements of the output vector into a fix range. Finally  $o_t$ , the output from output gate gets multiplied with output from the squashing function and produces the output for the current time step.

## 4.2. Gated recurrent unit

A gated recurrent unit (GRU) is another variation of RNN for modeling temporal dynamics of data. It works in the same way as LSTM, but has a simpler construction. Like LSTM it has the mechanism to control flow of information over time, but it lacks the separate memory cell. Figure 6 shows an example of GRU cell. Symbols represent similar things as described in section 4.1. To better understand GRU cell, first we have to look at the output of it. The output is formed by a weighted linear sum of two components  $\tilde{h}_t$  and  $h_{t-1}$  as seen from the Eq. 3. Here  $\text{concat}(x, y)$  means concatenation of input vectors  $x$  and  $y$ . Element wise multiplication and addition are represented by  $\otimes$  and  $\oplus$  respectively.

$$\begin{aligned} h_t &= (1 - z_t)h_{t-1} \oplus z_t\tilde{h}_t \\ z_t &= \sigma(W_z * \text{concat}(h_{t-1}, x_t)) \\ \tilde{h}_t &= \tanh(W_{\tilde{h}} * \text{concat}(r_t \otimes h_{t-1}, x_t)) \\ r_t &= \sigma(W_r * \text{concat}(h_{t-1}, x_t)) \end{aligned} \quad (3)$$

The term  $\tilde{h}_t$  represents newly proposed candidate embedding and  $h_{t-1}$  is simply the previous state. Candidate embedding  $\tilde{h}_t$  is a function of the current input and purified previous state. This purification is achieved by the reset gate  $r_t$  whose job is to decide what portion of previous state should contribute to candidate embedding  $\tilde{h}_t$ . Reset gate  $r_t$  itself is a separate neural network and its output close to 0 means network takes a small portion of previous state to form a current candidate state embedding. The update gate  $z_t$  decides what fractions of the candidate state and the previous state comprise the new state. It does so by taking a weighted sum of the candidate state and the previous state. The weights are  $z_t$  and  $1 - z_t$ . This whole new state is outputted where in case of LSTM, the output of the new state is controlled by an output gate as shown in Eq. 2.

## 4.3. Our methodology

ASL signs are usually performed using two hands, head movements and in some cases facial expressions. Motion dynamics of two hands are primarily responsible for creating different cues for various signs. Each joint in two hands and face area follows a specific sequential pattern for a particular gesture sign. For example, from Figure 2 we see the person is performing gesture for *air condition*. We can see that, for this sign first, one has to lift and show palm area of the left hand and then he has to lift both hands and make a back and forth motion with all fingers. This is basically a sequence of joints in both hands. Given the efficacy of variants of RNN in encoding sequential pattern in data, we want to see how those models perform on sequence data like ASL sign. To achieve this goal, we implemented LSTM and GRU model for ASL recognition problem. We noticed that, the joints located above the waist are mostly useful in case of sign language. That's why we took only 12 joints into consideration in our experiment. These joints come from both hands, head, neck and spine area. Figure 5 shows our implementation. For each time step, first, we concatenate 3D coordinates of 12 joints and then feed to our RNN network. Finally, we take the final state of the network and send it to a softmax layer and get prediction probabilities over our classes.

## 5. Experiments

In this section, we are going to describe all of our experimental designs. First we are going to explain how we evaluated model performances. We are also going to describe a calibration mechanism we proposed to find a good seed percentage of train data for test person. Then we are going to present implementation details and results we got from experiments.

### 5.1. Evaluation criteria

Our experiment on the dataset is two fold. First, we did experimentation with two deep learning sequence models, namely GRU and LSTM. Also, we proposed a calibration mechanism which essentially infers what fraction of data from an unseen test subject is needed to boost performance of our models. This calibration experiment is done in cross subject manner *e.g.* we train our models on some subjects and test on a different subject. This makes a lot of sense in real life situation where a trained model has to perform well on data it has not seen yet.

Table 2 represents some terminologies helpful to understand calibration mechanism described next.

We have data on four different subjects. First, we choose a test subject and take out 50% of data from this subject as test data say  $D_{test}$ . We use other 50% of data as seed data say,  $D_{seed}$  to train model along with the data from other

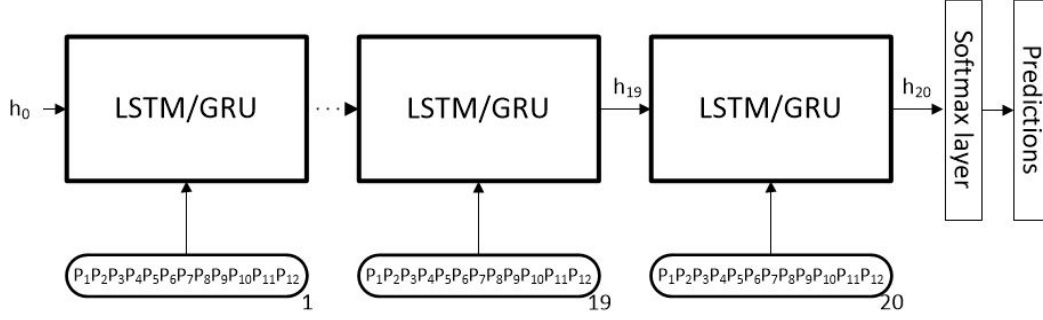


Figure 5: Implementation of RNN models to recognize ASL from skeletal data. Input at each time step is a concatenation of 3D coordinates of 12 joints.  $P_i$  is the  $(x, y, z)$  of  $i^{th}$  joint. Subscript of each input vector represents time step

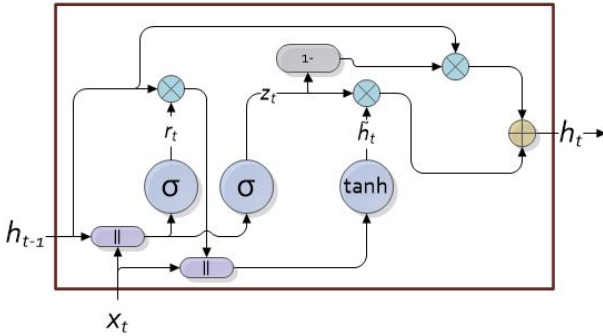


Figure 6: A GRU cell

Terms	Description
$D_{train}$	Training data, whole data from the subjects other than the test subject
$D_{test}$	Test data, 50% data from test subject used for testing models
$D_{seed}$	Seed data, 50% data from test subject used for boosting up model's performance on $D_{test}$

Table 2: Different terminologies regarding data split for one test subject

three subjects say,  $D_{train}$ . We first train the model only with  $D_{train}$  and test on  $D_{test}$ . Then we train our model with  $D_{train}$  and 10% from  $D_{seed}$  and again test on  $D_{test}$ . Then again we train the model with  $D_{train}$  and this time 20% of  $D_{seed}$  and test on  $D_{test}$ . As we increase the seed data in our training data, we should expect better test accuracy. We keep doing this until we use up to full  $D_{seed}$  for training with  $D_{train}$  and record different testing results on  $D_{test}$ . We repeat the whole procedure by taking each

subject as a test subject one at a time. We want to see here what fraction of seed data is required during training to get a good performance from our model. Figure 7 shows data splitting strategy for better understanding.

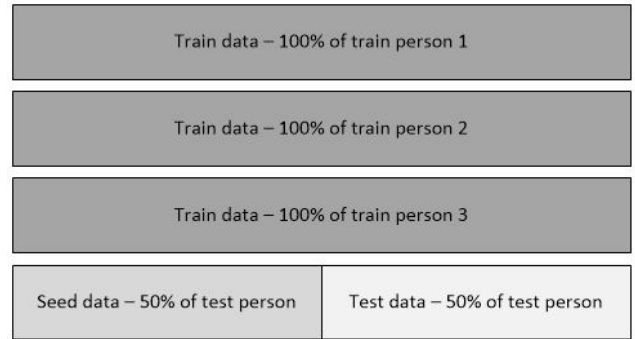


Figure 7: Train and test data split for a test person

## 5.2. Results

We got an almost similar types of results from both of LSTM and GRU implementation. However, we found that 30% to 50% data from the test subject as seed is good enough to give a performance boost. Table 3 shows the results for each combination of a test person, a % of seed data used in training for that person and a network type. We observe from the result that the highest accuracy is obtained using 30% or more seed data. This suggests that even if we have a small amount of seed data for a test persons, our system can boost up to a higher accuracy. In a practical scenario, system will be able to recognize new subjects correctly within a short period of subject's interaction with the system. Figure 8 shows plots of test accuracy vs percentage of data used in training. It shows the LSTM results at the top and the GRU results at the bottom. Plots are also suggesting that 30% or more of seed data is good enough for a

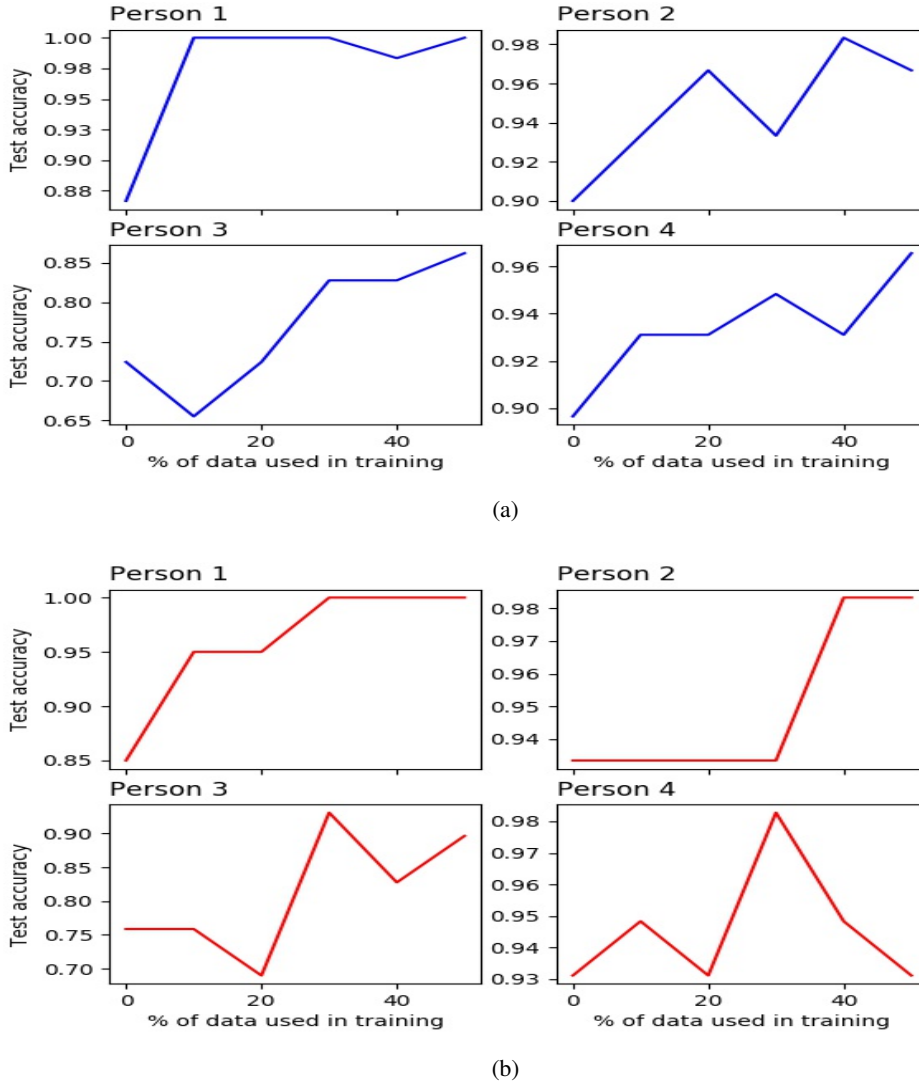


Figure 8: Test accuracy. X-axis represents percentage of data from the test subject used in training and Y-axis is the test accuracy on 50% hold out data for test subject. (a) Test accuracy for LSTM network. Each plot represents accuracy on test data for each subject in our dataset. (b) Test accuracy for GRU network in same way as (a).

performance increase to the model.

We also tested our model in conventional train test split manner. In that case, we took 30% of whole data as test data. Table 4 shows accuracy for both type of models after performing k-fold cross validation.

### 5.3. Implementation details

We used same value for hyper-parameters such as learning rate, state size of RNNs, batch size etc. for both compared models. We did a grid search for finding best hyper-parameters. Our experiments found the best learning rate as  $5e^{-4}$  and the state size for RNNs as 36. For each joint

state size is 3 which gives us the whole state size of 36 for 12 joints. We used Adam Optimizer for optimizing our networks. We used this optimizer because it has shown superior performance in training different types of network. It also adapts learning rate and decay rate as training continues. We used mini-batch size of 1 which means the network updates its parameter after seeing one training example. We keep training our network until a good accuracy is achieved. However, to battle over-fitting of the network, we followed early stopping strategy. If there is no increase in the testing accuracy for 40 epochs, we decided to stop training. We implemented our models using TensorFlow deep

	% seed	0%	10%	20%	30%	40%	50%
P1	LSTM	0.86	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	0.99	<b>1.00</b>
	GRU	0.85	0.95	0.95	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>
P2	LSTM	0.90	0.93	0.96	0.93	<b>0.98</b>	0.97
	GRU	0.93	0.93	0.93	0.93	<b>0.98</b>	<b>0.98</b>
P3	LSTM	0.72	0.65	0.72	0.82	0.82	<b>0.86</b>
	GRU	0.75	0.75	0.69	0.93	0.83	<b>0.90</b>
P4	LSTM	0.89	0.93	0.93	0.94	0.93	<b>0.96</b>
	GRU	0.93	0.95	0.93	<b>0.98</b>	0.95	0.95

Table 3: Test accuracies for each combination of test person, network type and % of data from test subject used in training.  $P_i$  denotes the  $i^{th}$  person

Model type	Test accuracy
LSTM	<b>0.94</b>
GRU	<b>0.91</b>

Table 4: Results from 70%–30% train-test split

learning library version 1.7 ???. Experiments were run on Tesla K80 GPU facilitated by ARGO research computing cluster ???.

## 6. Conclusion

Here goes conclusion

## References

- [1] 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. 2
- [2] 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. 2
- [3] L. Pigou, S. Dieleman, P.-J. Kindermans, and B. Schrauwen. Sign language recognition using convolutional neural networks. In L. Agapito, M. M. Bronstein, and C. Rother, editors, *Computer Vision - ECCV 2014 Workshops*, pages 572–578, Cham, 2015. Springer International Publishing. 2
- [4] 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. 2
- [5] 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. 2