

LABORATORY ORIENTED PROJECT

FINAL REPORT

Finger Spelling to Text Translation Using Deep Learning

Prof. In Charge: Dr.Surekha Bhanot

SAGNIK MAJUMDER | 2014A8TS0464P



Birla Institute of Technology and Science, Pilani - Pilani Campus

Contents

A. Introduction

B. Tools Used

C. Datasets

D. Work Pipeline

E. Final Results

A. INTRODUCTION

'Machine Translation'(MT) is in great need nowadays and there is a lot of work on it all around the globe. MT deals with the task of translating text from one input language into proper semantic text of another language i.e. meaningful text in an output language. MT is generally associated only with text to text translation keeping its purpose true to the actual meaning of translation i.e. translation from English to German etc. But, in this work, I give a new dimension to translation. I introduce here translation from a gesture or sign based language into a textual language. The scope of this work today is quite promising keeping in mind the fact that not much has been done in this particular domain.

The prospects of this work are many -

1. Sign or gesture based language, popularly known as sign language, is a relatively slow language in comparison to other verbal or textual language because of the time taken to produce signs. So, a good text generation model can be used to complete a sentence or generate further text from some starting primer in some sign language. This can speed up communication and can immensely help both the speaker.
2. When the speaker is using sign language, the listener may not always know it and may need the help of a translator for each sentence the speaker speaks. I can do away with the translator if I can use such a text generation model which can generate text and convey the meaning the speaker is trying to put through. The text generation can be used to complete the words from just a few initial characters. This can be rendered useful in a lot of basic conversation i.e. in case of the most common words that the verbally challenged speaker might use in his/her day to day communication.

There are many other uses of such a system which I will temporarily skip for the reason of brevity.

The fundamental tools behind the building of such a model are Convolutional Neural Networks (CNN) and Sequential Neural Networks (SNN). The CNNs are used to convert a particular sign or gesture into its corresponding character in some other verbal/textual language. The CNNs are basically used to generate the primer or the introductory text that will be used for further text generation. The SNNs take this primer text as the input and generate further text to complete the word or the sentence using Recurrent Learning. So, the model, in fullness, employs a CNN to convert images into

characters, to be used as the text primer, and the SNN will use the primer text to generate further text to complete a word or a sentence.

A Brief Overview of Deep Learning

Deep Learning is a subfield of machine learning concerned with algorithms inspired by the structure and function of the brain called artificial neural networks. It refers to artificial neural networks that are composed of many layers. It is just another name for artificial neural networks (ANN for short) but in a more refined and easier avatar. They have been existing for more than 40 years. In 2000's, NVIDIA accelerated ANN by bringing their chips for scientific computing and from then onwards, use of neural networks has picked up.

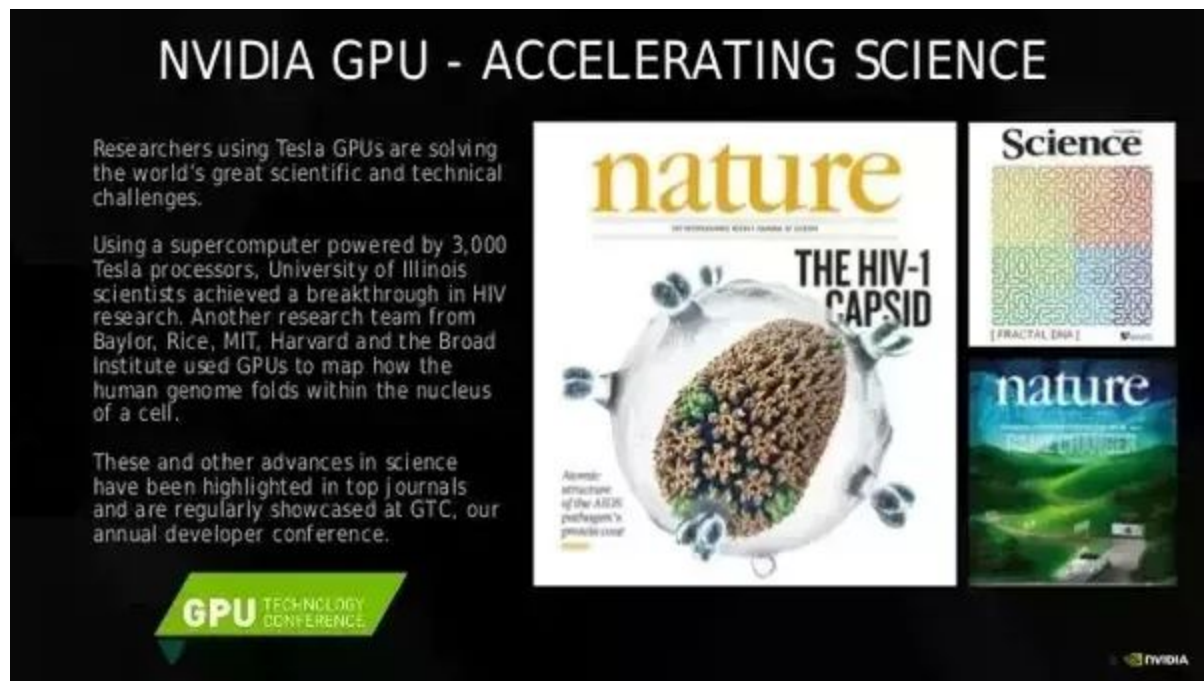


Figure 1: GPUs Have Helped Deep Learning Advance

It's a growing trend in ML due to some favorable results in applications where the target function is very complex and the datasets are large. For example in Hinton et al. (2012), Hinton and his students managed to beat the status quo prediction systems on five well known datasets: Reuters, TIMIT, MNIST, CIFAR and ImageNet. This covers speech, text and image classification - and these are quite mature datasets, so a win on any of these gets some attention. A win on all of them gets a lot of attention.

In machine learning and cognitive science, artificial neural networks (ANNs) are a family of models inspired by biological neural networks (the central nervous systems of animals, in particular the brain) and are used to estimate or approximate functions that

can depend on a large number of inputs and are generally unknown. Artificial neural networks are generally presented as systems of interconnected "neurons" which exchange messages between each other. The connections have numeric weights that can be tuned based on experience, making neural nets adaptive to inputs and capable of learning.

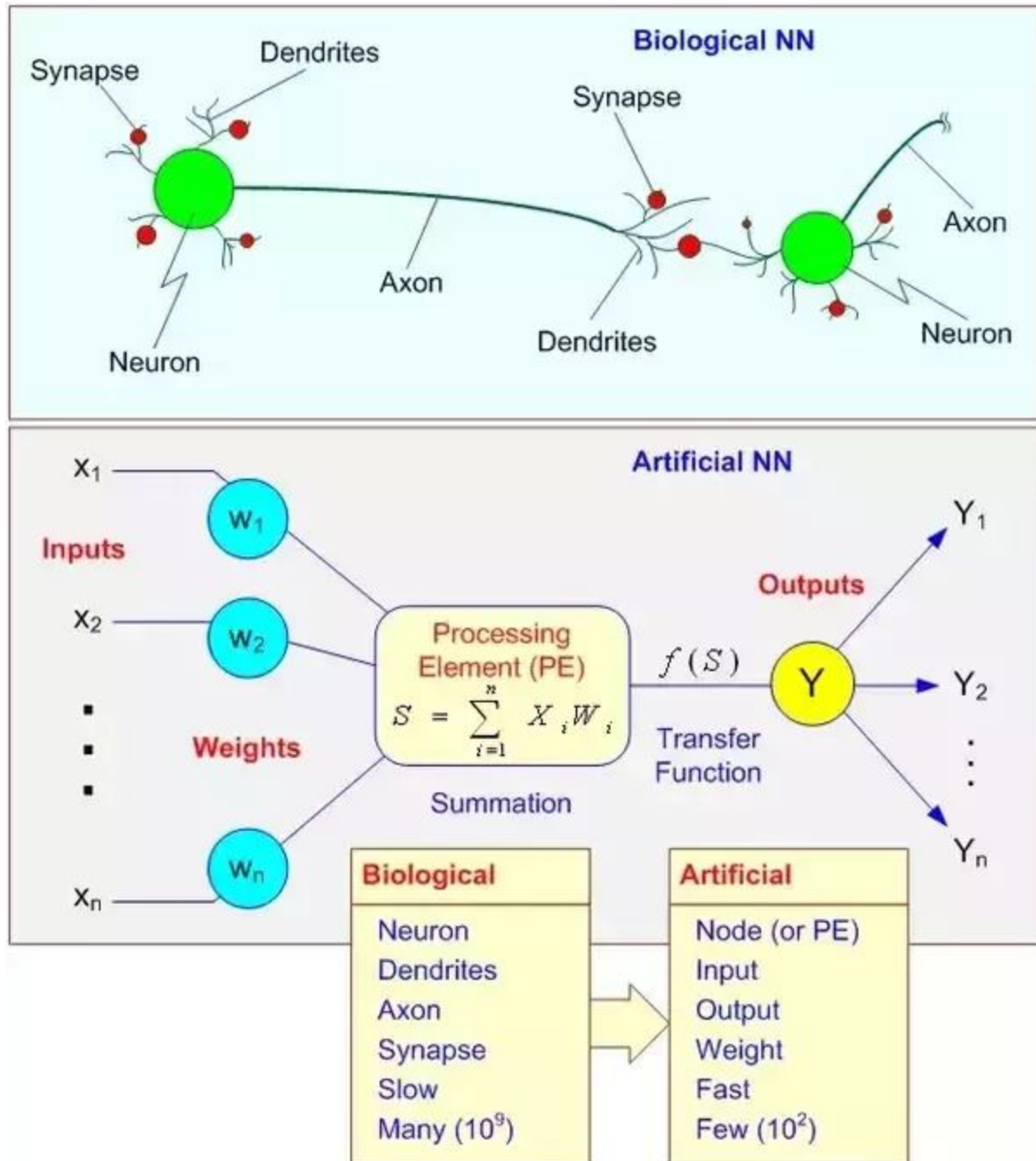


Figure 2: A Comparison Between Biological Neurons and an ANN in Deep Learning

Why 'Deep Learning' is called deep? It is because of the structure of ANNs. Earlier 40 years back, neural networks were only 2 layers deep as it was not computationally

feasible to build larger networks. Now it is common to have neural networks with 10+ layers and even 100+ layer ANNs are being tried upon.

One can essentially stack layers of neurons on top of each other. The lowest layer takes the raw data like images, text, sound, etc. and then each neuron stores some information about the data they encounter. Each neuron in the layer sends information up to the next layers of neurons which learn a more abstract version of the data below it. So the higher one goes up, the more abstract features the layers learn. One can see in the picture below has 5 layers in which 3 are hidden layers.

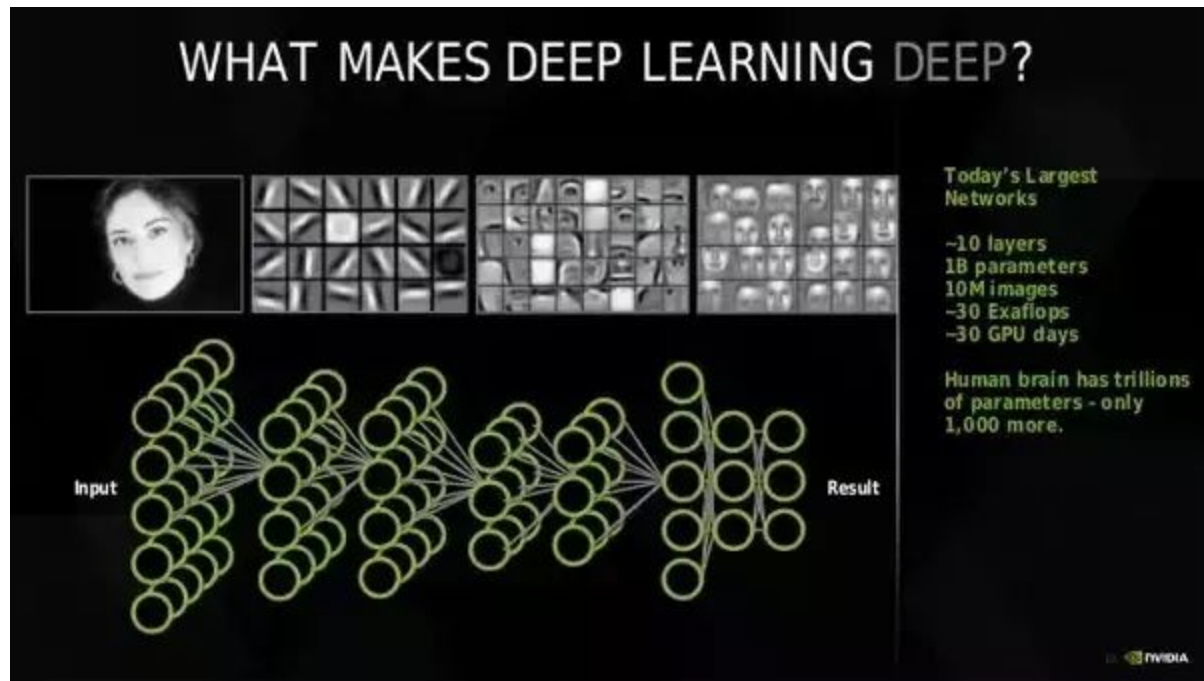


Figure 3: A Representative Deep Neural Network

Why feature engineering may turn obsolete?

Today as a data scientist, one may be learning feature engineering as a part of machine learning skill. In that, one need to transform one's data to the computer in a form that it can understand. For that one may use R or Python or spreadsheet software to translate the data. One may be converting one's data into a large spreadsheet of numbers containing rows and columns of instances and features. Then one feeds this data into the machine learning algorithm and it tries to learn from this data. As engineering the features is a time consuming task, we need to extract only the relevant features that improve our model. But as one is generally unaware of the usefulness of these features until one trains and tests the model, one is caught into the vicious cycle of developing new features, rebuilding the model, measure results, and repeat until one is satisfied with the results. This is very time consuming task and takes lot of your time.

How deep learning may save one's time?

In deep learning, ANNs are automatically extracting features instead of manual extraction in feature engineering. One can take the example of image as input. Instead of a user taking an image and hand compute features like distribution of colors, image histograms, distinct color count, etc., the user just have to feed the raw images in ANN. ANNs have already proved their worth in handling images, but now they are being applied to all kinds of other datasets like raw text, numbers etc. This helps the data scientist to concentrate more on building deep learning algorithms.

Big Data is required for deep learning

Soon, feature engineering may turn obsolete but deep learning algorithm will require massive data for feeding into our models. Fortunately, we now have big data sources not available two decades back – facebook, twitter, Wikipedia, project Gutenberg etc. However, the bottleneck remains in cleaning and processing these data into required format for powering the machine learning models. More and more big data will be made available for public consumption in near future.

Where can one get help in deep learning?

Open sourcing is predominant in deep learning, Tensorflow, Torch, Keras, Big Sur hardware, DIGITS and Caffe are some of the massive deep learning projects. In academic research, there are lots of papers with algorithm source code along with their findings. arXiv.org (a Cornell University Library) has open access to over 1 million papers in deep learning.

Does one need a computer expertise to learn deep learning?

One does not need a rigorous computer expertise to learn deep learning. The user's domain knowledge of the discipline will help him in building deep learning models. If one has learnt any data science language like 'R' or 'Python', spreadsheets like 'Excel' or any other basic programming and studied 'STEM', then one is proficient to delve into deep learning.

What are deep learning advantages?

One advantage that has already been explained is that the user doesn't have to figure out the features ahead of time. Another advantage is that one can use the same neural net approach for many different problems like Support Vector Machines, Linear Classifier, Regression, Bayesian, Decision Trees, Clustering and Association Rules. It is fault tolerant and scales well. Many perceptual tasks have been performed with help of CNNs. Some case studies are given below:

CNNs DOMINATE IN PERCEPTUAL TASKS

- Handwriting recognition MNIST (many), Arabic HWX (IDSIA)
- OCR in the Wild [2011]: StreetView House Numbers (NYU and others)
- Traffic sign recognition [2011] GTSRB competition (IDSIA, NYU)
- Asian handwriting recognition [2013] ICDAR competition (IDSIA)
- Pedestrian Detection [2013]: INRIA datasets and others (NYU)
- Volumetric brain image segmentation [2009] connectomics (IDSIA, MIT)
- Human Action Recognition [2011] Hollywood II dataset (Stanford)
- Object Recognition [2012] ImageNet competition (Toronto)
- Scene Parsing [2012] Stanford bgd, SiftFlow, Barcelona datasets (NYU)
- Scene parsing from depth images [2013] NYU RGB-D dataset (NYU)
- Speech Recognition [2012] Acoustic modeling (IBM and Google)
- Breast cancer cell mitosis detection [2011] MITOS (IDSIA)

Slide credit: Yann Lecun, Facebook & NYU

Figure 4: A Few Case Studies Showing the Vast Prospects of Deep Learning

How classical ML using feature engineering can be compared with deep learning using CNN? Well it can be seen with the help of these pictures.

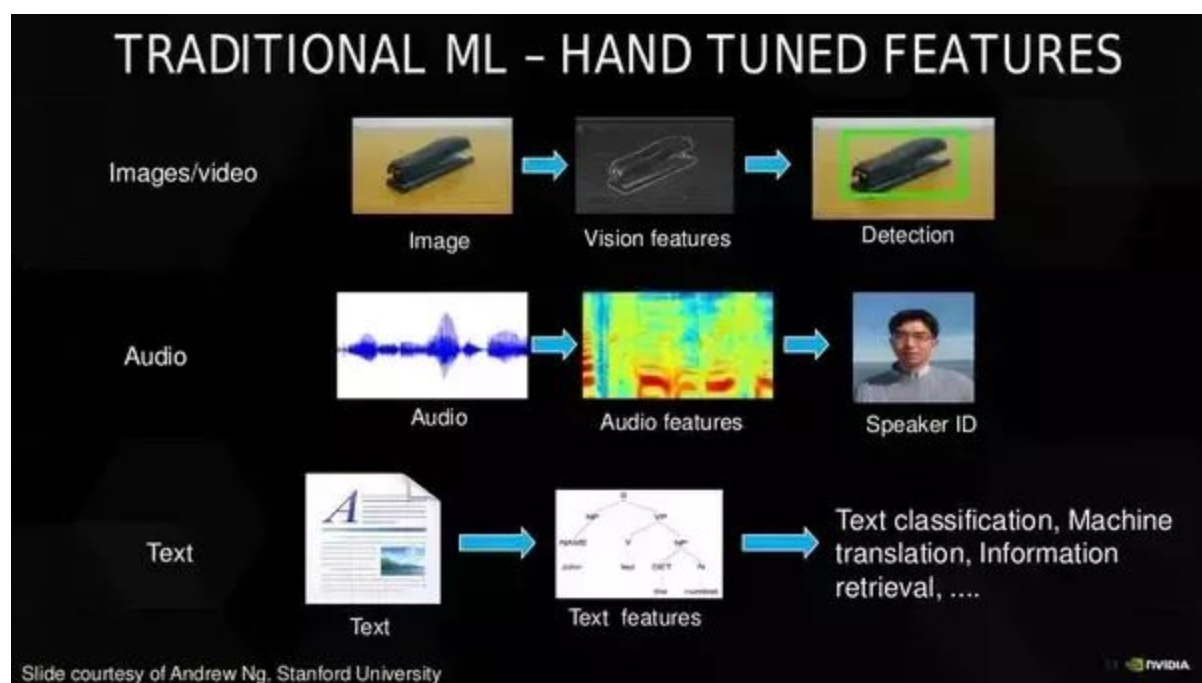


Figure 5 : Traditional ML

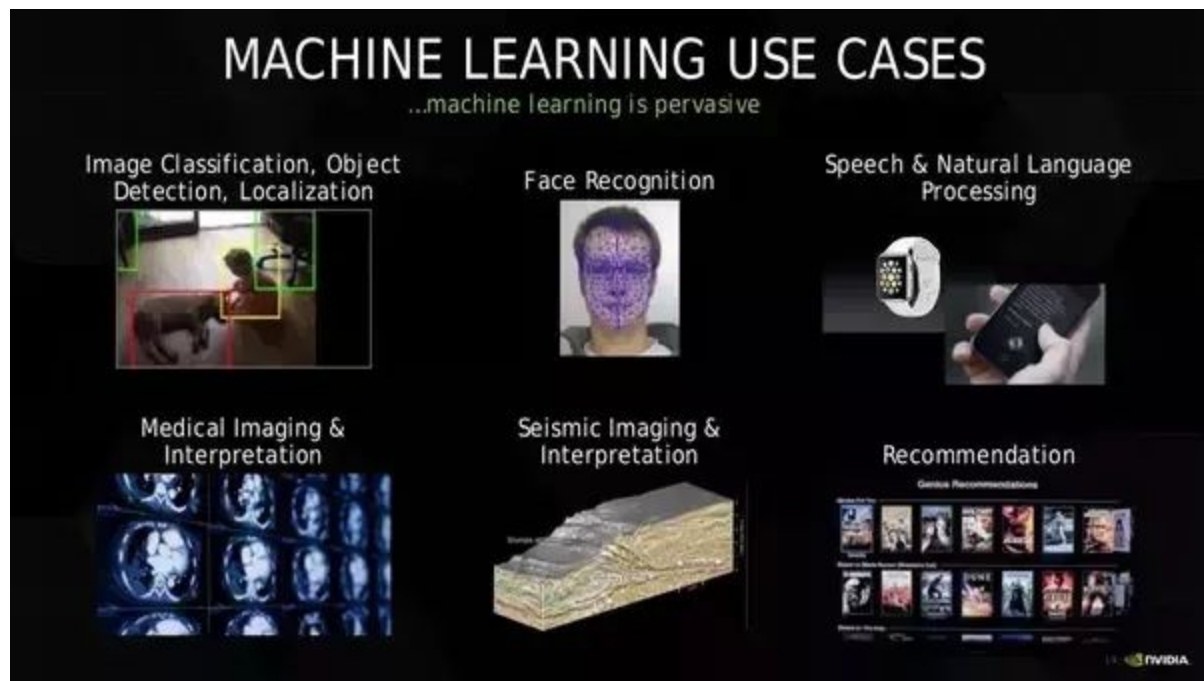


Figure 6 : Deep Learning in ML

From these pictures, one can see the wide applicability of deep learning in all aspects of life.

B. TOOLS USED

1. CONVOLUTIONAL NEURAL NETWORKS -

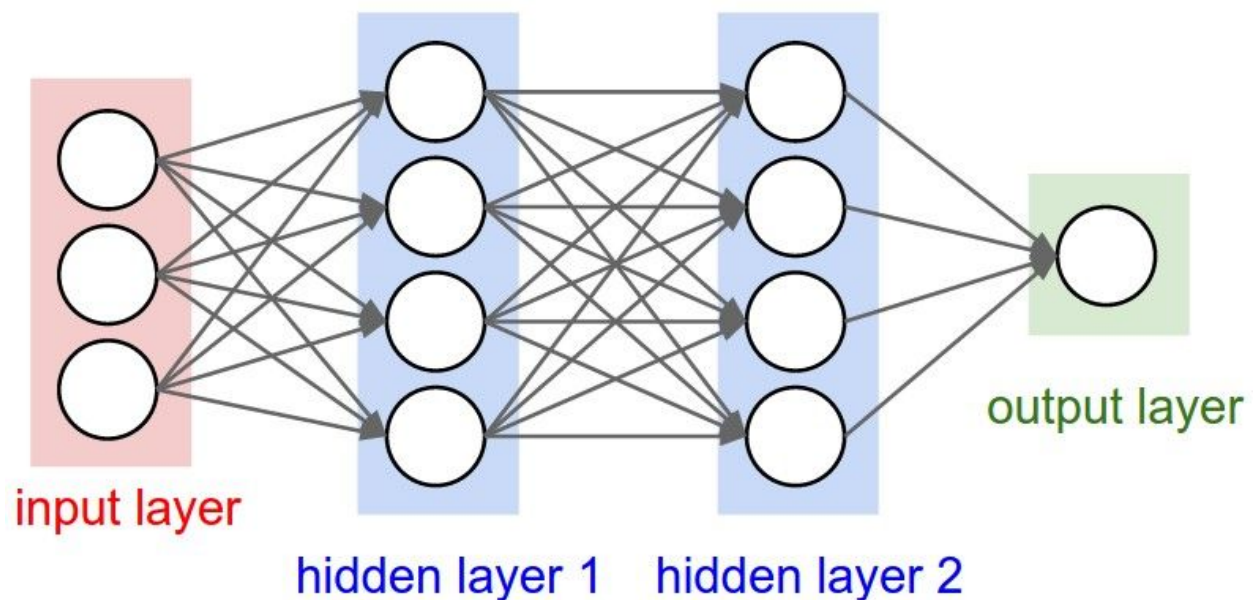
Convolutional Neural Networks are very similar to ordinary Neural Networks: they are made up of neurons that have learnable weights and biases. Each neuron receives some inputs, performs a dot product and optionally follows it with a non-linearity. The whole network still expresses a single differentiable score function: from the raw image pixels on one end to class scores at the other.

For eg.in CIFAR-10, images are only of size 32x32x3 (32 wide, 32 high, 3 color channels), so a single fully-connected neuron in a first hidden layer of a regular Neural Network would have $32 \times 32 \times 3 = 3072$ weights. This amount still seems manageable, but clearly this fully-connected structure does not scale to larger images. For example, an image of more respectable size, e.g. 200x200x3, would lead to neurons that have $200 \times 200 \times 3 = 120,000$ weights. Moreover, one would almost certainly want to have

several such neurons, so the parameters would add up quickly. Clearly, this full connectivity is wasteful and the huge number of parameters would quickly lead to overfitting.

3D volumes of neurons. Convolutional Neural Networks take advantage of the fact that the input consists of images and they constrain the architecture in a more sensible way. In particular, unlike a regular Neural Network, the layers of a ConvNet have neurons arranged in 3 dimensions: **width, height, depth**. (It is to be noted that the word *depth* here refers to the third dimension of an activation volume, not to the depth of a full Neural Network, which can refer to the total number of layers in a network.) For example, the input images in CIFAR-10 are an input volume of activations, and the volume has dimensions 32x32x3 (width, height, depth respectively). The neurons in a layer will only be connected to a small region of the layer before it, instead of all of the neurons in a fully-connected manner. Moreover, the final output layer would for CIFAR-10 have dimensions 1x1x10, because by the end of the ConvNet architecture it will reduce the full image into a single vector of class scores, arranged along the depth dimension.

Here is a visualization:



Left: A regular 3-layer Neural Network. Right: A ConvNet arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers. Every layer of a ConvNet transforms the 3D input volume to a 3D output volume of neuron activations. In this example, the red input layer holds

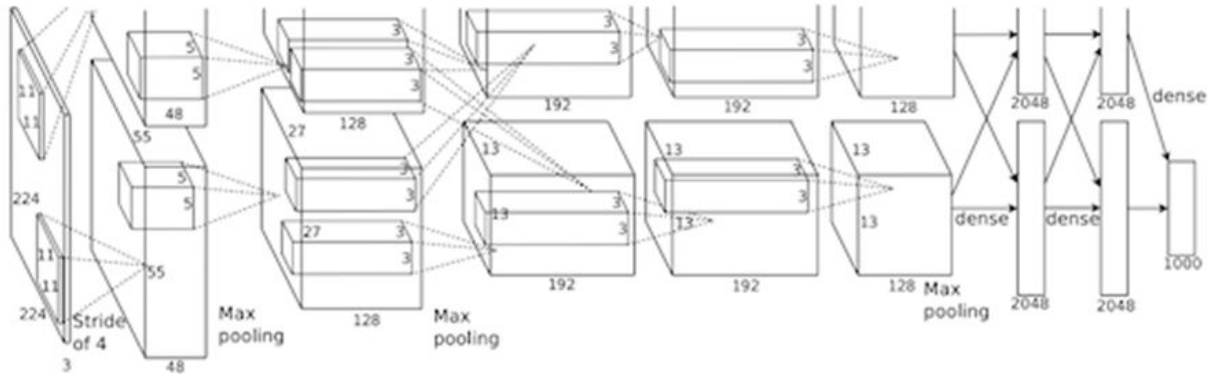
the image, so its width and height would be the dimensions of the image, and the depth would be 3 (Red, Green, Blue channels).

A ConvNet is made up of Layers. Every Layer has a simple API: It transforms an input 3D volume to an output 3D volume with some differentiable function that may or may not have parameters.

AlexNet -

The one that started it all (Though some may say that Yann LeCun's [paper](#) in 1998 was the real pioneering publication). This paper, titled "ImageNet Classification with Deep Convolutional Networks", has been cited a total of 6,184 times and is widely regarded as one of the most influential publications in the field. Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton created a "large, deep convolutional neural network" that was used to win the 2012 ILSVRC (ImageNet Large-Scale Visual Recognition Challenge). For those that aren't familiar, this competition can be thought of as the annual Olympics of computer vision, where teams from across the world compete to see who has the best computer vision model for tasks such as classification, localization, detection, and more. 2012 marked the first year where a CNN was used to achieve a top 5 test error rate of 15.4% (Top 5 error is the rate at which, given an image, the model does not output the correct label with its top 5 predictions). The next best entry achieved an error of 26.2%, which was an astounding improvement that pretty much shocked the computer vision community. Safe to say, CNNs became household names in the competition from then on out.

In the paper, the group discussed the architecture of the network (which was called AlexNet). They used a relatively simple layout, compared to modern architectures. The network was made up of 5 conv layers, max-pooling layers, dropout layers, and 3 fully connected layers. The network they designed was used for classification with 1000 possible categories.



AlexNet architecture (May look weird because there are two different "streams". This is because the training process was so computationally expensive that they had to split the training onto 2 GPUs)

Main Points

- Trained the network on ImageNet data, which contained over 15 million annotated images from a total of over 22,000 categories.
- Used ReLU for the nonlinearity functions (Found to decrease training time as ReLUs are several times faster than the conventional tanh function).
- Used data augmentation techniques that consisted of image translations, horizontal reflections, and patch extractions.
- Implemented dropout layers in order to combat the problem of overfitting to the training data.
- Trained the model using batch stochastic gradient descent, with specific values for momentum and weight decay.
- Trained on two GTX 580 GPUs for five to six days.

Why It's Important

The neural network developed by Krizhevsky, Sutskever, and Hinton in 2012 was the coming out party for CNNs in the computer vision community. This was the first time a model performed so well on a historically difficult ImageNet dataset. Utilizing techniques that are still used today, such as data augmentation and dropout, this paper really illustrated the benefits of CNNs and backed them up with record breaking performance in the competition.

2. SEQUENTIAL NEURAL NETWORKS -

There are many types of SNNs. SNNs are used for classification purposes in cases of sequential data where the data at a particular moment is dependent on the data from

the past i.e. temporal data. Examples of temporal data are videos or frame sequence, text from a book or a speech among other. SNNs deal with such data by storing its behavior in response to past data or the history of the network. The way to store the history varies from one type of SNN to the other and that's what differentiates one type of SNN from another. A few popular types of SNNs are Recurrent Neural Networks (RNN), Long Short Term Memory (LSTM) and Gated Recurrent Network (GRU).

RNN -

A recurrent neural network (RNN) is a class of artificial neural network where connections between units form a directed cycle. This allows it to exhibit dynamic temporal behavior. Unlike feedforward neural networks, RNNs can use their internal memory to process arbitrary sequences of inputs. This makes them applicable to tasks such as unsegmented, connected handwriting recognition or speech recognition. The governing equations of RNNs are

$$h_t = Wf(h_{t-1}) + W^{(hx)}x_t$$

$$\hat{y} = W^{(S)}f(h_t)$$

Here, 'W' is the weight matrices and 'h' is the state of the RNN. 'h' is a function of 't' and it is called the RNN state at time 't'. 'y' is the output of the the RNN at the particular time instant 't'.

A drawback of normal RNN is that due to the long term time dependence of the RNN state on past states, the gradients may die out during learning due to multiple steps of backpropagation. This is called the problem of 'vanishing gradients'. So, a better replacement for RNNs are LSTMs and GRUs.

LSTM -

LSTMs are invented as a solution to the 'vanishing gradients' of the RNNs. It exploits short term dependencies in place of long term dependencies which is the root cause of the vanishing gradients. The LSTMs use a memory cell denoted by 'c' which represents

the short term memory state of the LSTM at time 't'. The governing equation of the LSTMs are

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i)$$

$$\tilde{C}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c)$$

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f)$$

$$C_t = i_t * \tilde{C}_t + f_t * C_{t-1}$$

I have used LSTMs for the purpose of text generation given a primer or starting text. LSTMs are still a pretty popular choice for a lot of sequential operations because of its simplicity and convenience of training in comparison to the other better alternatives.

I have skipped describing the GRUs since they are not relevant for this particular task as of now though I plan to try with GRUs also once I am finished with modelling the whole thing

C. DATASETS -

1. Datasets for Text Generation -

a) The Penn Treebank, in its eight years of operation (1989-1996), produced approximately 7 million words of part-of-speech tagged text, 3 million words of skeletally parsed text, over 2 million words of text parsed for predicate argument structure, and 1.6 million words of transcribed spoken text annotated for speech disfluencies. The material annotated includes such wide-ranging genres as IBM computer manuals, nursing notes, Wall Street Journal articles, and transcribed telephone conversations, among others. This paper describes the design of the three annotation schemes used by the Treebank: POS tagging, syntactic bracketing, and disfluency annotation and the methodology employed in production. All available Penn Treebank materials are distributed by the Linguistic Data Consortium <http://www ldc upenn edu>. The whole dataset is split up into training, validation and testing where the validation set can be used to check the BLEU score and set the hyperparameters.

b) This is a full extract from a GitHub repository which I used for my initial trials with my model. I split it into training, validation and testing for the purpose. The splitting was done absolutely randomly.

2. Datasets for Image to Character Conversion

I used the American Sign Language (ASL) finger-spelling dataset for my work.

Finger-spelling refers to the system of signs or gestures where each gesture represents a single character. The datasets contain a set of RGB and depth images for each letter in the alphabet, organized by subject, for estimating generalization.

There are two sub-datasets to this dataset which are graded on the basis of difficulty of classification. They are as follows

a) The first dataset comprises 24 static signs (excluding letters j and z because they involve motion). This was captured in 5 different sessions, with similar lighting and background.

b) The second dataset (depth only) is captured from 9 different persons in two very different environments and lighting.

D. WORK PIPELINE

1. Language Model using Common Bag of Words Model -

There are an estimated 13 million tokens for the English language but are they all completely unrelated? Feline to cat, hotel to motel? I think not. Thus, I want to encode word tokens each into some vector that represents a point in some sort of "word" space. This is paramount for a number of reasons but the most intuitive reason is that perhaps there actually exists some N-dimensional space (such that $N \ll 13$ million) that is sufficient to encode all semantics of my language. Each dimension would encode some meaning that I transfer using speech. For instance, semantic dimensions might indicate tense (past vs. present vs. future), count (singular vs. plural), and gender (masculine vs. feminine).

So let's dive into my first word vector and arguably the most simple, the one-hot vector: Represent every word as an $R(\text{dimensions} : |V| \times 1)$ vector with all 0s and one 1 at the index of that word in the sorted english language. In this notation, $|V|$ is the size of my vocabulary. Word vectors in this type of encoding would appear as the following:

$$w^{aardvark} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, w^a = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, w^{at} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots w^{zebra} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

I represent each word as a completely independent entity. As I previously discussed, this word representation does not give us directly any notion of similarity. For instance,

$$(w^{hotel})^T w^{motel} = (w^{hotel})^T w^{cat} = 0$$

So maybe I can try to reduce the size of this space from $R(\text{dimensions} : |V|)$ to something smaller and thus find a subspace that encodes the relationships between words. Common Bag of Words (CBOW) model is an example of one such model and is quite commonly used for text generation in Natural Language Processing (NLP).

One approach is to treat {"The", "cat", "over", "the", "puddle"} as a context and from these words, be able to predict or generate the center word "jumped". This type of model I call a Continuous Bag of Words (CBOW) Model.

The setup is largely the same but I essentially swap my x and y i.e. x in the CBOW are now y and vice-versa. The input one hot vector (center word) I will represent with an x (since there is only one). And the output vectors as $y(j)$. I define V and U the same as in CBOW. I breakdown the way this model works in these 6 steps:

1. I generate my one hot input vector x
2. I get my embedded word vectors for the context $vc = Vx$
3. Since there is no averaging, just set $v^{\wedge} = vc$?
4. Generate $2m$ score vectors, $uc-m, \dots, uc-1, uc+1, \dots, uc+m$ using $u = Uvc$
5. Turn each of the scores into probabilities, $y = \text{softmax}(u)$
6. I desire my probability vector generated to match the true probabilities which is $y(c-m), \dots, y(c-1), y(c+1), \dots, y(c+m)$, the one hot vectors of the actual output.

As in CBOW, I need to generate an objective function for us to evaluate the model. A key difference here is that I invoke a Naive Bayes assumption to break out the probabilities. If you have not seen this before, then simply put, it is a strong (naive) conditional independence assumption. In other words, given the center word, all output words are completely independent.

$$\begin{aligned}
\text{minimize } J &= -\log P(w_{c-m}, \dots, w_{c-1}, w_{c+1}, \dots, w_{c+m} | w_c) \\
&= -\log \prod_{j=0, j \neq m}^{2m} P(w_{c-m+j} | w_c) \\
&= -\log \prod_{j=0, j \neq m}^{2m} P(u_{c-m+j} | v_c) \\
&= -\log \prod_{j=0, j \neq m}^{2m} \frac{\exp(u_{c-m+j}^T v_c)}{\sum_{k=1}^{|V|} \exp(u_k^T v_c)} \\
&= - \sum_{j=0, j \neq m}^{2m} u_{c-m+j}^T v_c + 2m \log \sum_{k=1}^{|V|} \exp(u_k^T v_c)
\end{aligned}$$

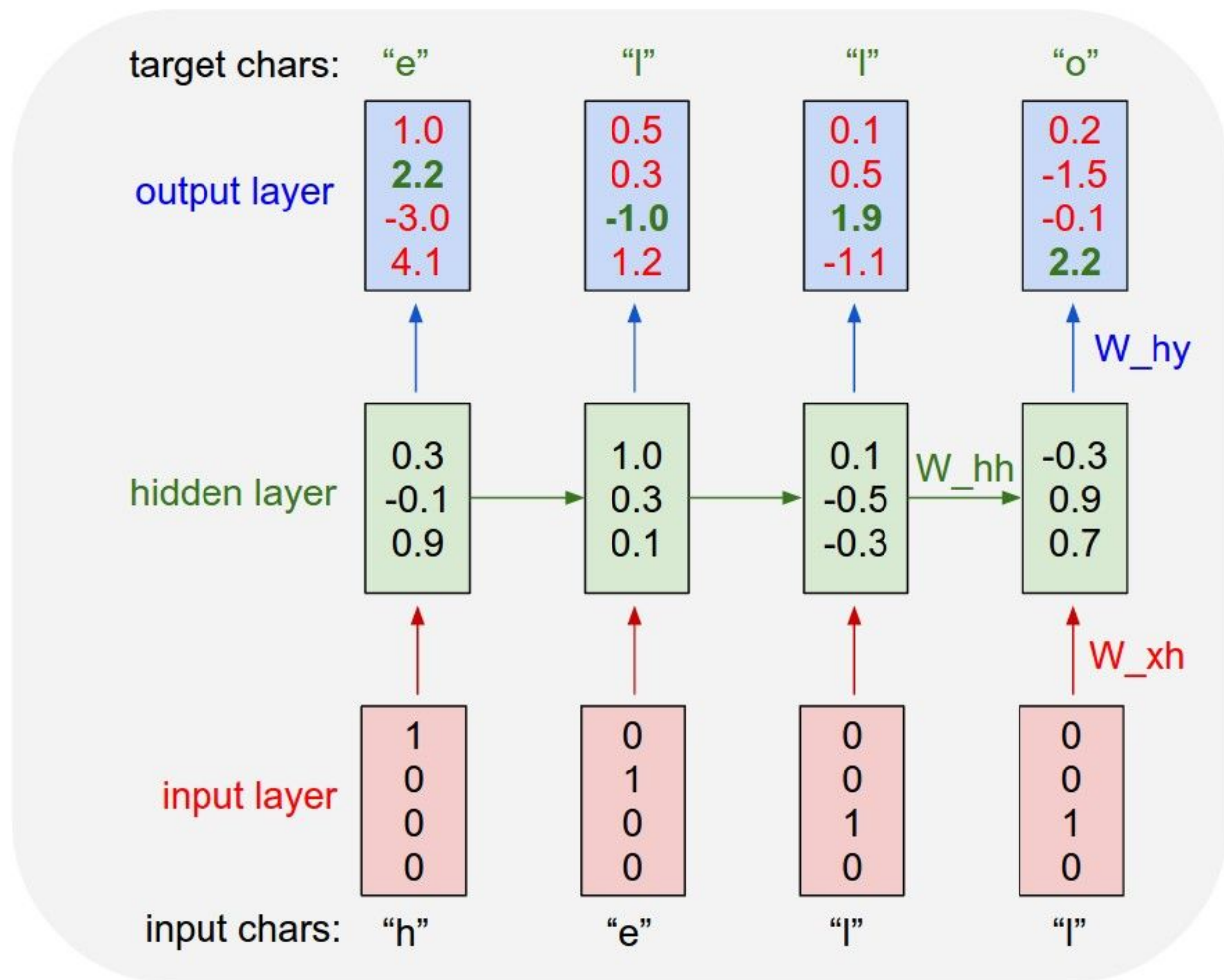
With this objective function, I can compute the gradients with respect to the unknown parameters and at each iteration update them via Stochastic Gradient Descent.

2. Text Generation using Character LSTM

I'll train RNN character-level language models. That is, I'll give the RNN a huge chunk of text and ask it to model the probability distribution of the next character in the sequence given a sequence of previous characters. This will then allow us to generate new text one character at a time.

As a working example, suppose I only had a vocabulary of fmy possible letters "helo", and wanted to train an RNN on the training sequence "hello". This training sequence is in fact a smyce of 4 separate training examples: 1. The probability of "e" should be likely given the context of "h", 2. "l" should be likely in the context of "he", 3. "l" should also be likely given the context of "hel", and finally 4. "o" should be likely given the context of "hell".

Concretely, I will encode each character into a vector using 1-of-k encoding (i.e. all zero except for a single one at the index of the character in the vocabulary), and feed them into the RNN one at a time with the step-function. I will then observe a sequence of 4-dimensional output vectors (one dimension per character), which I interpret as the confidence the RNN currently assigns to each character coming next in the sequence. Here's a diagram:



An example RNN with 4-dimensional input and output layers, and a hidden layer of 3 units (neurons). This diagram shows the activations in the forward pass when the RNN is fed the characters "hell" as input. The output layer contains confidences the RNN assigns for the next character (vocabulary is "h,e,l,o"); I want the green numbers to be high and red numbers to be low.

For example, I see that in the first time step when the RNN saw the character "h" it assigned confidence of 1.0 to the next letter being "h", 2.2 to letter "e", -3.0 to "l", and 4.1 to "o". Since in my training data (the string "hello") the next correct character is "e", I would like to increase its confidence (green) and decrease the confidence of all other letters (red). Similarly, I have a desired target character at every one of the 4 time steps that I'd like the network to assign a greater confidence to. Since the RNN consists entirely of differentiable operations I can run the backpropagation algorithm (this is just a recursive application of the chain rule from calculus) to figure out in what direction I should adjust every one of its weights to increase the scores of the correct targets

(green bold numbers). I can then perform a *parameter update*, which nudges every light a tiny amount in this gradient direction. If I re-feed the same inputs to the RNN after the parameter update I would find that the scores of the correct characters (e.g. “e” in the first time step) would be slightly higher (e.g. 2.3 instead of 2.2), and the scores of incorrect characters would be slightly lower. I then repeat this process over and over many times until the network converges and its predictions are eventually consistent with the training data in that correct characters are always predicted next.

A more technical explanation is that I use the standard Softmax classifier (also commonly referred to as the cross-entropy loss) on every output vector simultaneously. The RNN is trained with mini-batch Stochastic Gradient Descent and I like to use RMSProp or Adam (per-parameter adaptive learning rate methods) to stabilize the updates.

Notice also that the first time the character “l” is input, the target is “l”, but the second time the target is “o”. The RNN therefore cannot rely on the input alone and must use its recurrent connection to keep track of the context to achieve this task.

At test time, I feed a character into the RNN and get a distribution over what characters are likely to come next. I sample from this distribution, and feed it right back in to get the next letter. This process is repeated and text sampling is done. I have written the character-level LSTM in PyTorch.

3. Character Label Generation from ASL Images using CNN

The ASL data set images were fed resized to 32x32 images and were fed into a modified AlexNet, the simplest of standard baseline CNN architectures, which has 24 output neurons same in number as the number of classes in the ASL database. I used transfer learning to train an AlexNet, already trained on the ImageNet dataset, on the ASL dataset. I tried two versions of transfer learning: training only the last fully connected layer and training both fully connected layers at the end. I obtained the best results by training the last three fully connected layers since they contain the majority of the trainable parameters. I will describe transfer learning here in detail.

Transfer Learning -

Chen et al argue that a good training data must have width and depth, i.e. the dataset must contain a large number of subject classes for training and each individual class must have a good number of images in order for the network to train effectively. The ASL dataset from the Exeter university suffers severely from the problem of small depth.

Chen et al have resorted to the creation of a new dataset namely, Wide and Deep Reference Net (WDRef) in order to cope with this problem and later hope to transfer the learning acquired by the training of this model to the ASL Dataset. Essentially, the model trained on WDRef learns a high dimensional representation of faces towards the end of the network, which ASL did not allow due to its lack of depth. Training on ASL would hence be more convenient for the network, given that it has already learnt how to represent images from the previous ImageNet training.

I switched off training of all the other layers apart from the last few layers on which I am applying transfer learning by stopping the flow of gradients into those layers. This also speeds up training to a large extent.

E. FINAL RESULTS

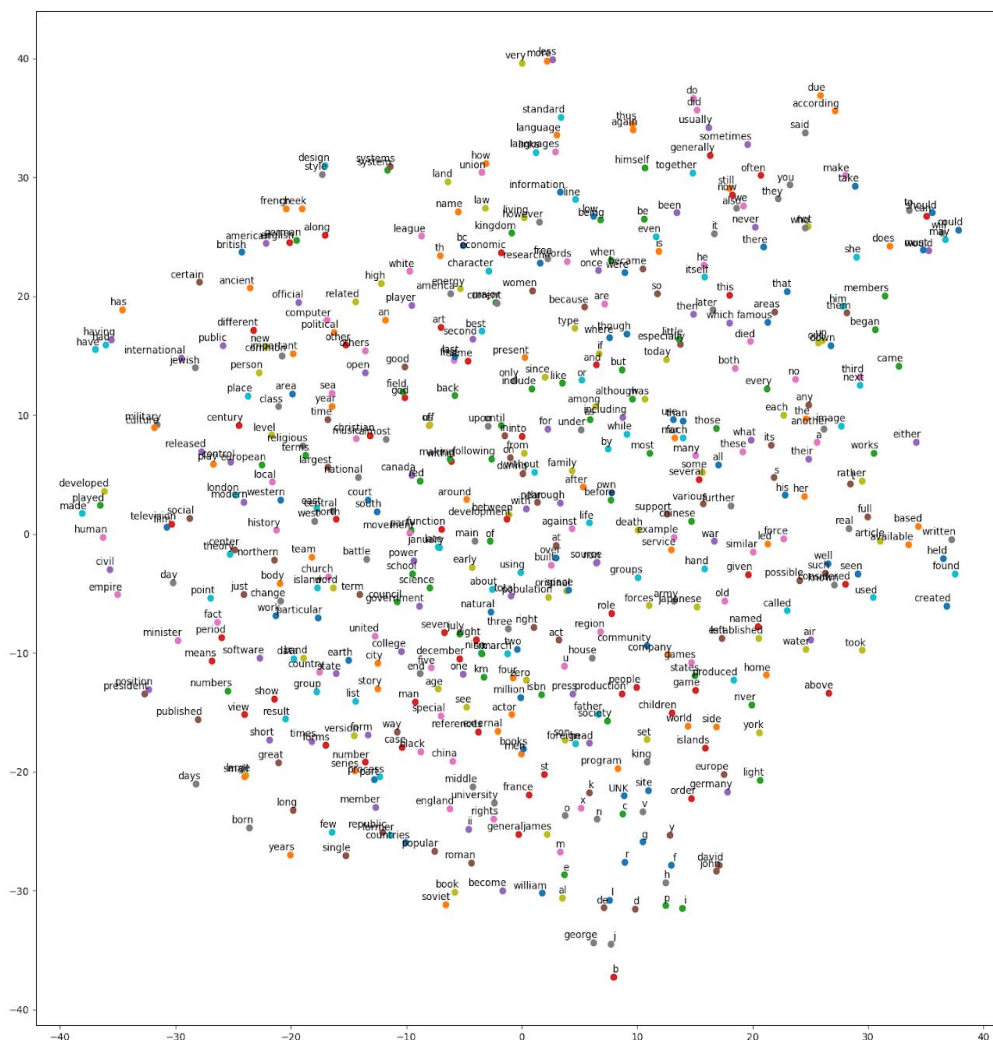


Fig. 1 - t-SNE plot of language model using CBOW

t-distributed stochastic neighbor embedding (t-SNE) is a machine learning algorithm for dimensionality reduction developed by Geoffrey Hinton and Laurens van der Maaten. It is a nonlinear dimensionality reduction technique that is particularly well-suited for embedding high-dimensional data into a space of two or three dimensions, which can then be visualized in a scatter plot. Specifically, it models each

high-dimensional object by a two- or three-dimensional point in such a way that similar objects are modeled by nearby points and dissimilar objects are modeled by distant points.

The t-SNE algorithm comprises two main stages. First, t-SNE constructs a probability distribution over pairs of high-dimensional objects in such a way that similar objects have a high probability of being picked, whilst dissimilar points have an extremely small probability of being picked. Second, t-SNE defines a similar probability distribution over the points in the low-dimensional map, and it minimizes the Kullback–Leibler divergence between the two distributions with respect to the locations of the points in the map. Note that whilst the original algorithm uses the Euclidean distance between objects as the base of its similarity metric, this should be changed as appropriate.

t-SNE has been used in a wide range of applications, including computer security research, music analysis, cancer research, bioinformatics, and biomedical signal processing. It is often used to visualize high-level representations learned by an artificial neural network.

Below are the graphs for training accuracy and validation accuracy of training and validation errors for my LSTM text generation model on the Shakespearian extract dataset.

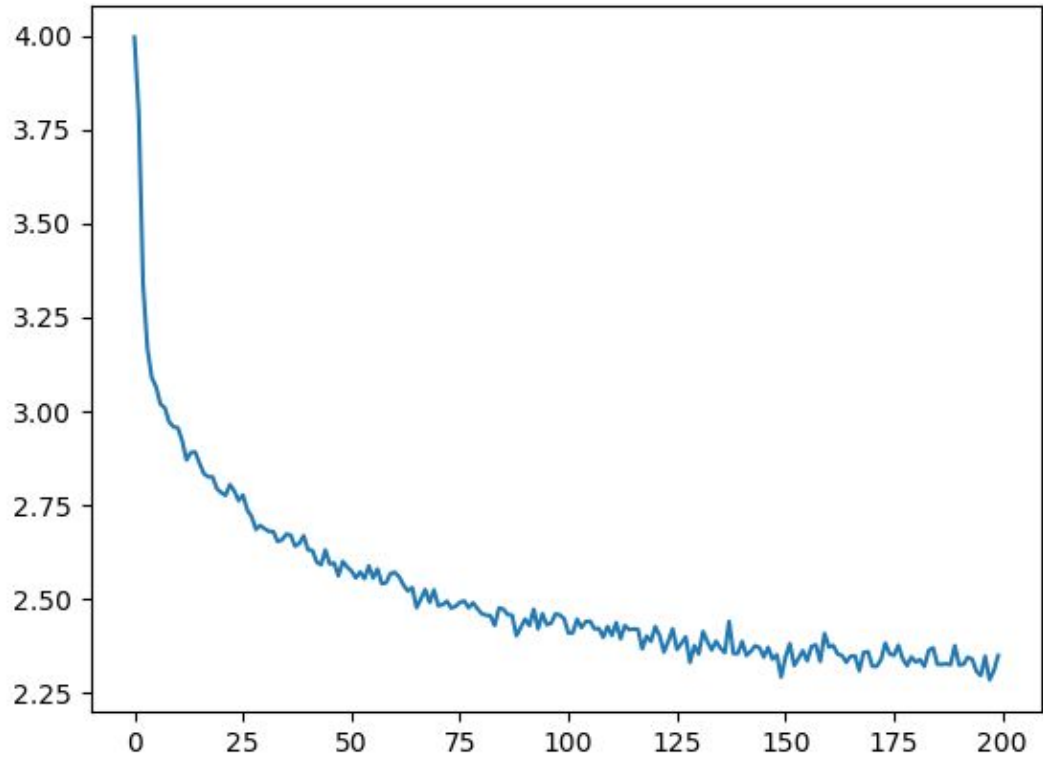


Fig. 2 - Training Error with Epochs for SNN

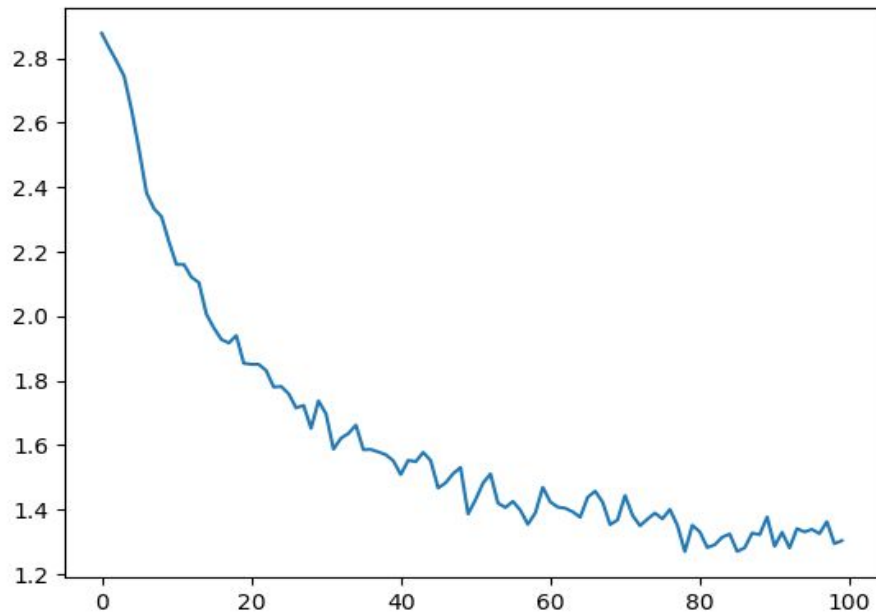


Fig. 2 - Validation Error with Epochs for SNN

Below are the graphs for training accuracy and validation accuracy of training and validation errors for my fine tuned AlexNet model with the last FC layer and the last two FC layers trained on the Exter University ASL dataset.

Training for ASL Dataset with 2 FC layer fine-tuned

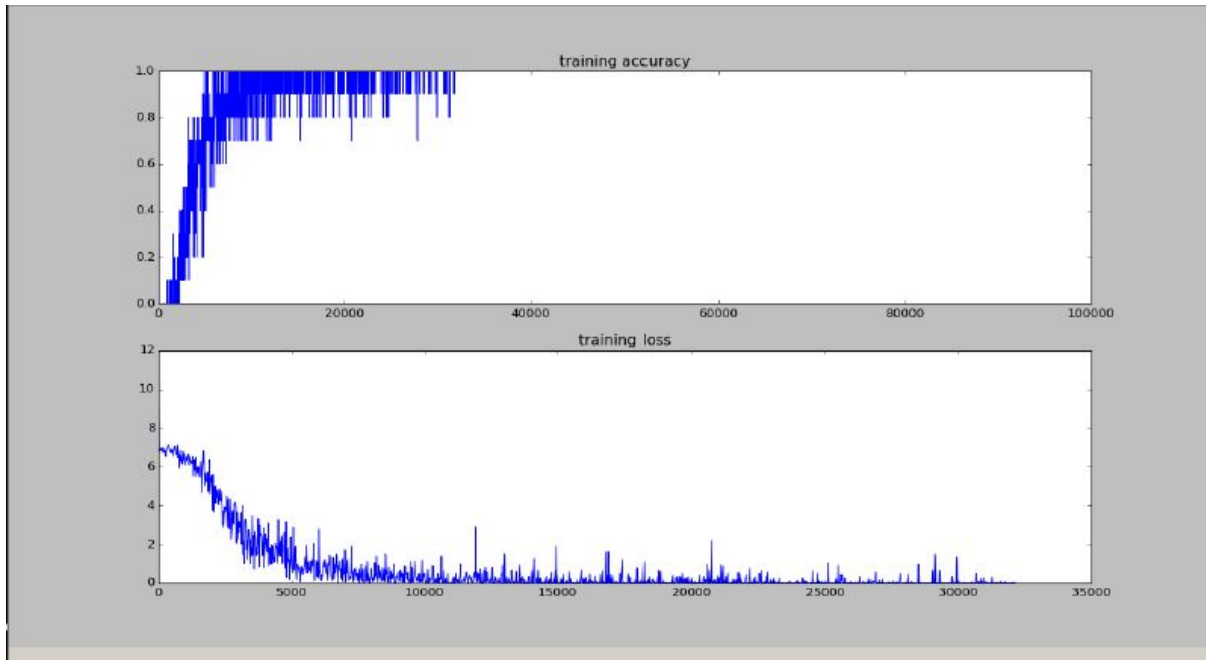


Fig. 4: Training Error and Accuracy on ASL Dataset with last 2 FC layers fine-tuned

Validation for ASL Dataset with 2 FC layer fine-tuned

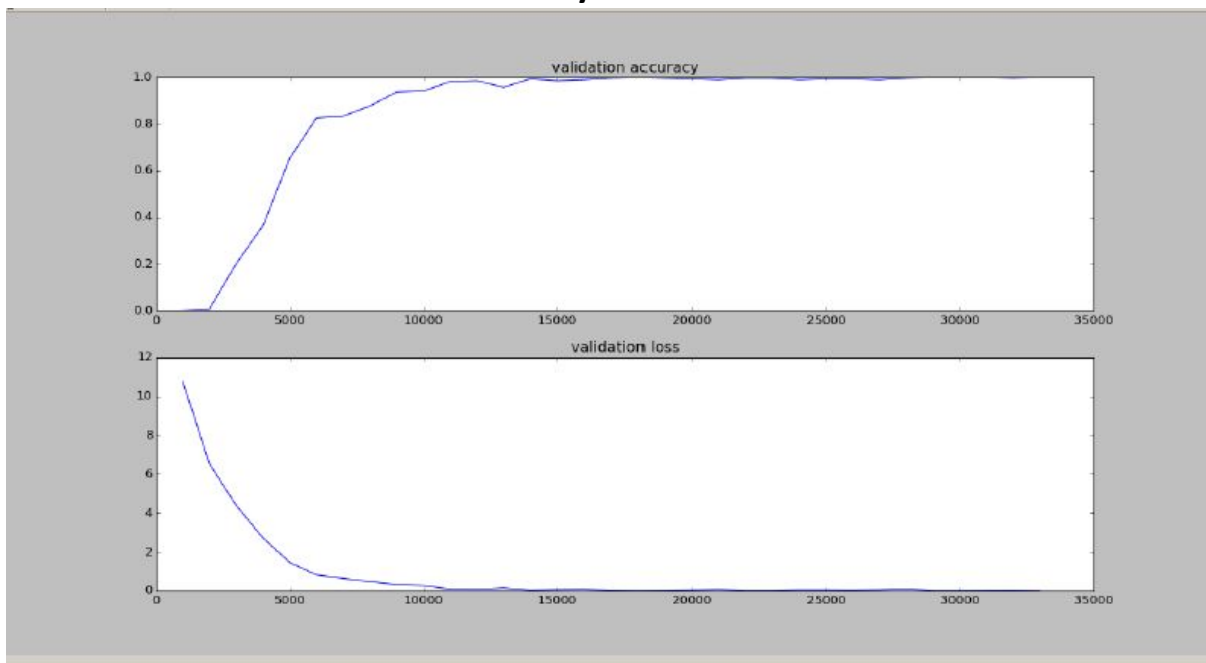


Fig. 5: Validation Error and Accuracy on ASL Dataset with last 2 FC layers fine-tuned

Training for ASL Dataset with 1 FC layer fine-tuned

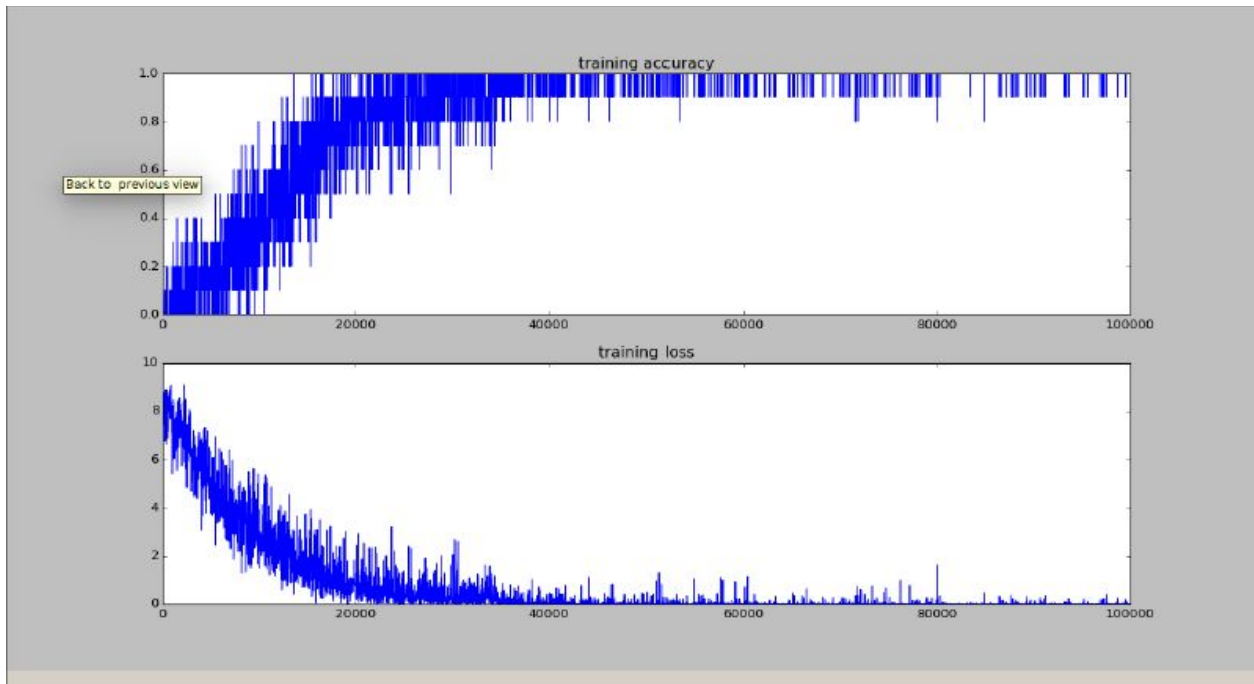


Fig. 6: Training Error and Accuracy on ASL Dataset with last 1 FC layers fine-tuned

Validation for ASL Dataset with 1 FC layer fine-tuned

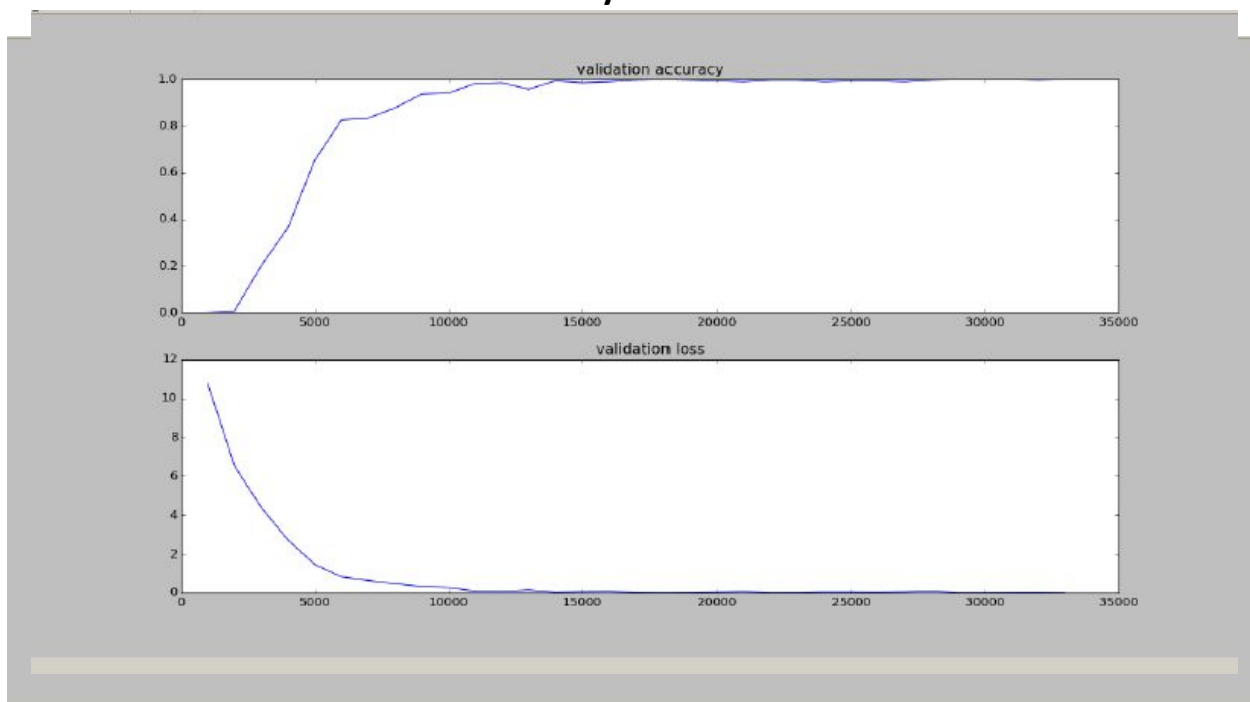


Fig. 7: Validation Error and Accuracy on ASL Dataset with last 1 FC layers fine-tuned

The final accuracy values for validation are as follows -

- Top-5 accuracy in percentage - 74.772
- Top-1 accuracy in percentage - 93.853

Please refer to the figure below.

```

abc-S@nik: ~/Desktop/LOP17_5M/ImageNet
Epoch: [942][174/180]   Time 1.357 (1.278)   Data 0.000 (0.006)   Loss 2.1374 (2.6478)   Prec@1 71.875 (73.400)   Prec@5 94.141 (92.536)
Epoch: [942][175/180]   Time 1.268 (1.278)   Data 0.000 (0.006)   Loss 1.7327 (2.6426)   Prec@1 80.859 (73.442)   Prec@5 95.312 (92.551)
Epoch: [942][176/180]   Time 1.217 (1.278)   Data 0.000 (0.006)   Loss 2.3592 (2.6410)   Prec@1 71.094 (73.429)   Prec@5 93.750 (92.558)
Epoch: [942][177/180]   Time 1.209 (1.277)   Data 0.000 (0.006)   Loss 3.0739 (2.6434)   Prec@1 71.094 (73.416)   Prec@5 89.062 (92.539)
Epoch: [942][178/180]   Time 1.335 (1.278)   Data 0.000 (0.005)   Loss 2.7524 (2.6440)   Prec@1 73.438 (73.416)   Prec@5 91.016 (92.530)
Epoch: [942][179/180]   Time 0.081 (1.275)   Data 0.000 (0.005)   Loss 2.0612 (2.6420)   Prec@1 77.987 (73.431)   Prec@5 94.340 (92.536)
Test: [0/39]   Time 2.117 (2.117)   Loss 1.4082 (1.4082)   Prec@1 75.391 (75.391)   Prec@5 91.797 (91.797)
Test: [1/39]   Time 1.213 (1.665)   Loss 1.2785 (1.3823)   Prec@1 84.766 (80.078)   Prec@5 94.141 (92.569)
Test: [2/39]   Time 1.195 (1.599)   Loss 1.2911 (1.3519)   Prec@1 88.281 (82.812)   Prec@5 96.094 (94.810)
Test: [3/39]   Time 1.190 (1.429)   Loss 0.2799 (1.0839)   Prec@1 94.141 (85.645)   Prec@5 98.438 (95.117)
Test: [4/39]   Time 1.205 (1.384)   Loss 1.0182 (1.0708)   Prec@1 82.422 (85.000)   Prec@5 98.047 (95.703)
Test: [5/39]   Time 1.236 (1.359)   Loss 4.0368 (1.5651)   Prec@1 53.906 (79.818)   Prec@5 87.891 (94.481)
Test: [6/39]   Time 1.205 (1.337)   Loss 2.5904 (1.7116)   Prec@1 68.359 (78.181)   Prec@5 94.141 (94.364)
Test: [7/39]   Time 1.229 (1.325)   Loss 2.0215 (1.7593)   Prec@1 81.644 (78.613)   Prec@5 94.922 (94.434)
Test: [8/39]   Time 1.203 (1.312)   Loss 3.7164 (3.0688)   Prec@1 55.078 (75.998)   Prec@5 86.719 (90.576)
Test: [9/39]   Time 1.253 (1.306)   Loss 2.9841 (2.0703)   Prec@1 64.062 (74.805)   Prec@5 90.234 (93.242)
Test: [10/39]   Time 1.344 (1.309)   Loss 3.8771 (2.2345)   Prec@1 67.578 (74.148)   Prec@5 89.844 (92.933)
Test: [11/39]   Time 1.234 (1.303)   Loss 1.5915 (2.1810)   Prec@1 75.000 (74.219)   Prec@5 94.141 (93.034)
Test: [12/39]   Time 1.333 (1.305)   Loss 1.5548 (2.1328)   Prec@1 72.266 (74.069)   Prec@5 96.484 (93.299)
Test: [13/39]   Time 1.243 (1.301)   Loss 2.1120 (2.1313)   Prec@1 66.406 (73.521)   Prec@5 97.656 (93.610)
Test: [14/39]   Time 1.219 (1.295)   Loss 1.7060 (2.1030)   Prec@1 74.009 (73.594)   Prec@5 96.484 (93.802)
Test: [15/39]   Time 1.329 (1.300)   Loss 0.9061 (2.0319)   Prec@1 85.156 (74.310)   Prec@5 97.656 (94.019)
Test: [16/39]   Time 1.226 (1.296)   Loss 1.3393 (1.9912)   Prec@1 79.297 (74.609)   Prec@5 95.703 (94.118)
Test: [17/39]   Time 1.245 (1.293)   Loss 2.7697 (2.0344)   Prec@1 63.281 (73.908)   Prec@5 96.094 (94.227)
Test: [18/39]   Time 1.203 (1.289)   Loss 3.5587 (2.1146)   Prec@1 40.047 (72.615)   Prec@5 96.484 (94.346)
Test: [19/39]   Time 1.205 (1.284)   Loss 1.0719 (2.0625)   Prec@1 83.594 (73.164)   Prec@5 94.531 (94.355)
Test: [20/39]   Time 1.194 (1.280)   Loss 1.0537 (2.0145)   Prec@1 84.375 (73.698)   Prec@5 96.094 (94.438)
Test: [21/39]   Time 1.223 (1.277)   Loss 1.2290 (1.9788)   Prec@1 85.547 (74.237)   Prec@5 95.312 (94.478)
Test: [22/39]   Time 1.305 (1.279)   Loss 1.3078 (1.9406)   Prec@1 80.469 (74.507)   Prec@5 98.047 (94.607)
Test: [23/39]   Time 1.267 (1.278)   Loss 2.6592 (1.9792)   Prec@1 70.312 (74.333)   Prec@5 94.922 (94.670)
Test: [24/39]   Time 1.209 (1.275)   Loss 4.2924 (2.0717)   Prec@1 50.203 (73.688)   Prec@5 87.891 (94.406)
Test: [25/39]   Time 1.204 (1.273)   Loss 2.4225 (2.0852)   Prec@1 78.516 (73.873)   Prec@5 92.969 (94.351)
Test: [26/39]   Time 1.239 (1.271)   Loss 0.4053 (2.0230)   Prec@1 93.750 (74.609)   Prec@5 99.219 (94.531)
Test: [27/39]   Time 1.400 (1.276)   Loss 0.3764 (1.9642)   Prec@1 92.188 (75.237)   Prec@5 98.828 (94.685)
Test: [28/39]   Time 1.290 (1.276)   Loss 0.4005 (1.9102)   Prec@1 92.969 (75.849)   Prec@5 98.828 (94.828)
Test: [29/39]   Time 1.240 (1.275)   Loss 0.1934 (1.8530)   Prec@1 97.266 (76.562)   Prec@5 99.009 (94.987)
Test: [30/39]   Time 1.229 (1.275)   Loss 1.1935 (1.8317)   Prec@1 83.203 (76.777)   Prec@5 95.703 (95.010)
Test: [31/39]   Time 1.337 (1.277)   Loss 3.2430 (1.8750)   Prec@1 60.938 (76.282)   Prec@5 85.156 (94.702)
Test: [32/39]   Time 1.244 (1.276)   Loss 1.3984 (1.8613)   Prec@1 78.516 (76.349)   Prec@5 97.656 (94.792)
Test: [33/39]   Time 1.292 (1.277)   Loss 1.8005 (1.8595)   Prec@1 73.828 (76.275)   Prec@5 94.922 (94.795)
Test: [34/39]   Time 1.241 (1.276)   Loss 2.9468 (1.8906)   Prec@1 68.750 (76.060)   Prec@5 88.672 (94.621)
Test: [35/39]   Time 1.299 (1.276)   Loss 4.7817 (1.9789)   Prec@1 53.516 (75.434)   Prec@5 83.984 (94.325)
Test: [36/39]   Time 1.241 (1.275)   Loss 4.7345 (2.0456)   Prec@1 57.812 (74.958)   Prec@5 84.706 (94.067)
Test: [37/39]   Time 1.192 (1.273)   Loss 3.0680 (2.0725)   Prec@1 68.359 (74.784)   Prec@5 89.062 (93.935)
Test: [38/39]   Time 0.637 (1.257)   Loss 2.6639 (2.0803)   Prec@1 73.846 (74.772)   Prec@5 87.092 (93.853)
* Prec@1 74.772 Prec@5 93.853
Epoch: [943][0/180]   Time 2.048 (2.048)   Data 0.059 (0.859)   Loss 2.5219 (2.5219)   Prec@1 73.438 (73.438)   Prec@5 94.141 (94.141)
Epoch: [943][1/180]   Time 1.195 (1.622)   Data 0.000 (0.430)   Loss 2.7019 (2.6119)   Prec@1 72.656 (73.047)   Prec@5 94.531 (94.336)
Epoch: [943][2/180]   Time 1.192 (1.478)   Data 0.000 (0.287)   Loss 2.7211 (2.6403)   Prec@1 69.922 (72.005)   Prec@5 93.750 (94.141)
Epoch: [943][3/180]   Time 1.194 (1.407)   Data 0.000 (0.215)   Loss 2.7313 (2.6690)   Prec@1 73.828 (72.461)   Prec@5 94.531 (94.238)
Epoch: [943][4/180]   Time 1.193 (1.364)   Data 0.001 (0.172)   Loss 2.6059 (2.6504)   Prec@1 74.009 (72.891)   Prec@5 96.025 (93.510)
Epoch: [943][5/180]   Time 1.193 (1.336)   Data 0.000 (0.144)   Loss 3.1905 (2.7454)   Prec@1 73.047 (72.917)   Prec@5 89.844 (92.904)
Epoch: [943][6/180]   Time 1.193 (1.315)   Data 0.000 (0.123)   Loss 3.9660 (2.9198)   Prec@1 69.141 (72.377)   Prec@5 89.062 (92.355)
Epoch: [943][7/180]   Time 1.190 (1.301)   Data 0.000 (0.100)   Loss 3.5384 (2.9971)   Prec@1 68.359 (71.875)   Prec@5 91.016 (92.188)
Epoch: [943][8/180]   Time 1.200 (1.289)   Data 0.000 (0.096)   Loss 2.4834 (2.9400)   Prec@1 74.219 (72.135)   Prec@5 89.453 (91.884)
Epoch: [943][9/180]   Time 1.193 (1.280)   Data 0.001 (0.086)   Loss 3.4876 (2.9948)   Prec@1 69.922 (71.914)   Prec@5 90.234 (91.719)
Epoch: [943][10/180]   Time 1.197 (1.272)   Data 0.000 (0.079)   Loss 2.6278 (2.9614)   Prec@1 70.703 (71.884)   Prec@5 92.969 (91.832)
Epoch: [943][11/180]   Time 1.193 (1.260)   Data 0.000 (0.072)   Loss 4.0925 (3.0557)   Prec@1 69.531 (71.615)   Prec@5 90.625 (91.732)
Epoch: [943][12/180]   Time 1.195 (1.260)   Data 0.000 (0.067)   Loss 3.2200 (3.0608)   Prec@1 69.922 (71.404)   Prec@5 92.188 (91.767)
Epoch: [943][13/180]   Time 1.210 (1.257)   Data 0.001 (0.062)   Loss 2.6322 (3.0377)   Prec@1 75.391 (71.763)   Prec@5 93.750 (91.908)
Epoch: [943][14/180]   Time 1.208 (1.254)   Data 0.000 (0.058)   Loss 3.0095 (3.0359)   Prec@1 75.000 (71.979)   Prec@5 89.844 (91.771)
Epoch: [943][15/180]   Time 1.310 (1.257)   Data 0.000 (0.054)   Loss 3.4840 (3.0639)   Prec@1 69.531 (71.826)   Prec@5 90.625 (91.699)

```

Also, a user-friendly interfacing code was written to input variable number of individual ASL images to create the primer text for the sentence completion. Kindly refer to the videos.