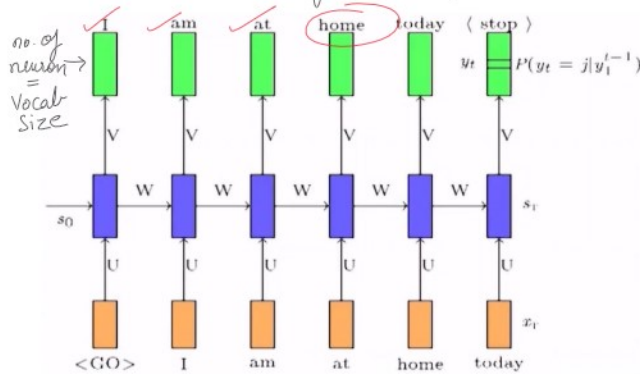


# RNN NLP Part 4

11 October 2023 11:46 AM

Let's revisit the task of auto completion, where in the task is to predict the next word given a word as input.



<GO>, <stop> → are special words to tell/signify the start and end respectively for text completion.

Vocab = { I, at, today, home, am }

Q. By looking at the vocab, tell how many neurons will be present in the outer layer?

Ans 5 neurons eg

j=0	0.1	I
j=1	0.1	at
j=2	0.2	today
j=3	0.1	home
j=4	0.5	am

{ Softmax Activation }

← find prediction as the probability is the highest

So at any time step 't', we want to compute:

$y_t \rightarrow$  output at  $t^{\text{th}}$  time step

$\arg \max_{j \in \mathcal{V}} P(y_t = j | y_{t-1}, y_{t-2}, \dots, y_1)$

max will only give the highest probability, arg max will give the value of 'j' for which the probability is highest.

j ∈ V  
↓  
word

Vocabulary

Eg  $P(y_t = \text{home} | \text{at}, \text{am}, \text{I})$

What is the probability that at time step 't' the word is j { j=0,1,2,3,4 } given all the previous words  $y_{t-1}, y_{t-2}, \dots, y_1$

and using RNN we are trying to model the above probability.

So using RNN,

{ Modelling choice : like in Machine Learning we assume that "y" is a linear composition of inputs, and using this we try to find the best model

$$P(y_t = j | y_{t-1}, y_{t-2}, \dots, y_1) = [\text{softmax}(v s_t + c)]_j$$

$s_t$  captures all the information until time step 't'.

Softmax in last layer.

Because we want to find probability of  $j^{\text{th}}$  neuron, which contains the  $j^{\text{th}}$  word.

finds this 'j' value

this entire thing

Vector

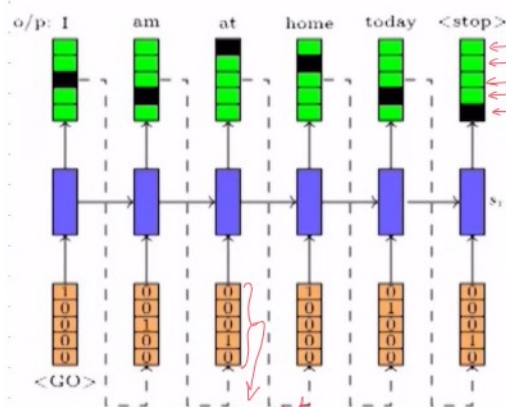
This vector is trying to predict

j=0	0.1	I
j=1	0.1	at
j=2	0.2	today
j=3	0.1	home
j=4	0.5	am

Now, we know that RNN has a limit and cannot just keep on accumulating information without the drawbacks

Here the RNN is trying to learn a probability vector over the entire vocabulary and it tries to model all that into the hidden state  $s_t$ .

## Input for Network



Every position is dedicated to a word and the word with highest probability is selected

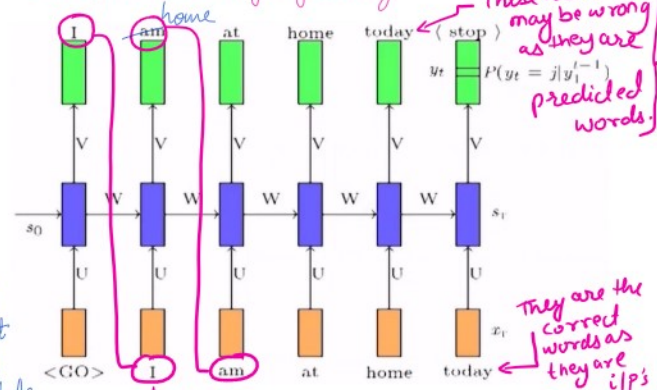
→ Let's say during the training time the word "home" is predicted in place of "am", then it doesn't mean we will put "home" in the next time step. "home" will

have a loss and will update its parameters during backpropagation to get the correct prediction but won't be used in next time step.

one hot encoding of words  
(option to use GloVe, word2vec embeddings)

However at test time 't', whatever the prediction it will be used in the next time step.

\* During training phase there might be a case where the prediction might go wrong

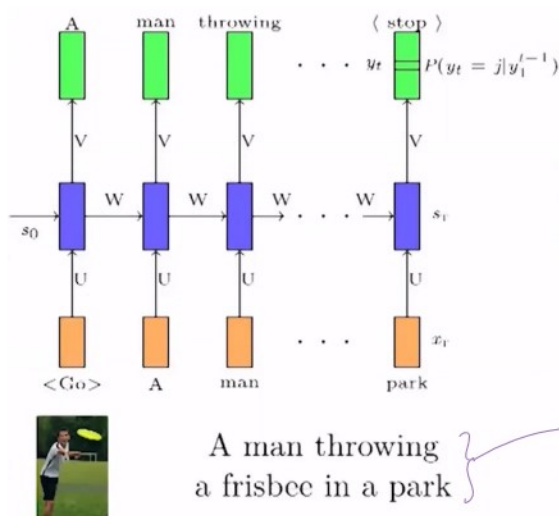


These words may be wrong as they are predicted words.

So there is no reason for these 2 words to be exactly same.

## Image Captioning

Task: Given an image as input, the task is to put a caption on the image describing the image.



In layman terms, we want to generate a sentence given an image.

We are interested in

$$P(y_t = j | y_{t-1}, y_{t-2}, \dots, y_1, I)$$

So basically, given the image information and all the previous words that have been predicted, what is the probability of next word.

Now the reason for doing this is because we are not going to predict the entire caption in one go, we do so word by word. because the natural language has a sequential structure i.e. the sequence of the word matters. for eg. a park throwing a frisbee in a man, has the same permutation of the vocab but is absolutely wrong in meaning as the sequence of the words is not correct. Basically whatever has been said / seen previously decides / determines the next word.

Q What kind of architecture can be used for the above task?

Ans Since caption generation problem has a sequential/temporal nature → RNN can be used

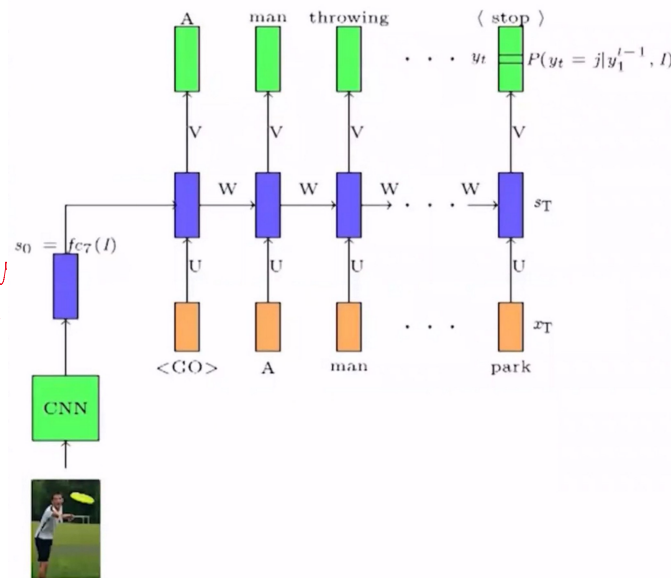
But, the input consists of an image: so to learn good feature representations → CNN can be used.

so, input → CNN  
output → RNN

**P.T.O.**

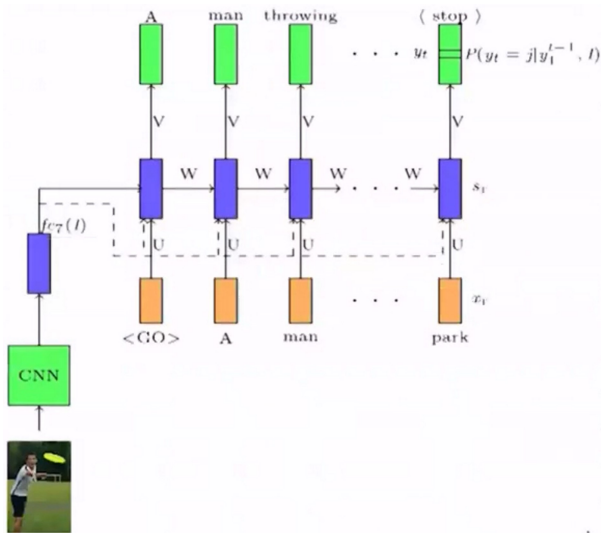
Since, CNN architectures are good for images, we utilize them to learn important features from the image and then pass it on to the RNN architecture/model

Rather than passing those features to CNN, it sends to RNN  
CNN learns the features of the input image



Option-1

→ Output of CNN is directly passed to time step-1



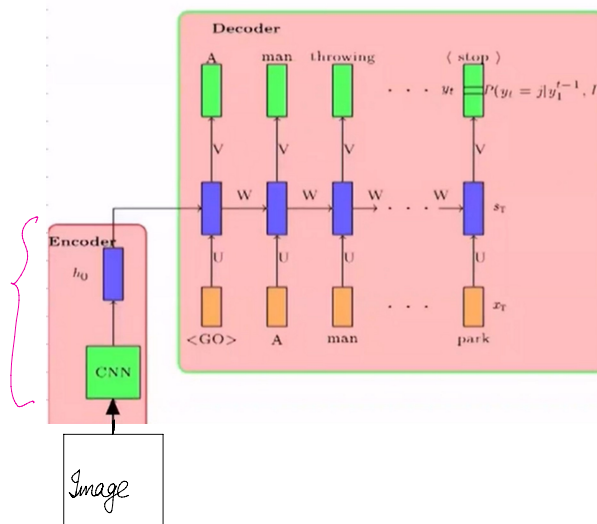
→ option-2

→ Output of CNN is directly passed to all time steps

So in general the architecture will look like ↓

### Encoder Decoder Architecture

learn the important and good representation of the input and pass it to decoder



The decoder receives the input from encoder and produces the output.



In the example discussed above: Encoder  $\rightarrow$  CNN  
 Decoder  $\rightarrow$  RNN } Now this is not all permanent, depending upon the task that we are trying to solve, we would like to use different encoder, decoder architecture.

In image captioning,

$\rightarrow$  A CNN is used to "encode" the image

$\hookrightarrow$  learn good feature representation of the input.

$\rightarrow$  An RNN is used to "decode" a sentence from this encoding.

Task: Image Captioning

Data:  $\{x_i = \text{image}, y_i = \text{caption}\}_{i=1}^N$

$\hookrightarrow$  Refers to a very deep network with many filters and pooling layers

Model: It has two parts: Encoder:  $S_0 = \text{CNN}(x_i)$

Decoder:  $S_t = \text{RNN}(S_{t-1}, y_{t-1})$  } We are using this to make a prediction over the probability of all possible words  
 Previous time Step.      previous predicted word

Parameters of the model:  $u, v, w, b, c$  and all weights and biases of CNN

and then, we train ALL the parameters of Encoder Decoder Together

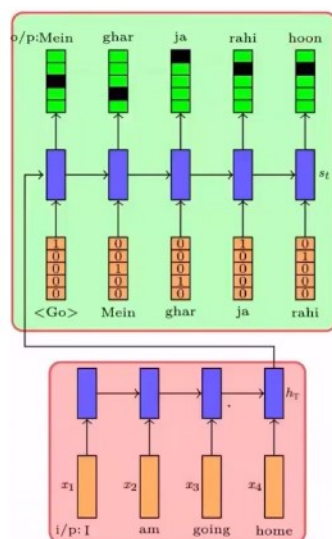
NOTE: The training of Encoder & decoder is not separate, the entire architecture is trained at once.

This type of model is also called as end to end model.

## # Real life Use Case of Encoder Decoder

$\rightarrow$  Machine Translation

Translate english to hindi



Encoder  $\rightarrow$  RNN  $\rightarrow$  To handle the text input  
 Decoder  $\rightarrow$  RNN  $\rightarrow$  To convert & give text output

Decoder { Using the representation, it tries to predict the output }

Encoder { It tries to make a good representation of the input by understanding the syntax & semantics and pass it to decoder }

## → Video Captioning

↳ Video is a sequence of images, and now for that entire sequence we want to put up a caption

↳ since CNN is good for images we use that, but since the images are in a sequence, i.e. the current frame influences the next frame, we put the entire thing in an RNN

⇒ Encoder = (CNN + RNN)

