# Hands on 5

## Transformers

**Task:-**

To Develop a transformer architecture with given configurations from scratch for language translation.

**Dataset:-**

The dataset, contains two data columns, one column has english words/sentences and other one has French words/sentences.

**Configurations:-**

a) A single attention head with 3 encoders and decoders respectively.

b) Multi attention heads of 8 with 3 encoders and decoders respectively.

c) Multi attention heads of 32 with 3 encoders and decoders respectively.

d) Multi attention heads of 8 with 5 encoders and decoders respectively.

a) Configuration - 1

This configuration relies on single mechanism to weight the importance of different words. It has singly struggled to capture complex dependences

But, by increasing number of encoder and decoder layer it has gained capacity to learn hierarchical representances of the input and output sequences.

Results :-    Training - loss = 0.21077
              val - accura =    88.73  %
Plot - 1      Test - accura =   88.60  %
              Test - BLEU- Score = 46.94  %

b) Configuration - 2
   increasing the multiple attention heads allows the model to attend to different part of input sentences simultaneosly.
   But, it has increased computational complexity.

   Here, on comparison with configuration '1', it takes more time to train and on average it gives similar results as of configuration '1'.

**Results:-** Training_loss = 0.17920

(Plot-2)  val_accuracy = 88.70 %

Test_accuracy = 88.60 %

Test_BLEU Score = 46.60

c) Configuration - 3

Multi head attention with 32, should further enhances the model, ability to capture diverse linguistic features and dependences.

But with given configuration of running '200' epochs has led my model to overfit and hence my model val_accuracy decreases in comparison to configuration one and two. Which should not be the ideal case.

Results!-      Training loss = 0.21347

val_accuracy = 88.55 %

(Plot-3)    Test_accuracy = 88.32 %

Test_BLEU Score = 43.37 %

d) Configuration - 4

with 8 multi head attentions and 5 each encoder and decoder blocks has basically enhances the model's ability to capture complex patterns and relationships in the data. The increased depth allows for more expressive representation.

On comparising it with configuration '2' it has given better results with better val_accuracy, test_accuracy. at captures the data complexity more accurately than config-2.

Results:-  Training - loss =  0.14989

(plot-4)  Val - accuracy :  89.91 %.

Test - accuracy :  89.77 %.

Test - BLEU - Score =  46.35

Strength:-

-> Configuration 2 and configuration 4, config.
has l~~ately~~  outperform  configuration
1  due  to  use  of  multi attention
head , enabling  better  capturing
of  contextual  information.

·) Configuration 4  with  increased  encoder
and  decoder  , has  excel  in  capturing
complex  linguistic  structures  and
nuances , especialy  in  longer  sentences.

weakness :-

-) Configuration 3, with  32  attention  head
had  suffered  from  increased computational
complexity, leading  to  longer  training
time  and  higher  resources  allocation.

-> Configuration 4 , while  potentially  more accurate
has  also  suffer  from  increased  complexity
making  it  harder  to  interpret
and  deeply  deploy  in  resource -
-constrained  setup.