



## Attention

y solution because context vectors forget earlier parts of sequence

## Attention Mechanisms

### Soft attention

y over each pixel, rating  $(0, 1)$   
continuous

### Hard

y over patches, rating  $\Sigma 0, 1$   
discrete

## Global Attention

y similar to soft attention

## Local attention

y local hard w/ global soft

y only compute attention over small window

## Self Attention

→ retrieval

given  $q$ , find keys & most similar to  $q$  and return the corresponding values  $v$

$$\text{Attention} = \sum \text{sim}(q, k) \cdot v$$

↑ weight

scaled  
dot product  
as similarity  
function

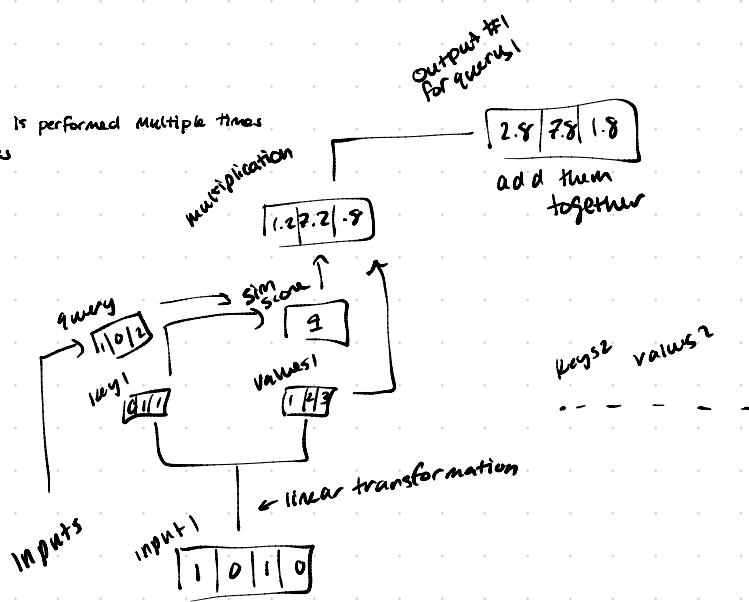
## Multihed self attention

y at every layer, self attention is performed multiple times

y multiple representation subspaces

y diff things/semantics

Head 1      Head 2  
↓            ↓  
POS          syntax



+ Positional Representations

+ non linearities

+ masking: parallelize w/out looking into the future  
used in decoder  
avoid data leakage

## Encoder - Decoder

2017 Vaswani:

translate eng to ita

encode entire sentence  
decoder token by token

### encoder

self attention

feed forward (nonlinear)

layer norm

residual connections b/wn encoder blocks

### decoder

y like encoder but w/ multihead attention also using output of encoders

y masking to make sure not looking into future

→ self attention  
is multihead  
attention from  
prev layers  
of respective  
encoder + decoder

## Advantages

- constant path length b/wn any two pos in seq

y no problem w/ long sequences

y wont loose context

## -parallelization

## Drawbacks

- quadratic time + space

y trying to make linear or quasi-linear

GPT-3

only decoder blocks

$$p(w_t | w_{1:t-1})$$

predict last token given last + tokens

For downstream tasks:

train classifier on last hidden state

## NLI, sentiment Classification

→ use pretrained network as generator

## In context learning

4 few shot settings

BERT

## Transformer architecture

↳ Only encoder  $\Rightarrow$  be bidirectional can't pretrain naive language tasks

Y sees entire text  $\rightarrow$  masked language modeling  
predict words based on context

↳ new training objective      context  
↳ given 2 chunks predict if 2nd follows first

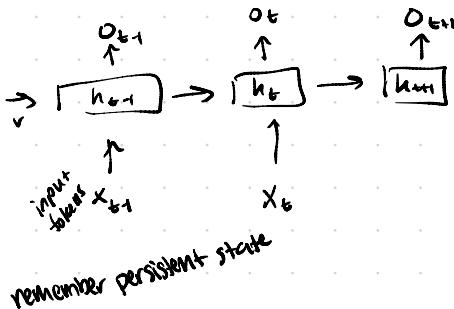
↳ downstream: additional classification layer

4 embeddings from last layer corresponding to special token [EOS]

## 3-Gram Model (1951)

$P(\text{word} \mid \text{two preceding words})$       lookup table

## Recurrent Neural Network (2011)



## Big LSTM (2016)

- ↳ built on top of RNN
- ↳ long term dependencies

## Transformer

- ↳ long term dependencies improved
- ↳ work = regardless of distance

## GPT-2 (2019)

## GPT-3 (2020)

Very big transformer

Auto regressive Generation Model

## Deep learning

### supervised

AlexNet  $\Rightarrow$  classify images to labels  
(2012)  
 $\Rightarrow$  minimal risk

### Unsupervised

$\Rightarrow$  unlabeled  
 $\Rightarrow$  no guarantees