

(1)

# Zhang - CHII - Attention & transform

- (1) Queries & keys & Values
  - (2) Attention Pooling by similarity
  - (3) Attention scoring
  - (4) Bahdanau Attention mechanism
  - (5) multi-head attention
  - (6) Self-Attention & Positional encoding
  - (7) transformer architecture
  - (8) transformers for vision
  - (9) large scale pre-train w/ transform
- 

Early years of deep learning boom were in NLP, convolutional & recurrent

there were plenty of methodological innovations but the architectures stayed mostly the same

ReLus, residual layers, batch norm, dropout, adaptive learning

(ev) (NLP)  
CNN & LSTM style architectures remained state-of-art

today transformers dominates NLP

<sup>NLP</sup>  
given a new task, the ~~dom~~ default approach is to grab a pre-trained transformer

BERT, Electra, RoBERTa, longformer

Adapt outlayers to task

& fine-tune on available data

(2)

transformers are based on the attention mechanism

original an enhancement to encoder-decoder RNNs

Recall that first Seq-2-Seq models the input was compressed into a single fix-length vector

the intuition behind Attention

- ↳ instead of compressing
- ↳ might be better for the decoder to revisit the input seq at every step

Also refer to ~~diff~~ points of the input @ ~~diff~~ points of the decode

Attention = means to dynamically attend to different parts of the input @ each decode step

high-level = encoder would produce a ~~rep~~ representation vector of equal length to the original input seq

decoder receives the input as a context vector — A weighted sum of each of the representations of the input @ each step

weights determine to extent to which a given portion of the input is considered

each step "focuses" on each token

key goal <sup>is</sup> for the weight assigning proc to be differentiable so it can be learned along w/ the weights & biases of a neural network

Attention was very successful @ outperforming the old RNN architectures

↳ Also Attention weighs provides insights of salient features

Vaswani (2017) proposed transformer arch dispensing w/ recurrent connections

↳ instead relying on cleverly arranged attention mechanisms to capture relationships between input & output tokens

Performed well & by 2018 = SOTA NLP

Moreover, same time the main practice became to pretrain large-scale models on enormous generic background corpora to optimize some self-supervised pretrain objective

↳ then to fine-tune using available data downstream

With this pretrain method, the gap between transformers & ~~the~~ original architectures grew

Large pretrained models are known as foundation models

### II. 1 Queries, keys & values

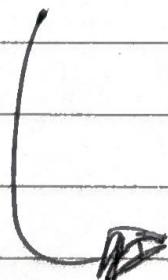
Prev archs such as CNNs & RNNs start inputs of defined size ( $224 \times 224$  e.g.)

Variable size has to be dealt w/ by processing one token @ a time

this can lead to significant problems when the input is of vary size & with vary input quality

in particular, for some long segs it can be difficult for the model to keep track of everything that has been used & generated by the network

### Attention Mechanism



## Attention Mechanism

$$\mathcal{D} = \{ (k_1, v_1), \dots, k_m, v_m \}$$

Define attention as:

$$\text{Att}(q, \mathcal{D}) \sum a(q, k_i) v_i$$

$q$  = query

$a(q, k_i)$  = Scalar attention weights

often referred to as attention pooling

if an  $a$  is high then the model pays particular attention to the terms

↳ this is where the name "attention" comes from

Attention over  $\mathcal{D}$  gets a lin comb of values contained in the database

if  $a(q, k_i) = 1$  then this is like same as a direct db &

normalize weights to ensure they add up to 1

$$a(q, k_i) = \frac{a(q, k_i)}{\sum a(q, k_j)}$$

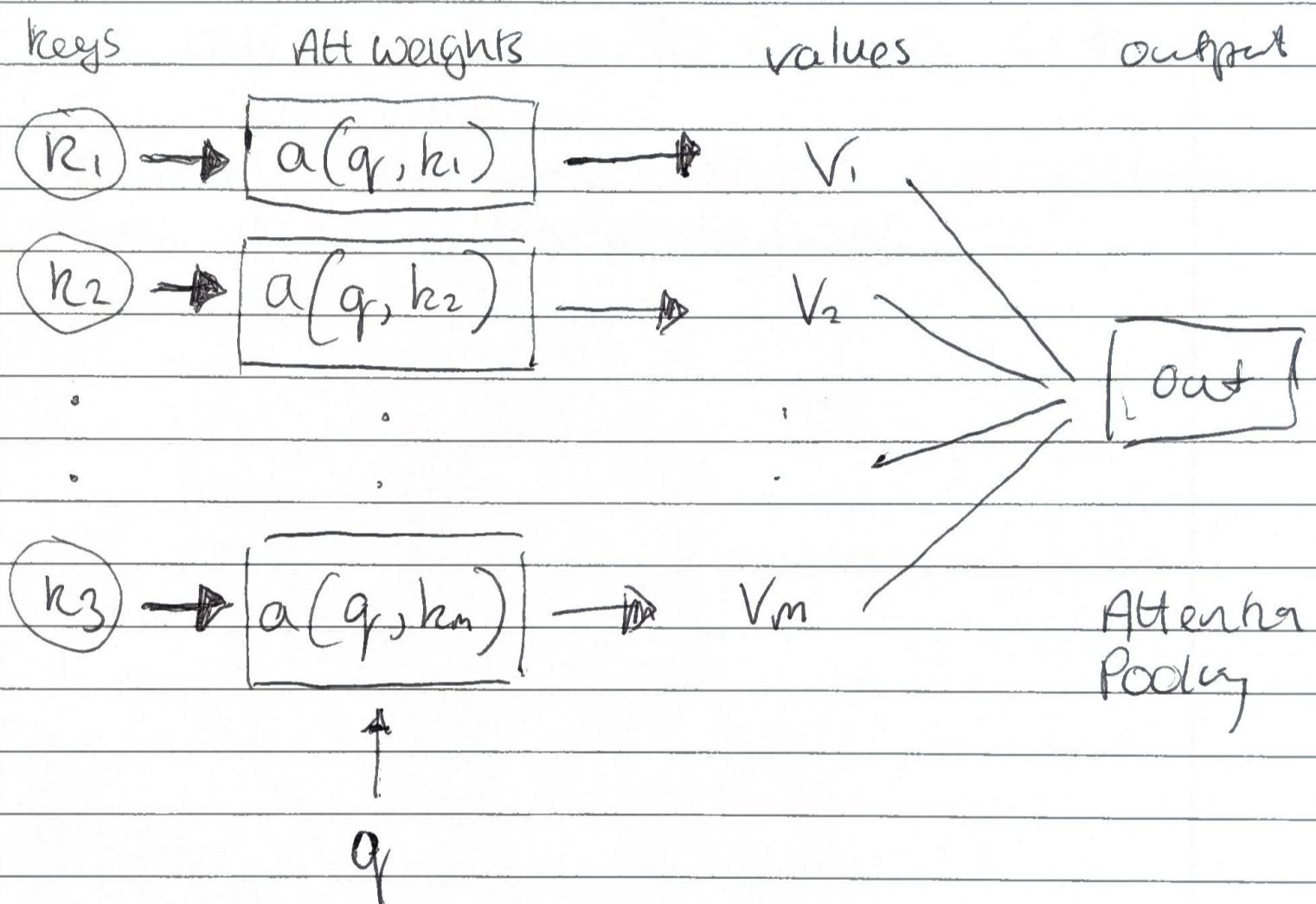
↳ Can also exp to ensure not neg

(6)

this operation is differentiable & never vanish

↳ Desirable properties of a model

Can create a non-Diff model that can be trained via re-enforcement learning



Att Mech computes a linear combination over values  $v_i$

weights are derived according to compact between  $q$  &  $k_i$

Attention mechanism allows us to aggregate data from many (key, val) pairs

ATL mech provides a differentiable means of control by which a NN can select elements from a set

to construct an assoc weighted sum over representations

## 11.2 Attention pooling by similar

Example applies Att in a reg & classif setting via kernel density estimation

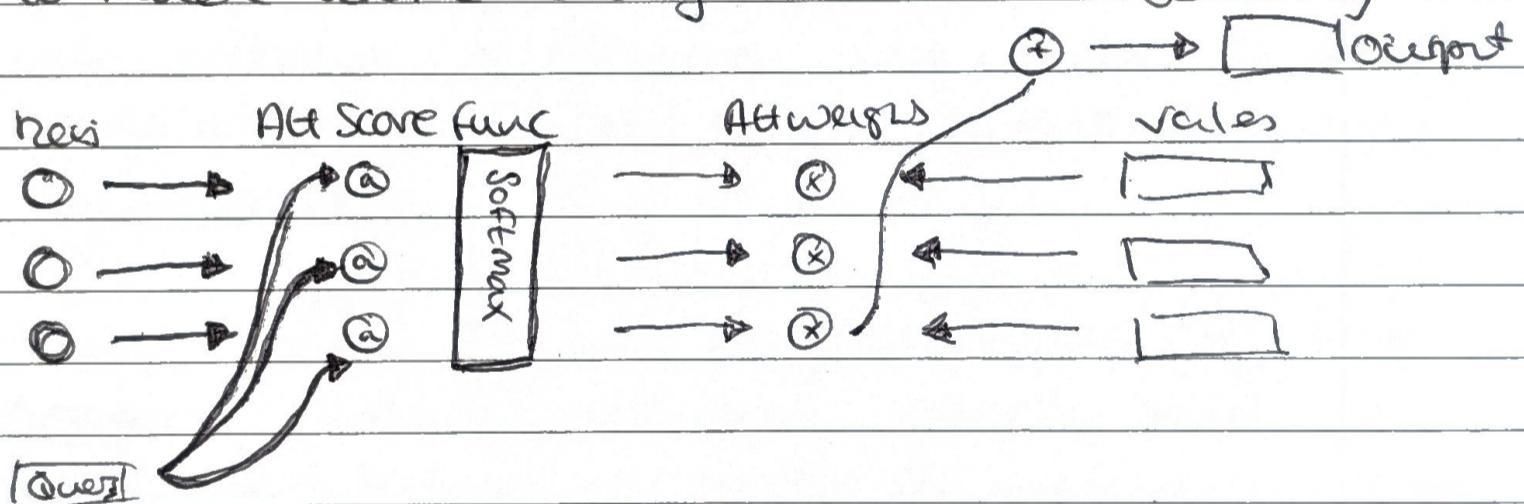
Skipped

## 11.3 Attention Score Functions

(a)

Distance functions are more expensive than dot prods

so much work has gone into Att Scoring



## Dot Product Attention

Skipped

## 11.4 Bahdanau Attention mechanism

In the previous chapter we looked @ encoder-decoder arch for seq2seq based on RNNs (Sutskever, 2014)

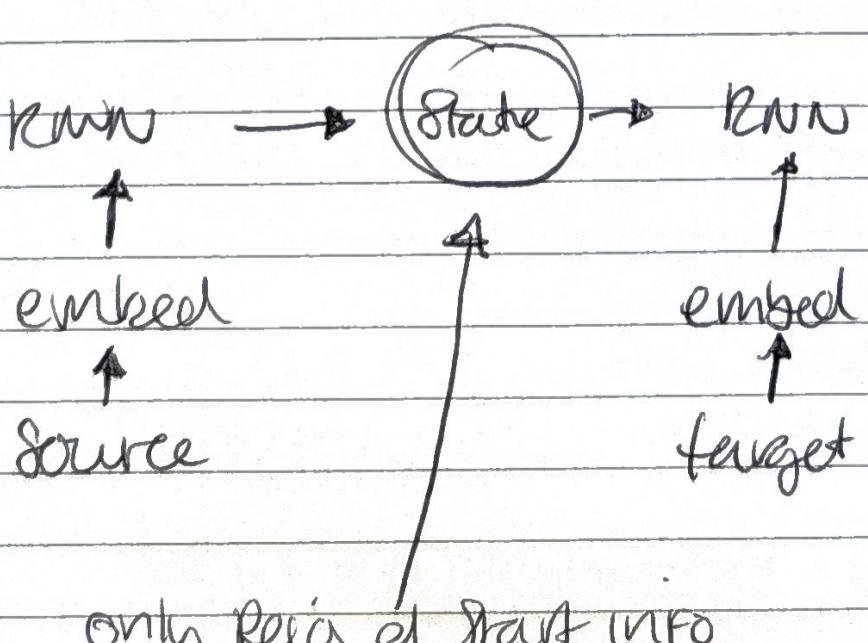
↳ encoder transforms variable length seq into fixed shape context var

↳ decoder generates the output seq token-by-token based on the generated tokens & context variable

### ~~Contextual~~

Conventionally, in RNNs all relevant info about a source sequence is ~~trans~~ translated into a fixed dimension vector by encode

this fixed state is picked up by the ~~encoder~~ Decoder as the single and exclusive source of info for generating the sequence



only relies on start info  
in encoder-decoder arch

the encode-Decode is reasonable for short seq  
but infeasible for long

Soon run out of "space" in the internal representation to store all important info

Consequently Decoder will fail to translate long & complex seqs

~~Bad Bahdanau et al (2014)~~

↳ when pred, if not all input tokens are relevant

the model attends to only part of the input sequence

that are Deemed relevant for prediction

used to update current state

before next generated token

key point =

$c$  = context vector state each

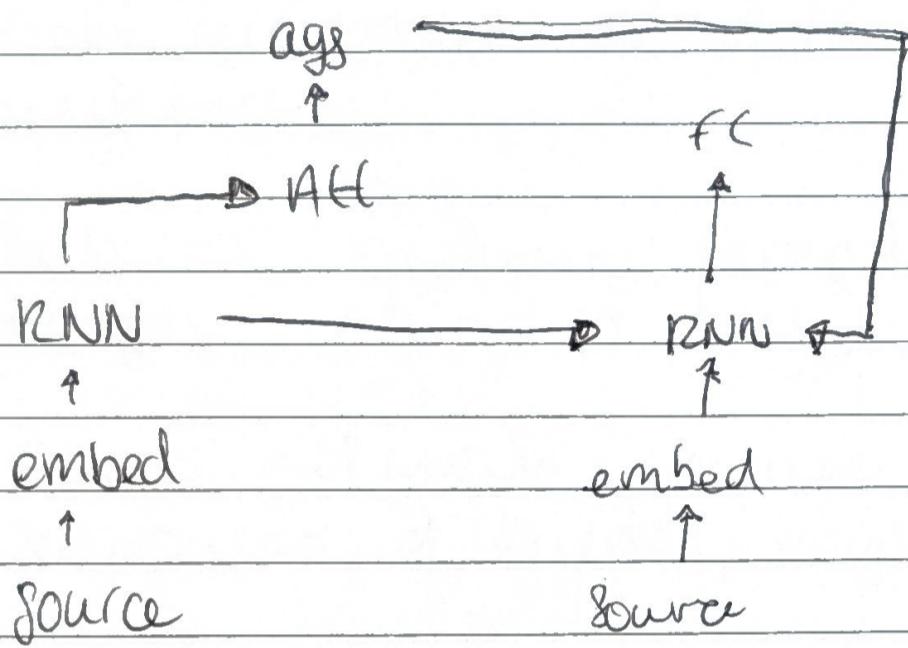
Dynamically update this @ $t$  token gen

↳ using (1) original text  $h_t$

(2) Decoder generated  $s_{t-1}$   
yields  $C$  of which is what Decoder uses

Content variable is the output of the attention pooling  $C_t$

$$C_t = \sum a(s_{t-1}, h_t) h_t$$



this model was later updated to include ATT at the generate step too

## Summary

when pred a token

↳ if not all input tokens are relevant

↳ RNN ↗ encoder-Decoder w/ ATT Mech

Selectively aggregates big parts of the input Seq

Achieved by treating the "state" as additive attention pooling

Decoder hidden state at previous time step  
= Query

encoder hidden state at all time steps  
as both keys & values

## 11.5 Multi-head Attention

given same set of queries, keys & values

we may want model to combine knowledge from diff behaviours of the same attention mechanism

such as capturing dependancies of various ranges (short vs long)

allow att mech to use diff representation subspaces of queries, keys & values

pool  
Instead of performing a single attention

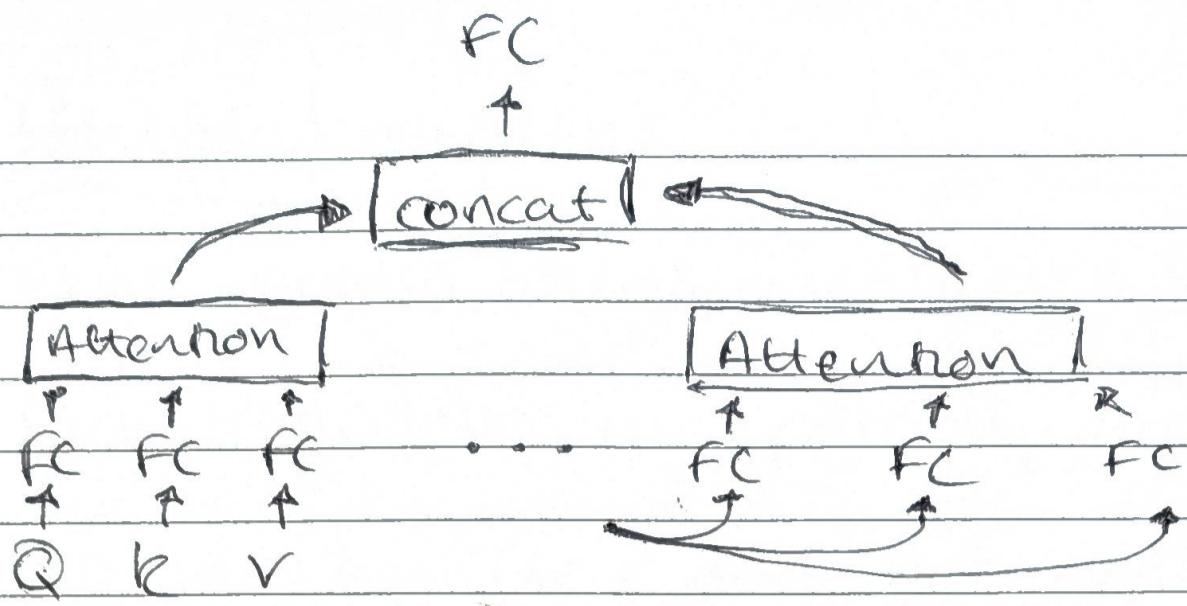
queries, keys and values can be transf w/ n indep learned linear projections

these are then fed to the attention pooling in parallel

the att pooling outputs are concatenated & transformed

↳ w/ another linear proj to form output

Multi Attention pool, where each Att pool output = A head



### What is Multi-Head Attention

#### 11.6 Self-Att & Positional encoding

Imagine imagine feeding sequence of tokens into an Att mech such that @ every step each token has its own query, keys & value

when comparing the value of a tokens representation @ the next layer

the token can attend (via its query) to any other tokens

using the full set of Query-key comp score

for each token, we can compute a rep by building the appr weighted sum over all the tokens.

this is diff to prev lect  
before Decoder attend to encoder step  
here each token attend to every other token

## Positional encoding

RUNS process tokens one-by-one

SELF-ATTENTION uses parallel computation

SELF-ATT BY ITSELF DOES NOT PRESERVE ORDER OF THE SEQUENCE

WHAT TO DO IF WE NEED ORDER?

REPRESENT THIS TO THE MODEL AS ADDITIONAL INPUT TO EACH TOKEN

THESE ARE CALLED POSITIONAL ENCODINGS

THESE CAN EITHER BE LEARNED OR FIXED

POS ENCODINGS CAN BE BASED ON SINE / COSINE FUNCTIONS FOR EXAMPLE

## 11.7 transformer Architecture

earlier Self-Att models rely on RNNs for input representations

the transformer model is solely based on attention mechanisms w/out any CNN or RNN layer

Originally proposed for Seq-to-Seq on text data

Later applied to wide areas (lang, vision, speech etc)

transformer like has an encoder-decoder structure

but there is no NN

just modules based on Self-Att & position encodings

Skipped as rest has no connection to NNs

## 11.8 Vision transformer

Skipped

## 11.9 Large Scale Pretrain w/ transformers

All examples so far has includes Coding  
models from scratch to perform specific task

each Model becomes a specific expert  
that is sensitive to even slight shift  
in the data sets

for better generalized models Pretrain  
on large data has become increasingly  
common