

Summary:

single self attn:

$$Z = \text{softmax}(QK^T) \cdot V$$

$$Q = XW^Q$$

$$K = XW^K$$

multi-head attn: $(z_1, \dots, z_h) W^O = Z$

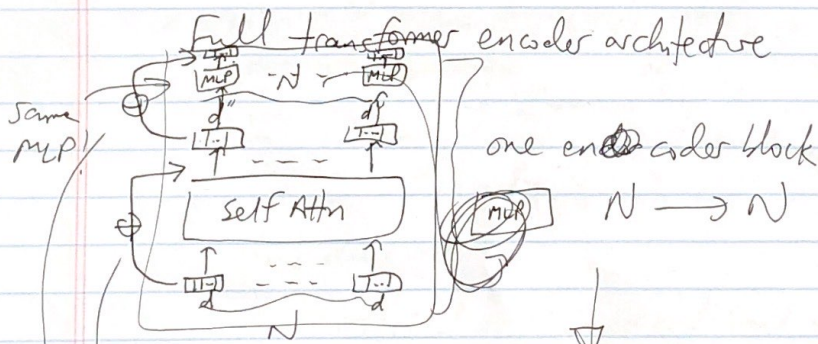
$$V = XW^V$$

- can check: perm equivariant! interchanging order of rows of X will result in same output w/ output rows interchanged.

- "MLP w/ input-dependent weights"

$$Z \sim (XW^Q W^K^T X) \cdot XW^V$$

this is like ~~an~~ X-dep weight



Can stack!

(original "Attn is all you need" had to encoder blocks)

+ residual/skip connections!

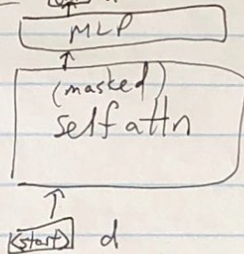
The encoder produces a transformed embedding that takes into acct entire input w/ all its ^{inter}relations

Transformer Decoder

To generate ~~text~~, need the decoder architecture

(original transformer paper had encoder-decoder pair w/ add'l encoder att'n layer...)

Same structure as encoder but w/ "future masking"



output is vector z w/ same dim as original embedding

interpreted as "score" for N_{vocab} words

Multiply into vocabulary/dictionary embedding matrix (e.g. for English)

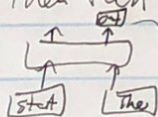
$$z \cdot D \leftarrow d \left\{ \begin{pmatrix} | & | & \dots & | \end{pmatrix} \right.$$

N_{vocab} dim'd vector

↳ interpret as score for each word in vocab
choose word w/ highest score, or top-k etc.

↓
e.g. "The"

Then feed back into decoder



e.g. "The" → "Thing" | use ~~prev~~ embedding of last word
(which knows about entire preceding sequence) to predict next word

encoders useful for learning embeddings

↓
classification / categorization
sentiment analysis

⋮
BERT example

"Bidirectional Encoder Representations from Transformers"
Google (2018)

QK^T

- ~~Future~~ Future masky: don't want attention scores to depend on future words

$$\text{softmax}(\text{mask: } \begin{matrix} \langle \text{start} \rangle \\ I \\ an \\ for \end{matrix} \begin{matrix} \langle \text{start} \rangle & I & an & for \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{matrix} \cdot QK^T)$$

- This is an example of "autoregressive model"

$$\text{Learning } P(x_1, \dots, x_n) = P(x_1) P(x_2 | x_1) P(x_3 | x_1, x_2) \dots P(x_n | x_1, \dots, x_{n-1})$$

- GPT is an example of this ~~encoder~~ decoder only model
Generative Pre-trained Transformer OpenAI (2018)

GPT \rightarrow BERT \rightarrow GPT2 \rightarrow ...

now GPT-family (decoder only) are state of the art!

- Can also give "prompt" \rightarrow just an initial sequence instead of $\langle \text{start} \rangle$

~~Can~~ prompt called diff language \rightarrow translation!

- Both BERT & GPT use ^{self-supervised} pre-training to learn attn model

"masked language modeling"

"next word prediction"

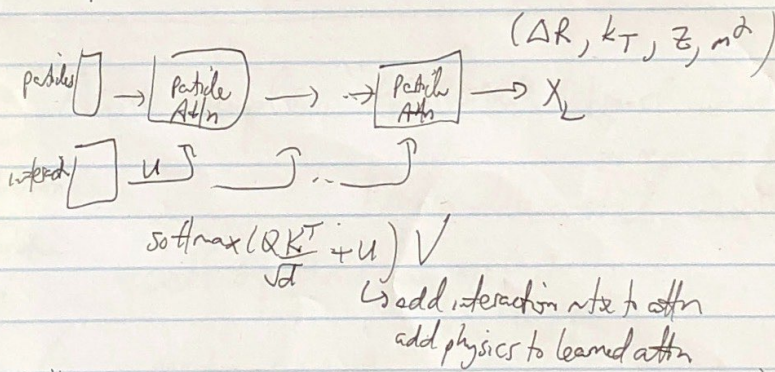
using huge amount of data & grant model
fine tuning on "downstream" ~~(task)~~ "Foundation model" "backbone"
tasks (eg translation, classification, ...)

Examples of applications of transformers

Qu, Li, Qian 2020.03772 "Particle Transformer for Jet Tagging"¹

- Introduced permutation equiv. wch. for jet tagging based on transformers
- Also introduced new dataset "JetClass" for training, it
 (previous datasets ~1M) 100M jets \approx 10 types \times 10M each
 - 4 rec
 - particle ID (chg hadron, neutral hadron, e, μ , γ)
 ↳ not truth, recs!
 - displacement

• particle features & pairwise "interaction" features N²N²



• "class attn" (don't completely understand this)

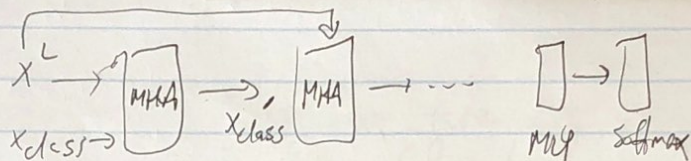
extract classifier output - standard trick in vision transformer³

randomly initialized token

accumulates ~~the~~ info from event from attention to other tokens

$$\begin{cases} Q = W_Q x_{\text{class}} \\ K = W_K [x_{\text{class}}, x_L] \\ V = W_V [x_{\text{class}}, x_L] \end{cases}$$

Class attn (from vision transformer lit)



starts off random but accumulates info

- Part achieved SOTA results on all jet tagging tasks!

example of
"foundation
model"

- Also interesting: pretrain + fine train \gg train from scratch on smaller dataset
- Also transformer benefited more from pretrain than GNN!

- Astro example: "ASTROMER transf. based embeddy for repr of light curves"
Domoso-Oliva et al 2205.01677

cf TimeModAttn
example from
Astro presentation
- pretrain
on sim

- self supervised, like NLP etc while Part
- positional encoding (not perm equiv)

→ masked light curve modeling
encodes-decoder

- Data R-band light curves of MACHO survey (Galactic Bridge & LMC) for pretrain - 1.5M light curves
- 20k labeled variable stars for f.t. & eval
Also OGLE-II 360k labeled; ATLAS 420k

pretraining + finetuning
 Also outperforms classifiers trained on ^{just} ~~large~~ datasets!

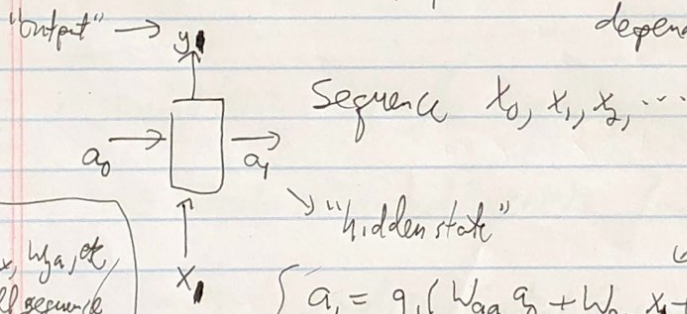
Although transformers have "taken" over, prior approaches to sequence modeling still useful to have in toolbox (as in ASTROMER example) — RNN, LSTM, GRU

Ref: stanford CS-231n "cheatsheet" for RNNs

Brief overview of RNN, LSTM, GRU architectures

— very sequential, not at all perm equiv.

— not feed forward per se — or like FF w/ # layers depends on length of sequence



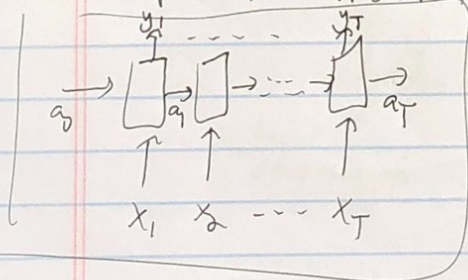
$W_{aa}, W_{ax}, W_{ya}, etc.$
 same for all sequence
 "weight sharing"

$$\begin{cases} a_1 = g_1(W_{aa} a_0 + W_{ax} x_0 + b_a) \\ y_0 = g_2(W_{ya} a_1 + b_y) \end{cases}$$

depends on prev hidden state + current input

depends on current hidden state

repeat for each element of sequence



produces transformed sequence + hidden state that encodes full sequence

• Can also use as $1 \rightarrow \text{many}$

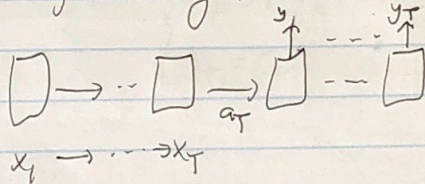
$x_1 \rightarrow y_1 \rightarrow y_2 \rightarrow \dots$ (music, text gen)

• $\text{many} \rightarrow 1$: just keep y_T ~~seq~~ \rightarrow classification

• $\text{many} \rightarrow \text{many}$: name entity recog (classify each token according to type)
esp named entities (person, corp, ...)

seq

• $\text{many} \rightarrow \text{many}$: translation



\rightarrow also: stores to train since sequential

short term memory problem

• vanilla RNNs struggle w/ vanishing/exp gradients - long sequences information is lost

\rightarrow introduce "gated" RNNs (analog of residual/skip connectors!)

Two popular examples: Gated Rec. Unit (GRU)

Long Short-Term Memory (LSTM)

GRU \subset LSTM (special case, LSTM more general)

hidden + cell state

Local \leftarrow
short term info

\rightarrow global info
long term