

# Transformers & Attention Mechanism

- Used in SOTA NLP like ChatGPT
  - Used in pretty much everything now
  - Permutation-equivariant architecture!
  - actually a type of GNN
  - like an MLP w/ input-dependent weights  
$$X \rightarrow W(X) \cdot X$$
- much more powerful in principle!

{ - Key idea: (self) attention mechanism

- Key idea: learned embeddings

↳ learn relations btw elements of sequence  
relations expressed by dot product  
in embedding space

## NLP example:

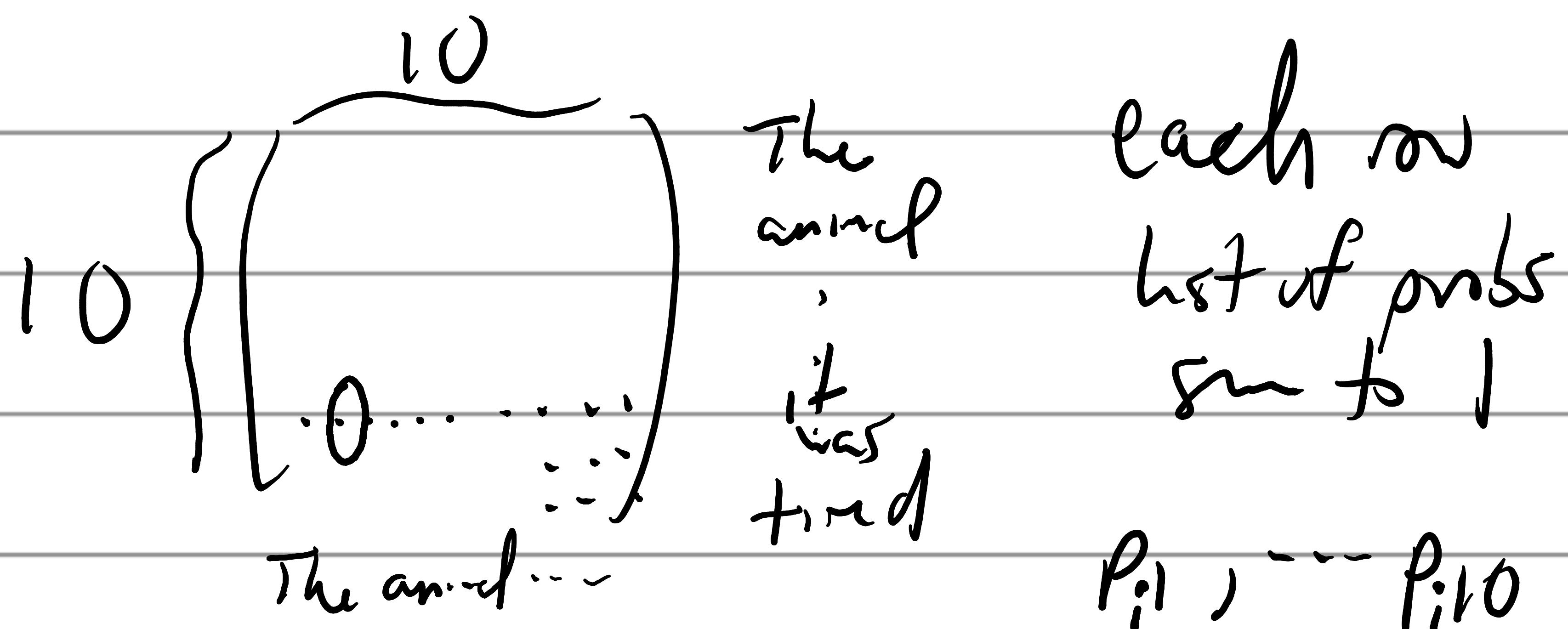
"The animal didn't cross the street because it was tired."

- What does "it" refer to?

Clearly "animal" not "street"!

- earlier NLP models had trouble w/ this  
(trouble being long-term relations)

- Attention matrix



each row  
list of probs  
sum to 1

$$P_{i1}, \dots, P_{i10}$$

$$\sum P_i = 1$$

→ Maybe set  $P_{72}$  and  $P_{77}$   
large, others small

# The Attention Mechanism (Bahdanau 2014)

Want to map input

sequence to embedding space

$$\mathbb{R}^d \rightarrow \mathbb{R}^{d'}$$

so that relations are captured by dot product

$$y_i \cdot z_j \quad y, z \in \mathbb{R}^{d'}$$

want probabilities, use softmax

$$p_{ij} = a(y_i \cdot z_j) \quad \text{"attention matrix"}$$

if  $y, z$  come from same input  $\rightarrow$  "self-attention"

) Simplest map: linear

$$\left\{ \begin{array}{l} X \rightarrow XW^Q \leftarrow d \times d' \text{ matrix} \\ \quad = Q \text{ "query"} \end{array} \right.$$

$$\left\{ \begin{array}{l} X \rightarrow XW^K = K \text{ "key"} \end{array} \right.$$

We can apply attention matrix  
to input sequence to produce a new sequence

$$\left( \sum_{i=1}^N p_i; \vec{x}_i \right)$$

$$\vdots$$
  
$$\left( \sum_{i=1}^N p_i; \vec{x}_i \right)$$

If we were in embedding space

can learn another map

$$x \rightarrow xW^V = V \text{ "values"}$$

Transformer architecture

Sequence  $\rightarrow$  sequence mapping

Recap:

Given sequence of vectors  $\in \mathbb{R}^d$

$$X = \begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vdots \\ \vec{x}_N \end{pmatrix} \in \mathbb{R}^{N \times d}$$

Learn 3 maps to embedding space  $\mathbb{R}^{d'}$

Query  $Q = X.W^Q$

Key  $K = X.W^K$

Value  $V = X.W^V$

$d \times d'$  matrices  
Learned weights

self attn matrix:

$$P = \alpha(Q.K^T)$$

$N \times N$  matrix

Transformer output:

$$Y = P.V \in \mathbb{R}^{N \times d'}$$

\* Note: permutation equivariant!

$$X = \begin{pmatrix} \vec{x}_1 \\ \vdots \\ \vec{x}_n \end{pmatrix} \in \mathbb{R}^{d} \rightarrow Y = \begin{pmatrix} \vec{y}_1 \\ \vdots \\ \vec{y}_n \end{pmatrix} \in \mathbb{R}^{d'}$$

interchange  $\vec{x}_i \leftrightarrow \vec{x}_j \longleftrightarrow$  interchange  $\vec{y}_i \leftrightarrow \vec{y}_j$

can check ✓ indices

$$Y_{id} = \sum_a (X_{ia} W_{a\beta}^Q (W^K)^T_{\beta b} (X^T)_{bj}) X_{jc} W_{cd}^V$$

perm acts here    sum perm inv't

In fact can argue for this architecture  
as simplest non-linear permutation-equiv.  
seg → seg map

## "Multi-headed Attn"

in example "it" also relates to "find"  
how to capture that?

could overload style after matrix

but more powerful to learn another one!

$$\cdot W^{Q,K,V}_{r=1, \dots N_h}$$

analogous to multiple filters  
in CNN!

• each gives  $N \times d'$  embedded sequence

$$z_r$$

• Concatenate  $(z_1, \dots, z_{N_h})$   $\underbrace{N \times (d' \cdot N_h)}$

• feed it all together w/ final matrix

$$W^{Q,K,V} \quad (d' \cdot N_h) \times d''$$

$$z_{\text{out}} = z^T W^{Q,K,V}$$

## "Positional Encodings"

- in NLP & other sequences, position does  
conta info!
- add back in position info w/ positional encodings

$$\vec{x}_i \rightarrow \vec{x}_i + \vec{\text{pos}}(i) \quad \text{common example}$$

$\sin(\omega_a i)$

$\cos(\omega_a i)$

- transformer still perm equiv  
but not wrt original input sequence

## "Tokenization"

In NLP need to map words to vectors

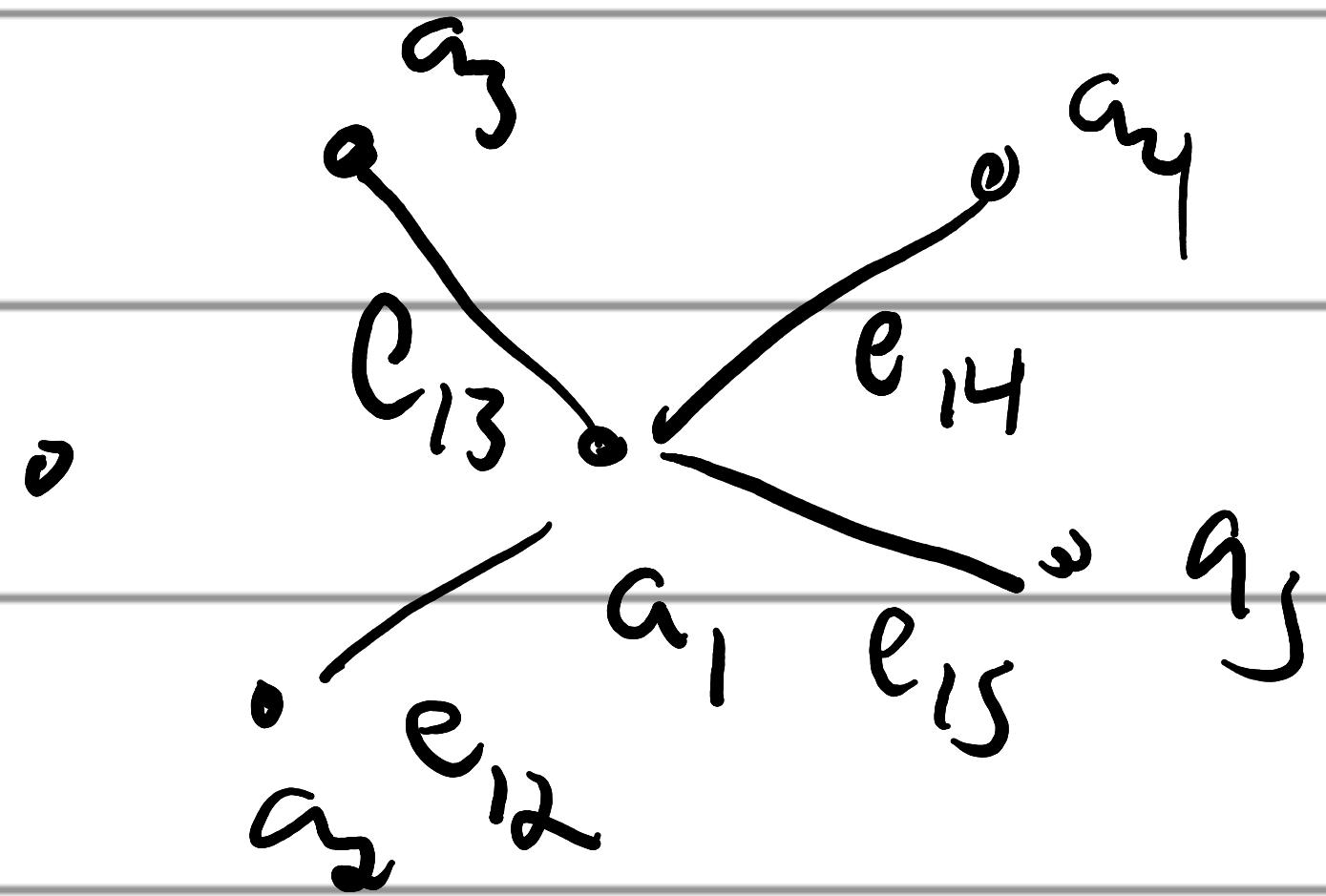
"tokens"

many tokenization algorithms

offering ML - pretraining task

# Transformers as GNNs

Recall GNN



MLPs

$$a_i' = \varphi^{\text{node}}(a_i, \rho(e'_i; j_i; e_{ij}))$$

$$e_{ij}' = \varphi^{\text{edge}}(a_i, a_j, e_{ij})$$

Transformer

$$x_i' = \cancel{\text{Attn}} \ p_{ij} x_j \leftarrow \begin{array}{l} \text{this is like } e_{ij} \\ \text{this is } a_j' \\ \hookrightarrow \text{attn}(x_i, x_j) \end{array}$$

~ Fully connected graph

nodes are  $x_i$   
edges are  $\text{Attn}(x_i, x_j)$

} multiplying  
not usually  
considered

So Transformers still go beyond standard GNNs

Summary:

single self attn:

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

$$Q = XW^Q$$

$$K = XW^K$$

$$V = XW^V$$

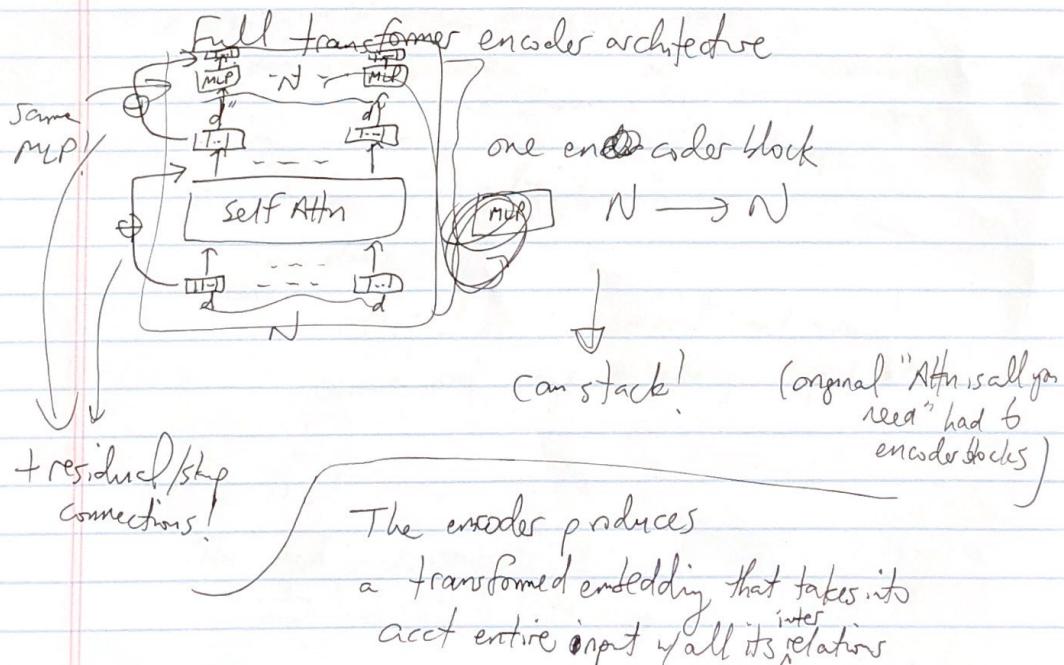
$$\text{multi-head attn: } (z_1, z_h) W^O = z$$

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

- "MLP w/ input-dependent weights"

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

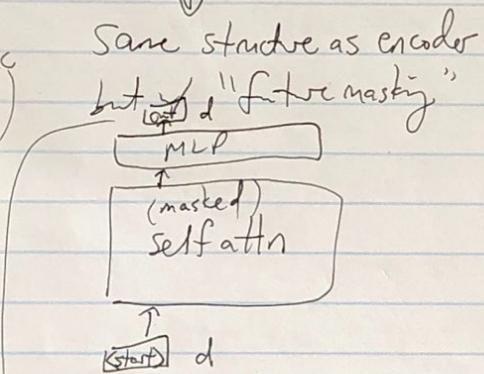
this is like ~~one~~  $X$ -dep weight



## Transformer Decoder

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

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



→ output is vector, <sup>z</sup> w/ same dim as original embedding.

$$z \cdot D \leftarrow_d \{ (| | | \dots |) \}$$

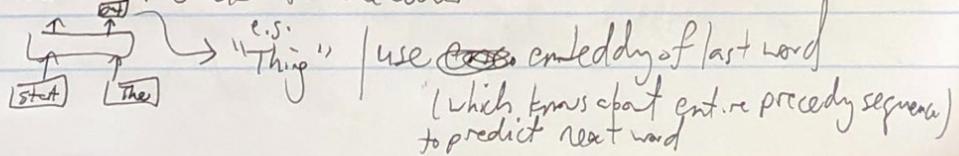
Nocab don't vector

Multiplify into vocabulary/dictionary  
Nocab embedding matrix  
(~50k for English)

↳ 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



encoders useful  
for learning embeddings

classification / categorization  
sentiment analysis

BERT example

"Bidirectional Encoder Representations from Transformers"  
Google (2018)

$QKT$

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

$$\text{softmax}(\text{mask}: \begin{matrix} (\text{start}) & I & \text{am} & \text{fan} \\ I & \vdots & \ddots & \ddots \\ \text{am} & \vdots & \ddots & \ddots \\ \text{fan} & \vdots & \vdots & \vdots \end{matrix}, QKT)$$

- This is an example of "autoregressive model"

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 ~~one~~ decoder-only model

Generative Pre-trained Transformer OpenAI (2018)

GPT → BERT → GPT2 → ...

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

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

~~the~~ prompt can be diff language → translation!

Both BERT & GPT use <sup>self-supervised</sup> pre-training to learn attn model

"masked language  
modeling"

"next word prediction"

using huge amount of data & great 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"

- Introduced permutation equiv. arch. for jet tagging based on transformers
- Also introduced new dataset "JetClass" for training it

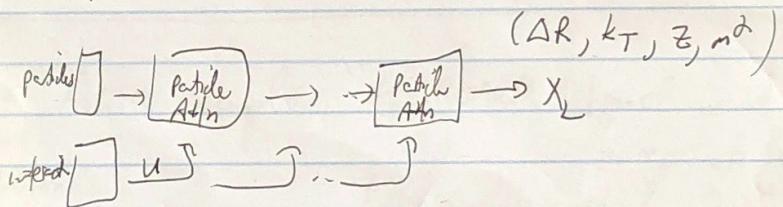
(previous datasets ~1M) 10M jets  $\Rightarrow$  10 types  $\times$  10M each

- 4 vec

$\xrightarrow{\text{new to JetClass}}$  - particle ID (chg hadron, neutral hadron, e,  $\mu$ ,  $\tau$ )  
↳ not truth, reco!

- displacement

- particle features & pairwise "interaction" features NNMT



$$\text{softmax}\left(\frac{QK^T}{\sqrt{d}} + U\right) V$$

↳ add interaction info to attn  
add physics to learned attn

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

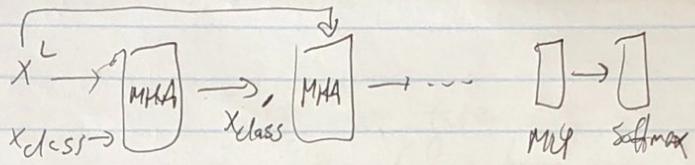
extract class feature output - standard trick in vision transformer<sup>3</sup>

randomly initialized token

accumulates ~~attn~~ 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 <sup>• Also interesting: pretrain + fine tune >> + train from scratch</sup>  
 "foundation model" <sup>on smaller dataset</sup>
- Also transformer benefited more from pretrain than GNN!

- Astro example: "ASTROMER transf. based catalog

~~cf TimeMHA~~  
 example from  
 Astro presentation  
 - pre-trained  
 on sim

for repr of light curves  
 Donoso - Ohwa et al 2205.01677

- self supervised, like NLP etc  
<sup>while Part</sup>
- positional encoding (not per query)

→ masked light curve modeling  
 encoder-decoder

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