

Lecture 8(8); Self Attention, Transformer NN

- demo for HW6 pb 3,4
- PROJECTS due Sunday 8/10
- HW6 due Fri 8/15 (incl. demo)
- last class Wed 8/13 (word-to-rec) NLP

Transformer

Movie concept \Rightarrow transformer NN

- NO CLASS FRI 8/15
use time for PROJ, HW6
- TRACE finals
- grades due Mon 8/18



Self-attention

- each token/step t attend each step (prev) s

- transf seq (encoder) \rightarrow new seq (enc)
NOT enc-dec mech like before

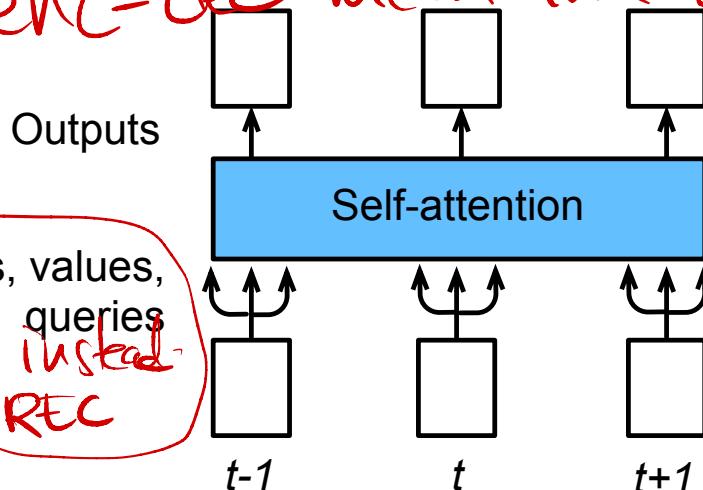
- To generate n outputs with n inputs, we can copy each input into a key, a value and a query

- No sequential \Rightarrow view SEQ information is preserved

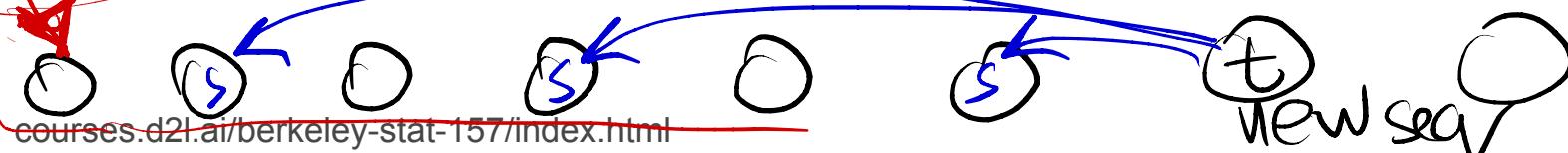
- Run in parallel \Rightarrow AI CAP RACE

Long Sequence

ATTN

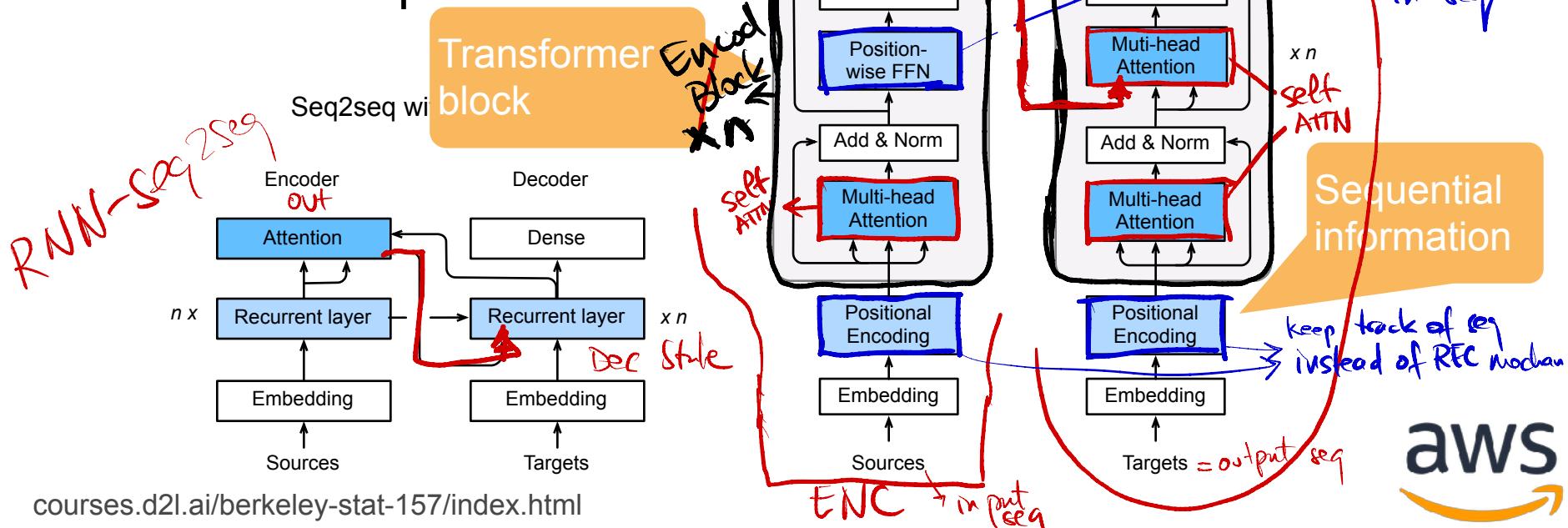


ENCODER has self attn



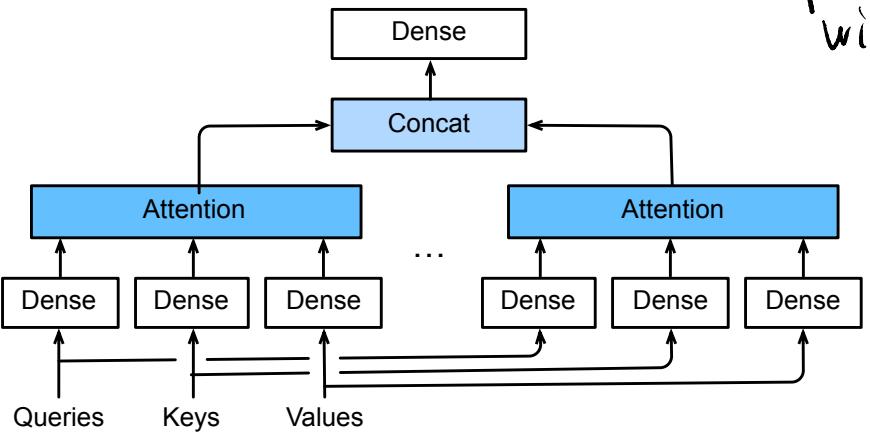
Transformer Architecture

- It's an **encoder-decoder** arch
- Differ to seq2seq with attention in 3 places



Multi-head Attention

Add & Norm → layer normalization (after comp block)
 - normalizes output at flat layer
 - prevent problems with computation activations
 with gradients

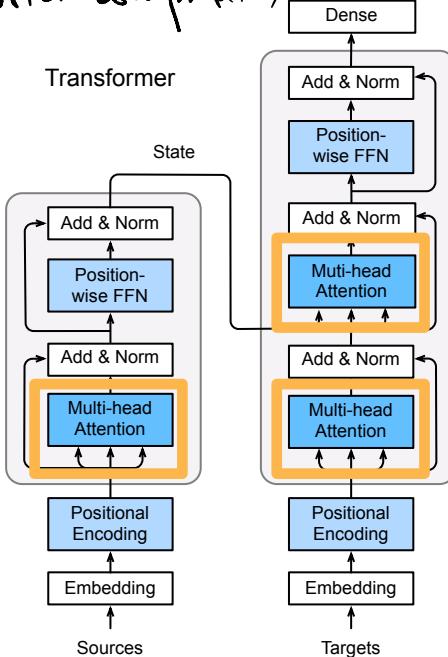


$$\mathbf{W}_q^{(i)} \in \mathbb{R}^{p_q \times d_q}, \mathbf{W}_k^{(i)} \in \mathbb{R}^{p_k \times d_k}, \text{ and } \mathbf{W}_v^{(i)} \in \mathbb{R}^{p_v \times d_v}$$

$$\mathbf{o}^{(i)} = \text{attention}(\mathbf{W}_q^{(i)} \mathbf{q}, \mathbf{W}_k^{(i)} \mathbf{k}, \mathbf{W}_v^{(i)} \mathbf{v})$$

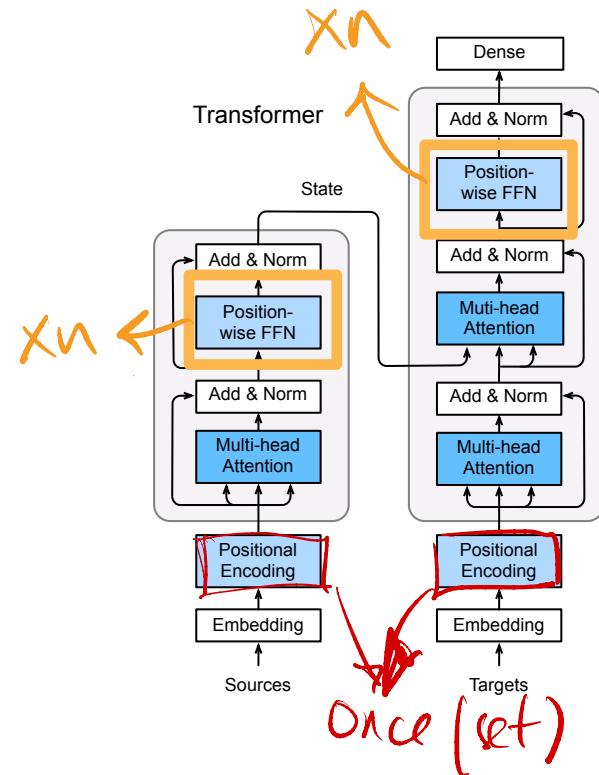
for $i = 1, \dots, h$

$$\mathbf{o} = \mathbf{W}_o \begin{bmatrix} \mathbf{o}^{(1)} \\ \vdots \\ \mathbf{o}^{(h)} \end{bmatrix}$$



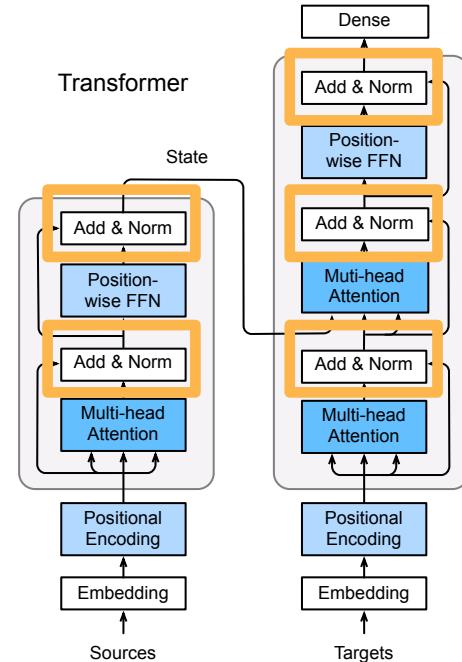
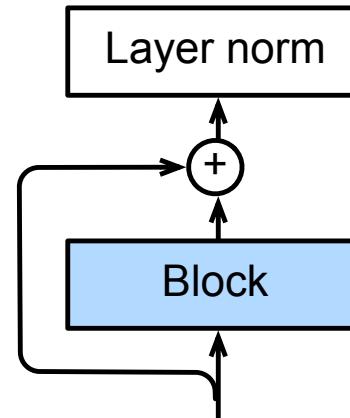
Position-wise Feed-Forward Networks

- Reshape input (batch, seq len, fea size) into (batch * seq len, fea size)
- Apply a two layer MLP
- Reshape back into 3-D
- Equals to apply two (1,1) conv layers



Add and Norm

- norm per layer*
- Layer norm is similar to batch norm → norm per batch
 - But the mean and variances are calculated along the last dimension
 - `X.mean(axis=-1)` instead of the first batch dimension in batch
norm `X.mean(axis=0)`



Positional Encoding

- Assume embedding output $X \in \mathbb{R}^{l \times d}$ with shape (seq len, embed dim)
- Create $P \in \mathbb{R}^{l \times d}$ with

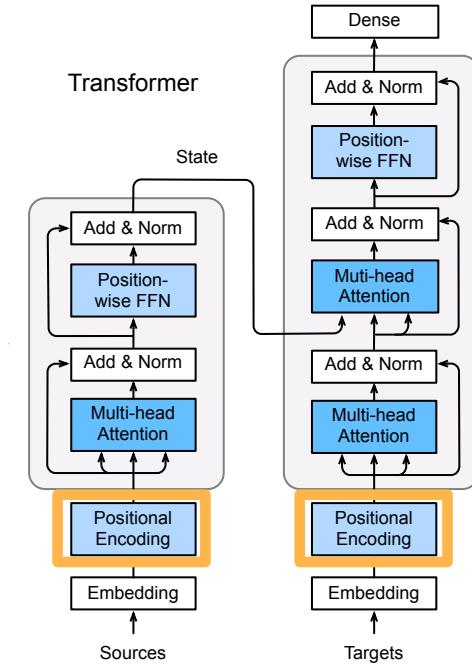
global
pos

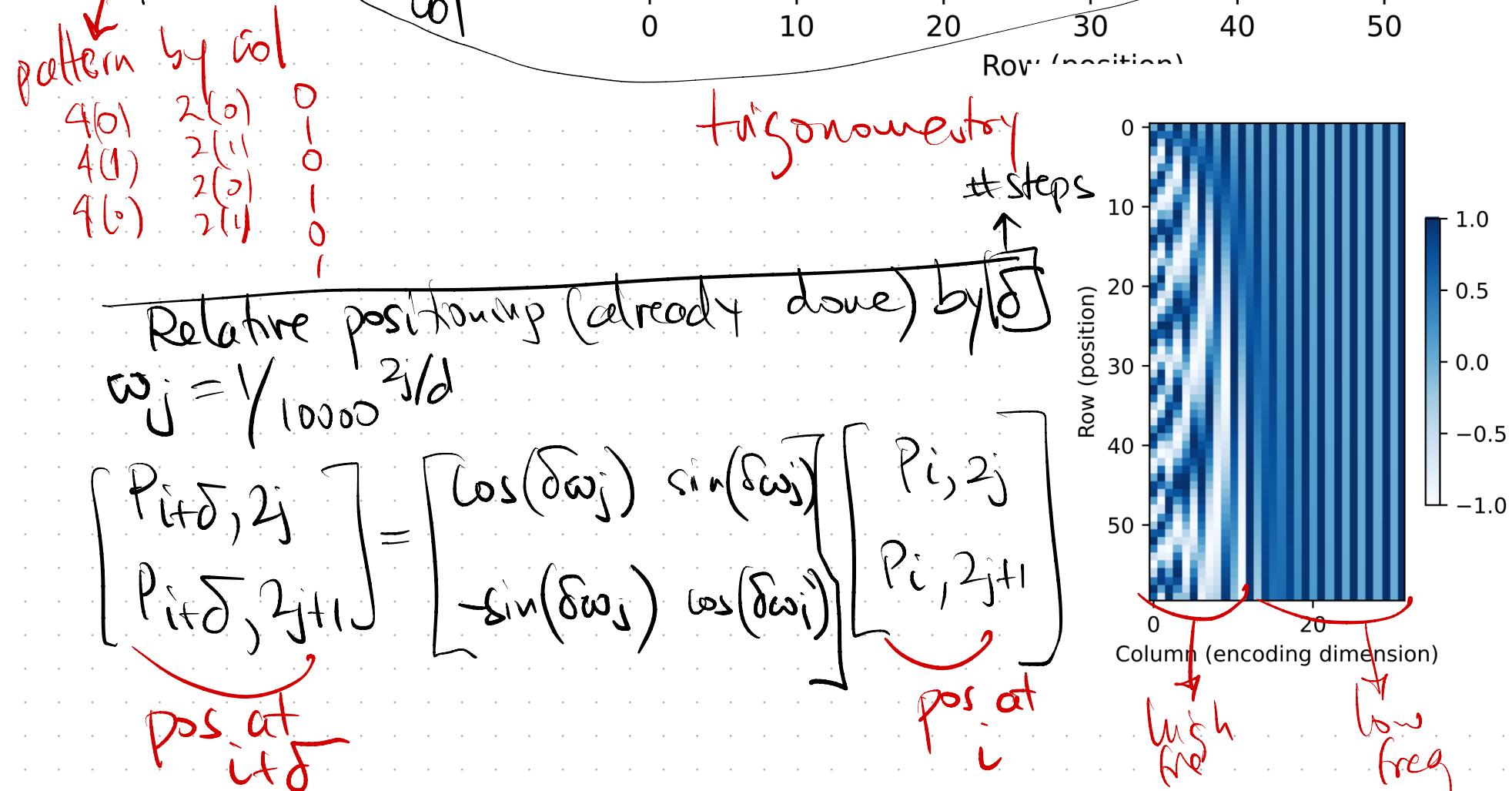
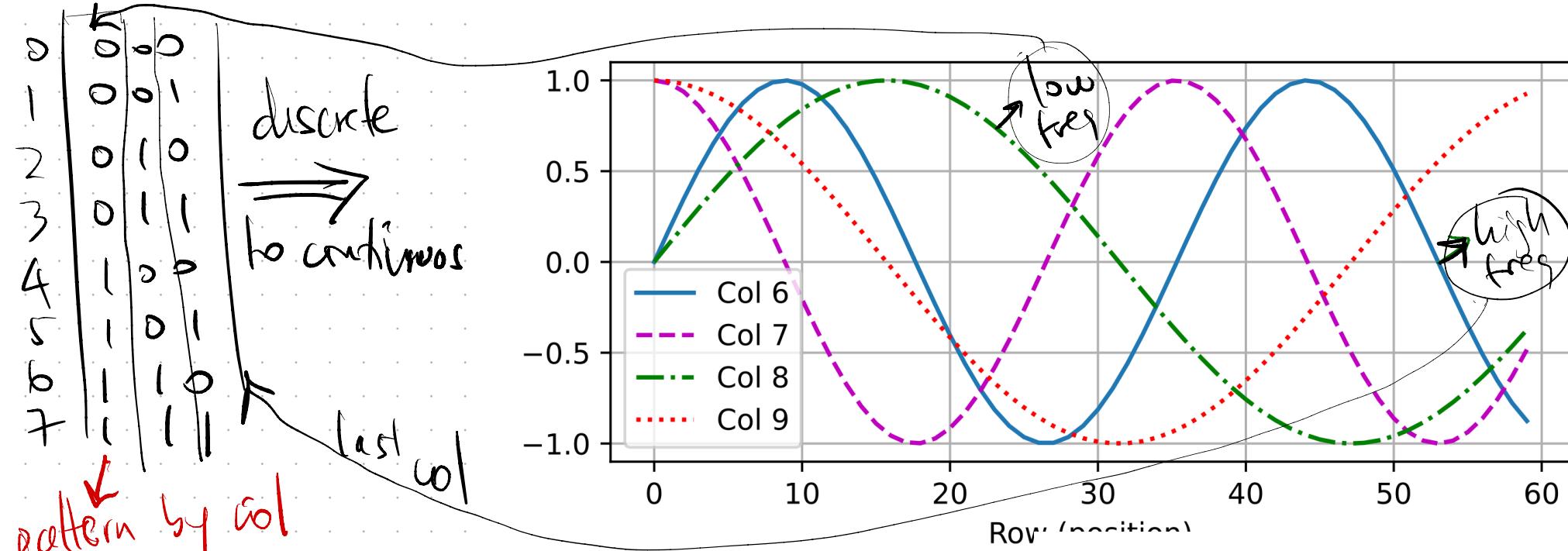
$$\begin{cases} P_{i,2j} = \sin(i/10000^{2j/d}) \\ P_{i,2j+1} = \cos(i/10000^{2jd}) \end{cases}$$

trigonometry
- periodicity
- separability
- analogy with
binary bits

- Output $X + P$

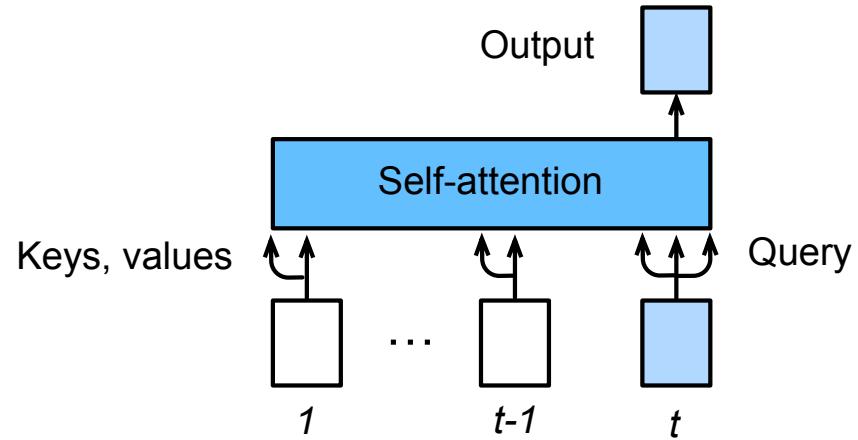
$X + P$
data seq positions





Predicting

- Predict at time t :
 - Inputs of previous times as keys and values
 - Input at time t as query, as well as key and value, to predict output



Code...

BERT



Transfer Learning in NLP

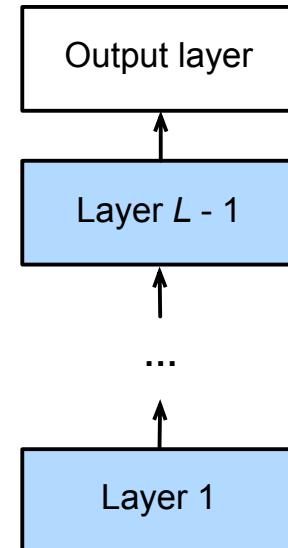
- Use pre-trained models to extract word/sentence features for the new task
 - E.g. word2vec or language model
- Often don't update the pre-trained models
- Need to construct a new model to capture the information needed for the new task
 - Word2vec ignores sequential information, language model only looks in a single direction

Motivation of BERT

- A fine-tuning based approach for NLP
- The pre-trained model capture sufficient data information
- Only need to add a simple output layer for a new task

NLP

Positive



CV



Classifier

Feature
extractor

I love this movie

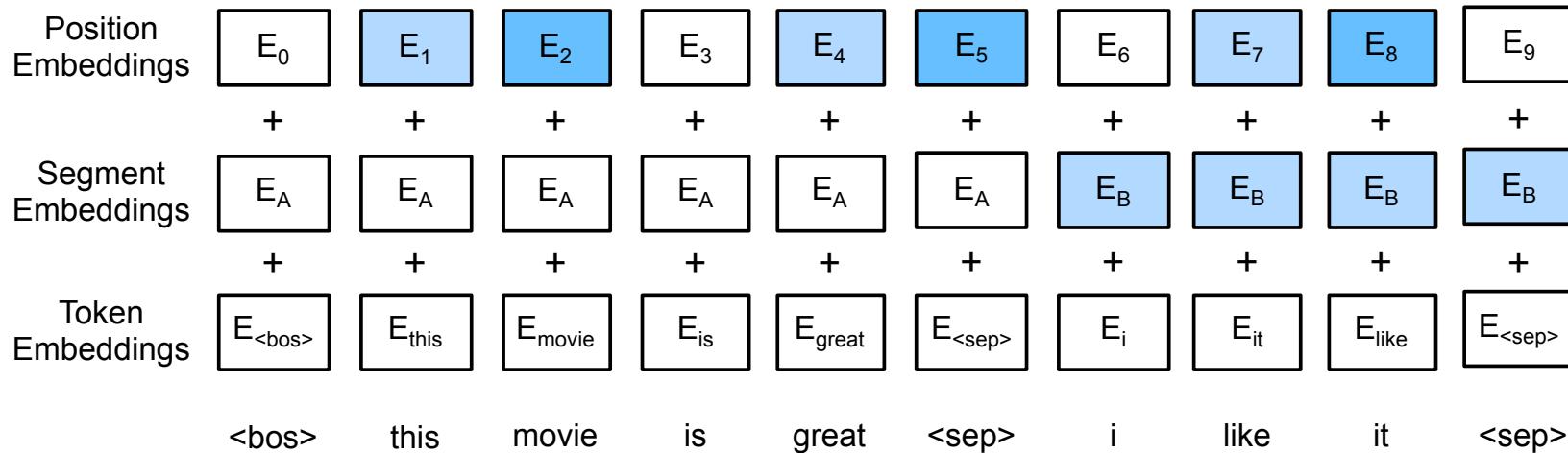


BERT Architecture

- A (big) Transformer encoder (without the decoder)
- Two variants:
 - Base: #blocks = 12, hidden size = 768, #heads = 12, #parameters = 110M
 - Large: #blocks = 24, hidden size = 1024, #heads = 16, #parameters = 340M
- Train on large-scale corpus (books and wikipedia) with > 3B words

Modification of inputs

- Each example is a pair of sentences
- Add an additional segment embedding



Pre-training Task 1: Masked Language Model

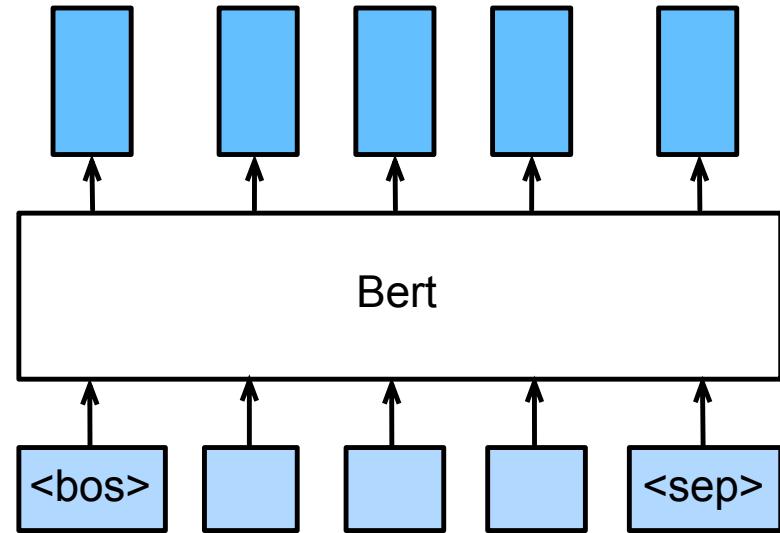
- Randomly mask (e.g. 15%) tokens in each sentence, predict these masked tokens
 - Transformer is bidirectional, which breaks the unidirectional limit of standard LM
- No mask token (<mask>) in fine tuning tasks
 - 80% of the time, replace selected tokens with <mask>
 - 10% of the time, replace with randomly picked tokens
 - 10% of the time, keep the original tokens

Pre-training Task 2: Next Sentence Prediction

- 50% of time, choose a sequential sentence pair
 - <bos> this movie is great <sep> i like it <sep>
- 50% of time, choose a random sentence pair
 - <bos> this movie is great <sep> hello world <sep>
- Feed the Transformer output of <bos> into a dense layer to predict if it is a sequential pair

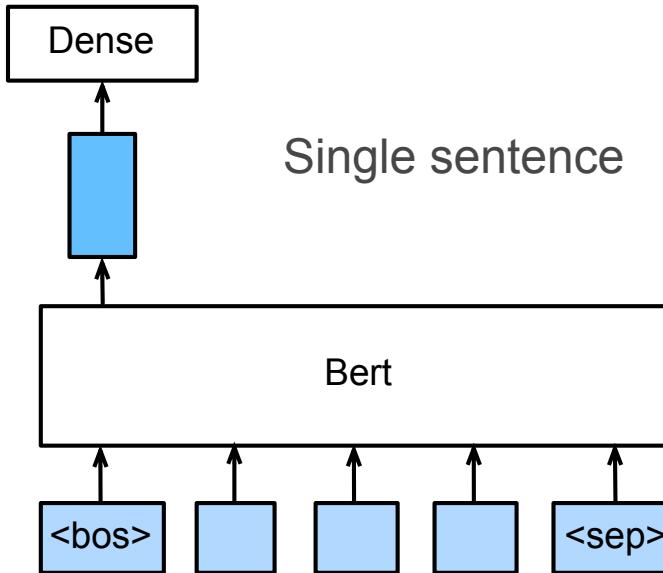
Bert for Fine Tuning

- Bert returns a feature vector for each token that captures the context information
- Different fine-tuning tasks use a different set of vectors

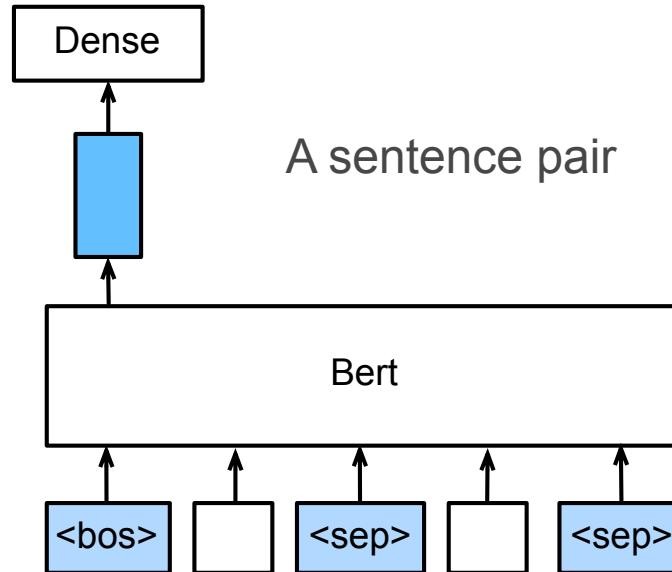


Sentences Classification

- Feed the <bos> token vector into a dense output layer



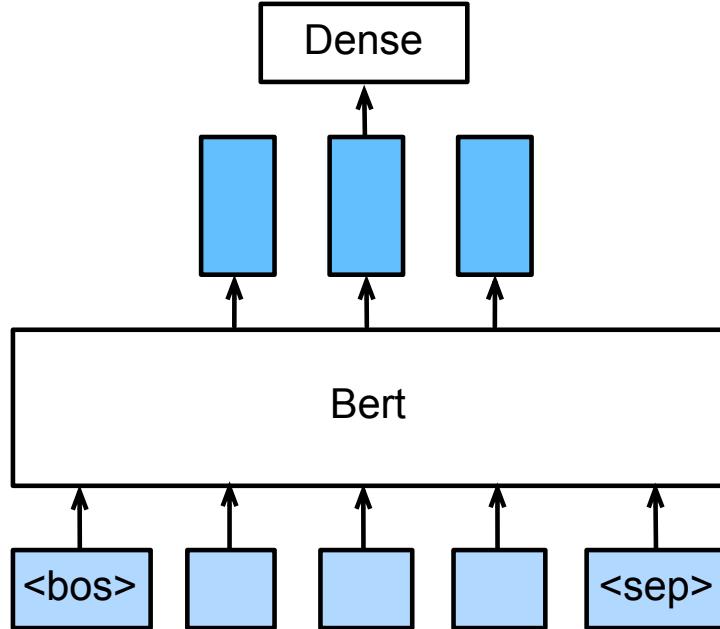
Single sentence



A sentence pair

Named Entity Recognition

- Identify if a token is a named entity such as person, org, and locations...
- Feed each non-special token vector into a dense output layer

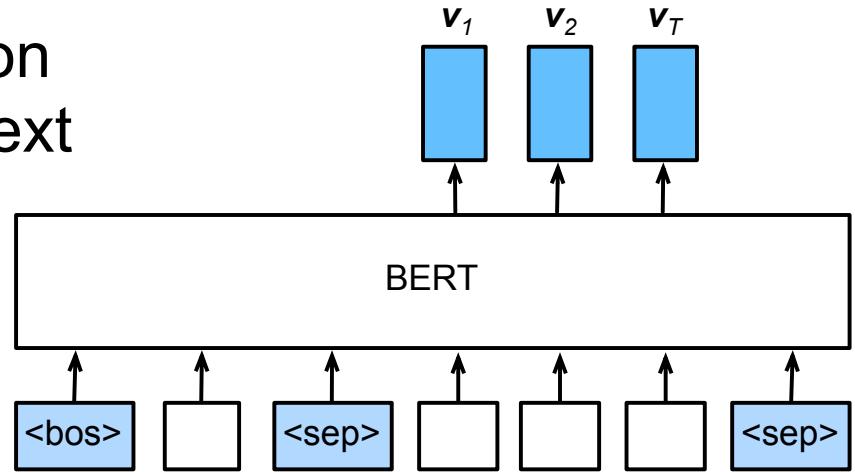


Question Answering

- Given a question and a description text, find the answer, which is a text segment in the description
- Given p_i the i -th token in the desperation, learn s so that

$$p_1, \dots, p_T = \text{softmax}(\langle s, v_1 \rangle, \dots, \langle s, v_T \rangle)$$

p_i is the probability i -th token is the segment start. Same for the end



Check GluonNLP for code...