

Introduction to Deep Learning

22. Encoder-Decoder, Seq2seq

STAT 157, Spring 2019, UC Berkeley

Alex Smola and Mu Li

courses.d2l.ai/berkeley-stat-157

Lecture 8/6 : Attention

Next - Multi Attention

- Transformers Architect
- use Attention instead of Convolutions, Rec

Attention

"Attention is All you need"
(2017)

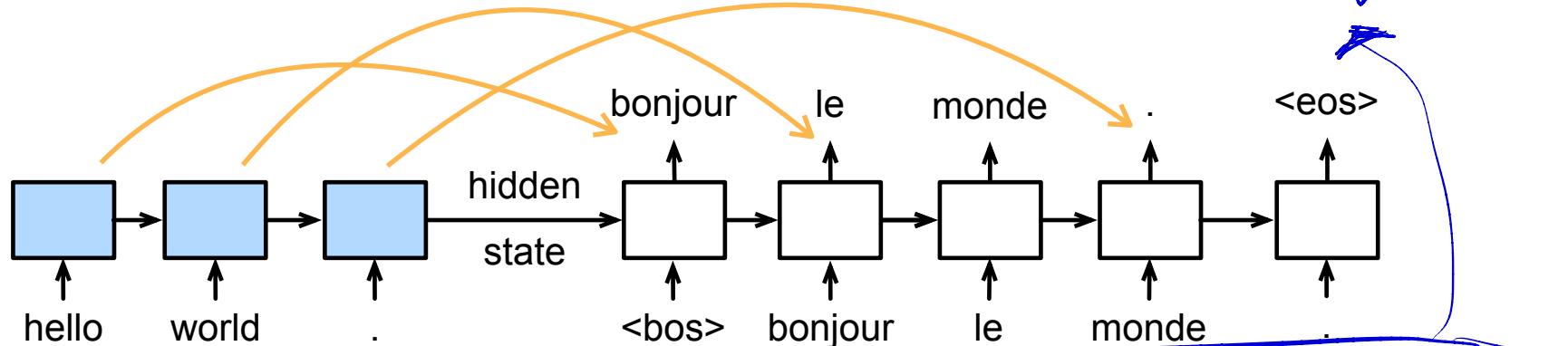


Motivation

Attention = a new layer in NN

- Each generated token might be related to different source tokens

from RNN parts
K? q? v?



- `Attention(q, k, v, alpha, h...)`
- math formulas
- code (d2L)

• how these play a role
inside RNN (NN for seqence)
(2017)

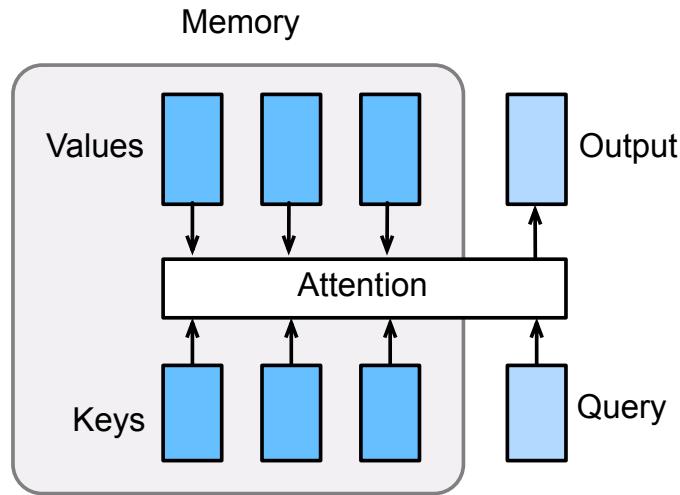
Attention Layer : q =queries ; k =keys ; v =values $\alpha = \text{sim}(k, q)$
dim ; H =hidden ; $h = \# \text{attention heads}$ "kernel"
(multi head attention).

- An attention layer explicitly select related information

- An attention layer explicitly select related information
 - Its memory consists of key-value pairs
 - The output will be close to the values whose keys are similar to the query

current step, query to the values whose keys are similar to the query

Attention (q_V) $D^{(k, v)}$) =
 $= \sum_{i \in D} \alpha(q_V, k_i) v_i$ "Aggregate"
prev steps



weighted-avg (v_i)
with weights $\alpha(q_j k_i)$

Flashback: $\alpha(u, v) = \text{similarity kernel}(x_u, x_v)$

Weighted avg by similarity: HWS

$$\text{pred}(z, Y) = \frac{\sum_{i=1}^N \alpha(z, x_i) \cdot Y_i}{\sum \alpha(z, x_i)} = \frac{\text{weighted avg}(Y_i)}{\text{sum } \alpha(z, x_i)}$$

by weights

$\alpha(z, x_i) = \text{high} \Rightarrow x_i \text{ is important for } z$

$\alpha(z, x_j) = \text{low} \Rightarrow x_j \text{ not important for } z$

$\text{pred}(z)$ should pay **ATTENTION** to x_i , but
not much to x_j

RNN Attention properties

- $\sum \alpha(q_i, k_i) = 1$ (distribution) norm $\alpha(q_i, k_i) = \frac{\alpha(q_i, k_i)}{\sum_j \alpha(q_i, k_j)}$

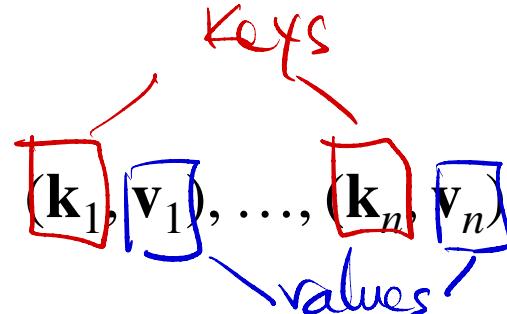
- $\alpha(q_i, k_i) \geq 0 \Rightarrow$ use softmax instead of norm

$$\underbrace{\alpha(q_i, k_i)}_{\text{output}} = \frac{\exp(\text{linear score}(q_i, k_i))}{\sum_j \exp(\text{linear-score}(q_i, k_j))}$$

Attention Layer

queries

- Assume query $\mathbf{q} \in \mathbb{R}^{d_q}$ and the memory with $\mathbf{k}_i \in \mathbb{R}^{d_k}$ $\mathbf{v}_i \in \mathbb{R}^{d_v}$



- Compute n scores a_1, \dots, a_n with $a_i = \alpha(\mathbf{q}, \mathbf{k}_i)$
- Use softmax to obtain the attention weights

$$b_1, \dots, b_n = \text{softmax}(a_1, \dots, a_n)$$

- The output is a weighted sum of values

$$\mathbf{o} = \sum_{i=1}^n b_i \mathbf{v}_i$$

Vary α to get different attention layers

Dot Product Attention

linear similarity = "linear kernel"
score (q, k) = $\langle q \cdot k \rangle$

- Assume the query has the same length as values $q, k_i \in \mathbb{R}^d$

- Vectorized version

- m queries $Q \in \mathbb{R}^{m \times d}$ and n keys $K \in \mathbb{R}^{n \times d}$

$$\alpha(q, k) = \langle q, k \rangle / \sqrt{d}$$

SCALE DOT PRODUCT

scale - by length

comp.
irrelevant

(normalization, scaling)

$$\alpha(Q, K) = QK^T / \sqrt{d}$$

matters

$$\text{Euclidean? } \alpha(q, k) = -\frac{1}{2} \|q - k\|^2 = q^T k - \frac{1}{2} \|k\|^2 - \frac{1}{2} \|q\|^2$$

$$\text{Softmax(Euclid)} = \exp\left(-\frac{1}{2} \|q - k\|^2\right)$$

aws

- Attention "Pooling" as similarity computation

$$\left\{ \begin{array}{l} \alpha(q_i, k) = \frac{\langle q_i, k \rangle}{\sqrt{d}} \text{ linear} \\ \end{array} \right.$$

important for parallel computation

$$-\alpha(q_i, k) = \exp\left(-\frac{1}{2}\|q_i - k\|^2\right) \text{ Gaussian}$$

$$-\alpha(q_i, k) = \exp\left(-\frac{1}{2} \frac{\|q_i - k\|^2}{2\sigma^2}\right) \text{ Gaussian scaled for width adapted}$$

$$-\alpha(q_i, k) = \begin{cases} 1 & \text{if } \|q_i - k\| \leq \text{threshold} \\ 0 & \text{if not} \end{cases} \text{ "Boxcar"}$$

$$-\alpha(q_i, k) = \max\{0, 1 - \|q_i - k\|\}$$

Attr score = $\sum_i [V_i] \cdot \frac{\alpha(q_i, k)}{\text{normalized}}$

values

- Masked Softmax : sequences of different lengths.

$$\begin{matrix} q^T \cdot k \\ \downarrow \sqrt{d} \\ \text{query key} \end{matrix} \rightarrow q^T \boxed{M} \cdot k$$

masked matrix (all zero dimensions)

$d_{q, k}$

Multilayer Perception Attention

- Learnable parameters $\mathbf{W}_k \in \mathbb{R}^{h \times d_k}$, $\mathbf{W}_q \in \mathbb{R}^{h \times d_q}$, and $\mathbf{v} \in \mathbb{R}^h$

$$\alpha(\mathbf{k}, \mathbf{q}) = \mathbf{v}^T \tanh(\mathbf{W}_k \mathbf{k} + \mathbf{W}_q \mathbf{q})$$

Annotations: \mathbf{v}^T is labeled "values", $\mathbf{W}_k \mathbf{k}$ is labeled "keys", and $\mathbf{W}_q \mathbf{q}$ is labeled "queries".

- Equals to concatenate the key and query, and then feed into a single hidden-layer perception with hidden size h and output size 1

Bahdanau : important

RNN $\xrightarrow{\text{step}}$ Transformer
enc-dec

Bahdanau

$s = \text{index}_{\text{prev}}$
hidden rec step

$$t = \sum_{s=1}^S d_{ts} \cdot h_s$$

$$d_{ts} = \text{sm}(t, k_s)$$
$$k_s = v_s = h_s$$

= key $k_s = \text{val}$ $v_s = h_s = \text{Encoder hidden state}$

= how useful
is $\text{enc}(x_s)$ at step t

Decoder.

query
 $q_t = \text{step } t \text{ hidden state}$
of decoder

$$\text{score}(s_t, h_s) = V_t \cdot \tanh(W_s \cdot h_s + W_t s_t)$$

s token in
input sequence

dec
hidden
state
at
step t
aws

Similarity

$$\alpha(t, s) = \text{softmax}(\text{score}(s_t, h_s))$$

At decoder step/state

from past sequence (x)

t: What encoder steps h_s are most relevant?

Seq2seq with Attention

\Rightarrow give those steps high weights
 $\alpha(t, s)$

Multihed Attention

- Multiple Attn model = sorting at same time
 $H_1 H_2 H_3 \dots H_h$ $h = \# \text{ heads (Attn size)}$

- Computation simplified, parallel \Rightarrow larger scale with GPU
- q, k, v = linear function of input \Rightarrow much faster

- replace REC, CONV layers with particular attn

- no REC \Rightarrow diff sequence mechanism?
step 1, 2, ...
in sequence
- more general / can be repurposed easily
not common in RNN

$$\text{Score} = \frac{Q^T}{\sqrt{d}} \cdot V \quad ; \quad \alpha = \text{softmax(score)}$$

dot prod easy to compute. but possible

$Q = X \cdot W_q$ queries $K = X \cdot W_k$ keys $V = X \cdot W_v$ values

prev enc(l) of x to work with q

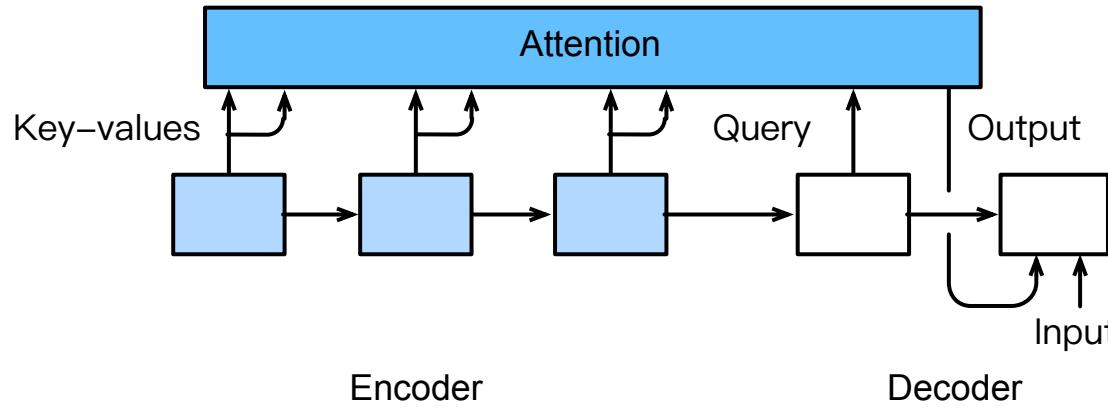
Current decoder demand

learn W_q, W_k, W_v, W_o output

Multihead (Q, K, V) = Concat (H_1, H_2, \dots, H_h) $\cdot W_o$

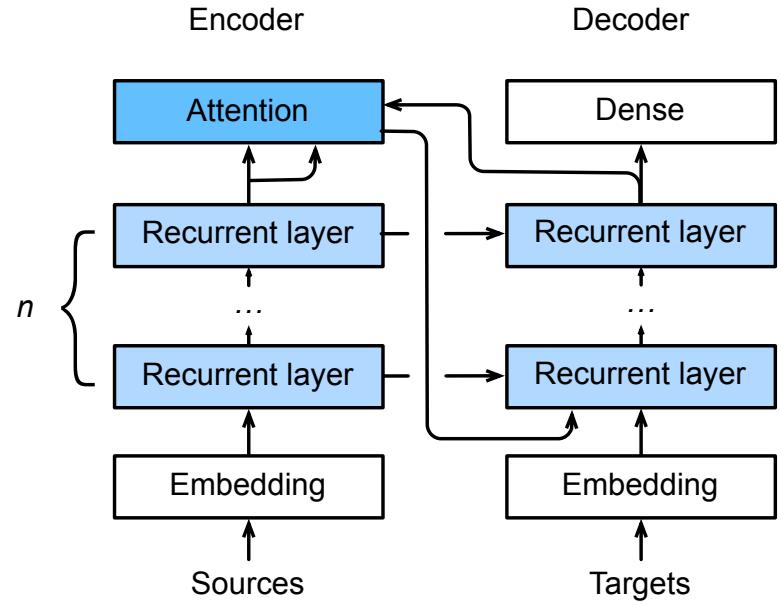
Model Architecture

- Add an additional attention layer to use encoder's outputs as memory
- The attention output is used as the decoder's input



Encoder/Decoder Details

- The output of the last recurrent layer in the encoder is used
- The attention output is then concatenated with the embedding output to feed into the first recurrent layer in the decoder



Next five: transformers

due
Fr 8/11 ← HW6 demo part 2.

Code...

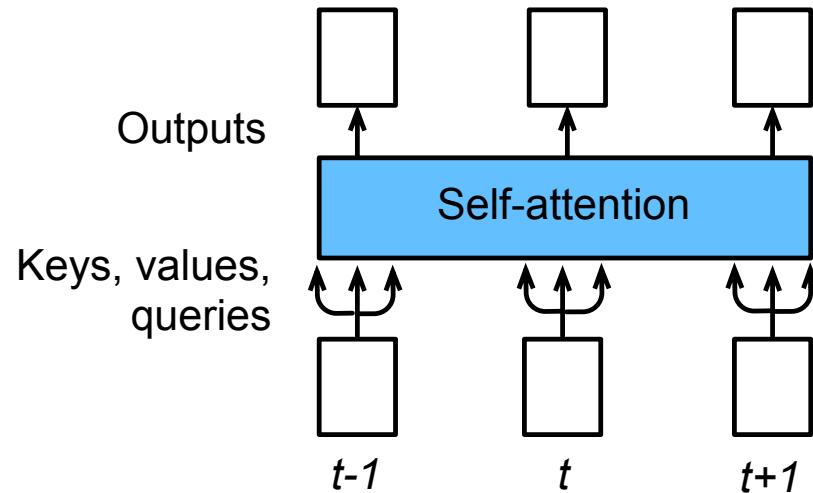
projects due 8/10, keep
updating chat with TA

Transformer



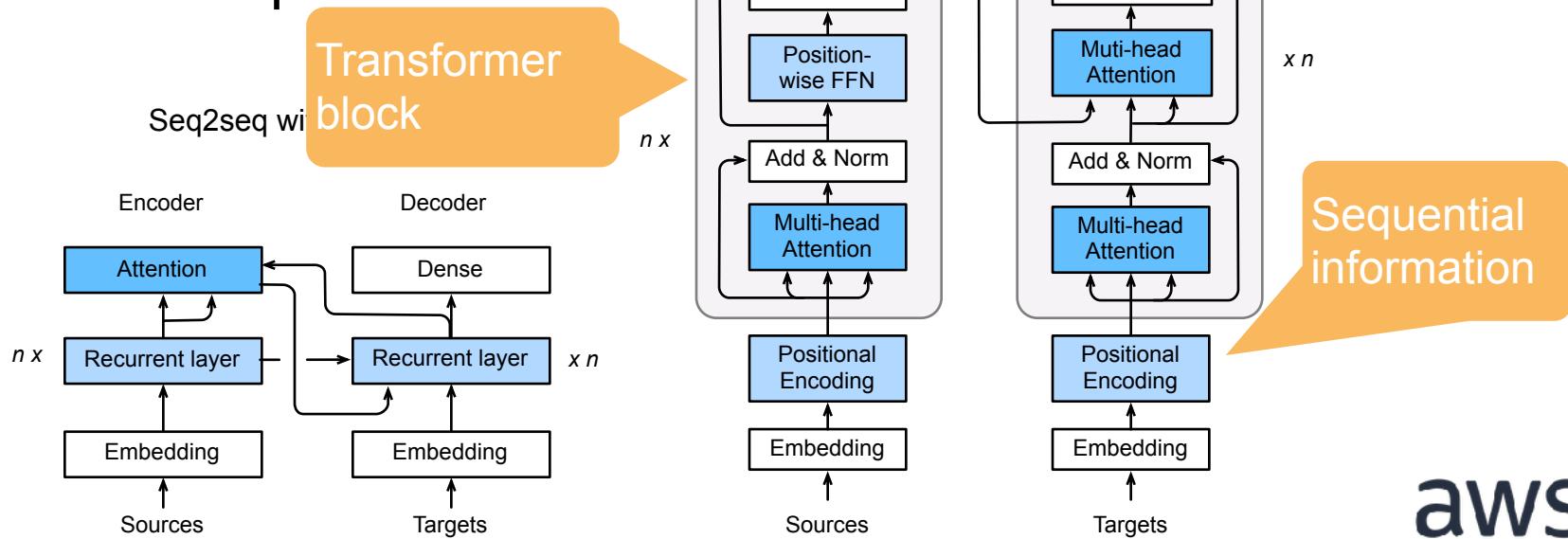
Self-attention

- To generate n outputs with n inputs, we can copy each input into a key, a value and a query
- No sequential information is preserved
- Run in parallel

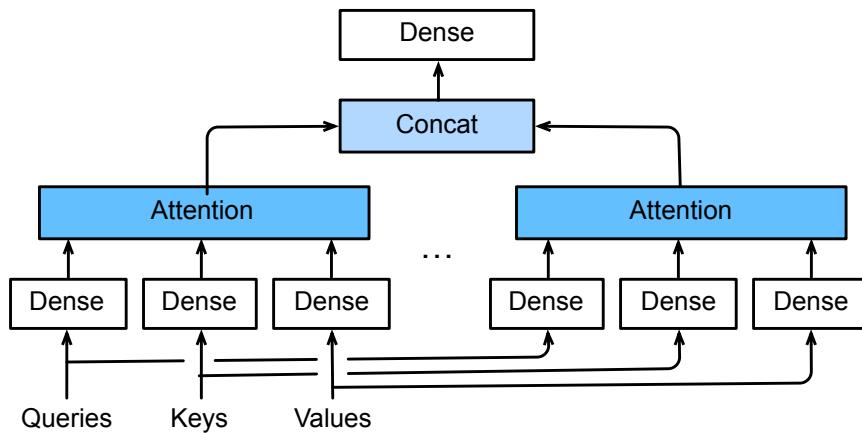


Transformer Architecture

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



Multi-head Attention

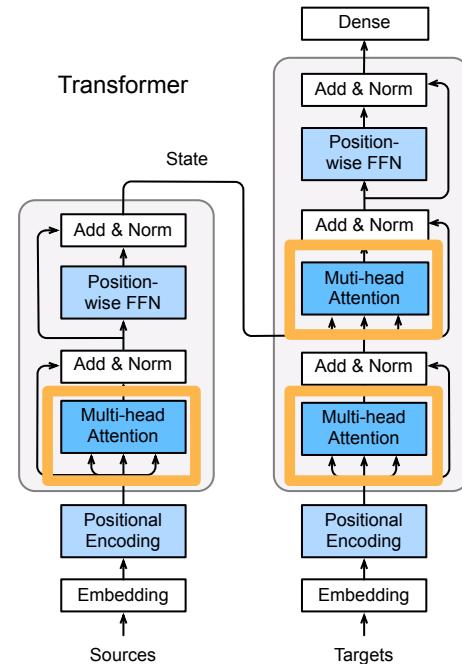


$$\mathbf{W}_q^{(i)} \in \mathbb{R}^{p_q \times d_q}, \mathbf{W}_k^{(i)} \in \mathbb{R}^{p_k \times d_k}, \text{ and } \mathbf{W}_v^{(i)} \in \mathbb{R}^{p_v \times d_v}$$

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

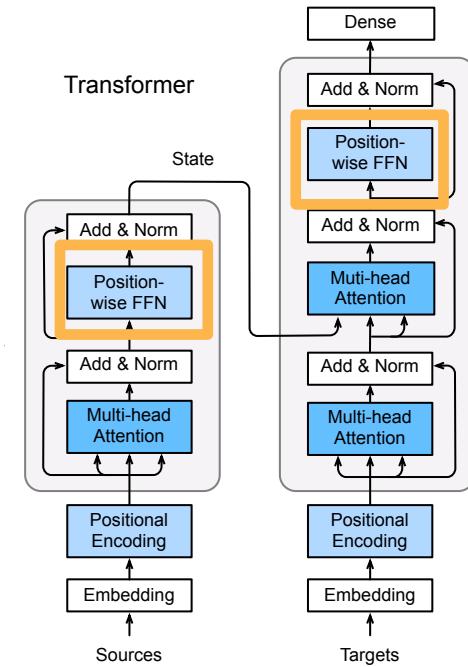
for $i = 1, \dots, h$

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



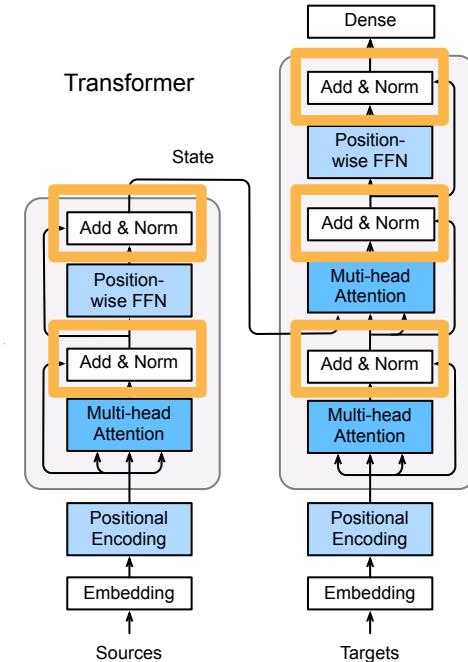
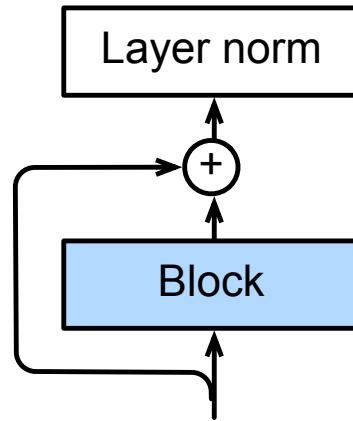
Position-wise Feed-Forward Networks

- 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

- Layer norm is similar to batch norm
- 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)`



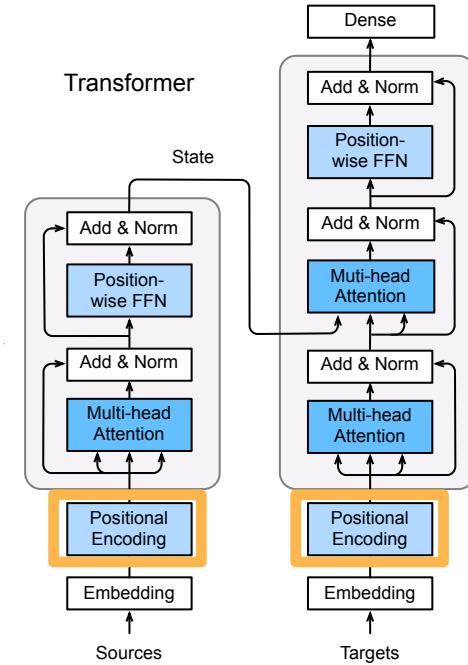
Positional Encoding

- 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

$$P_{i,2j} = \sin(i/10000^{2j/d})$$

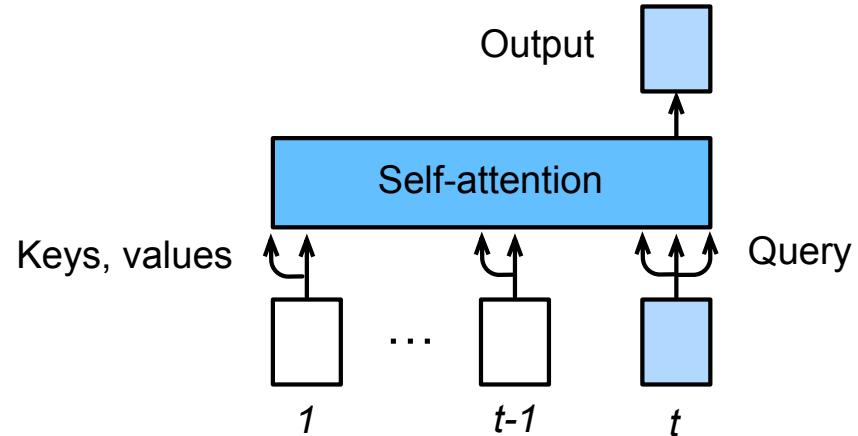
$$P_{i,2j+1} = \cos(i/10000^{2jd})$$

- Output $X + P$



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

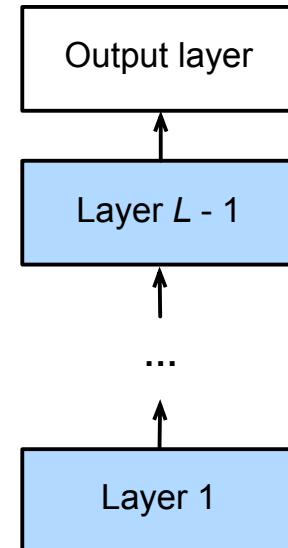
- Use pre-trained models to extract word/sentence features for the new task
 - E.g. word2vec or language model
- Often don't update the pre-trained models
- Need to construct a new model to capture the information needed for the new task
 - Word2vec ignores sequential information, language model only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

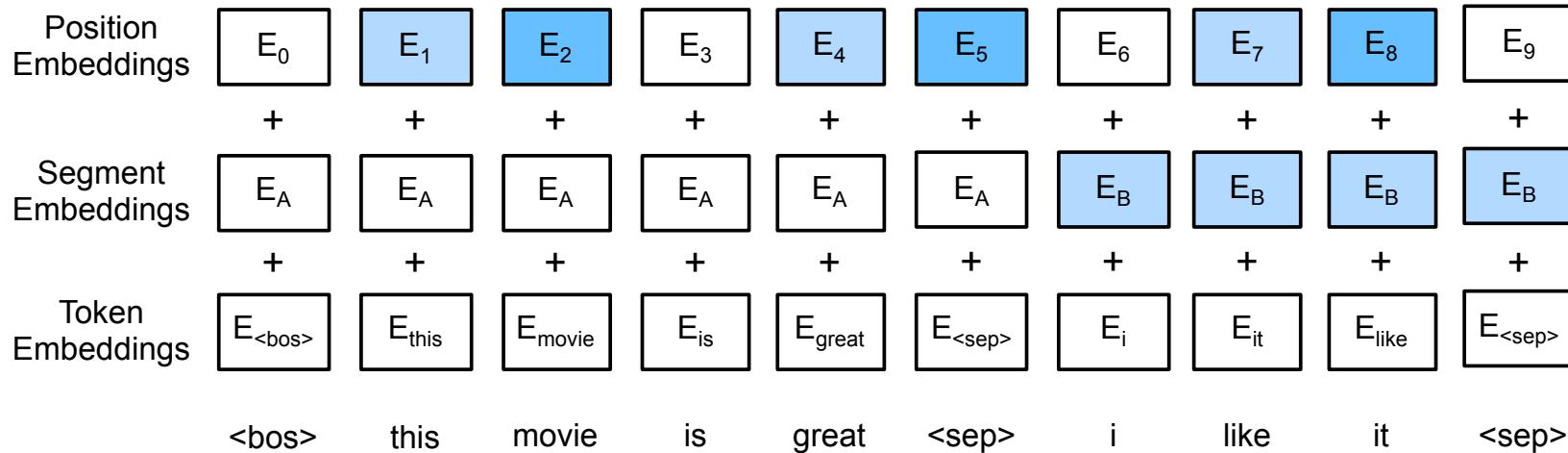


BERT Architecture

- 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

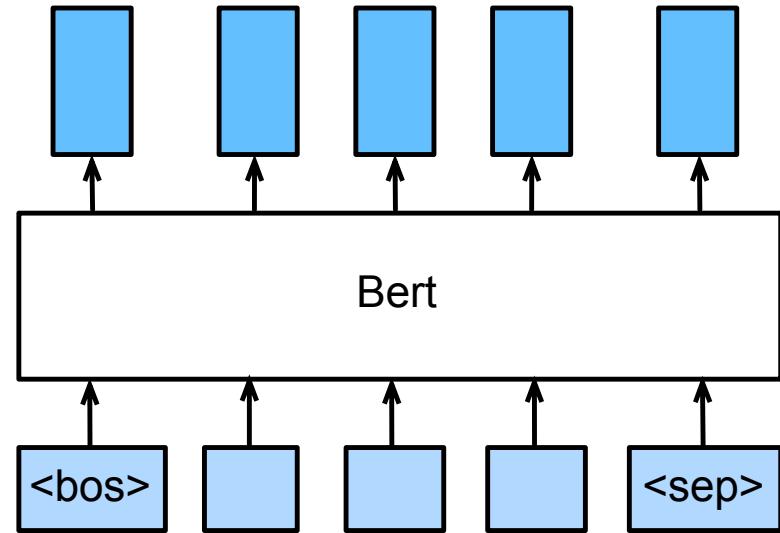
- Randomly mask (e.g. 15%) tokens in each sentence, predict these masked tokens
 - Transformer is bidirectional, which breaks the unidirectional limit of standard LM
- No mask token (<mask>) in fine tuning tasks
 - 80% of the time, replace selected tokens with <mask>
 - 10% of the time, replace with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

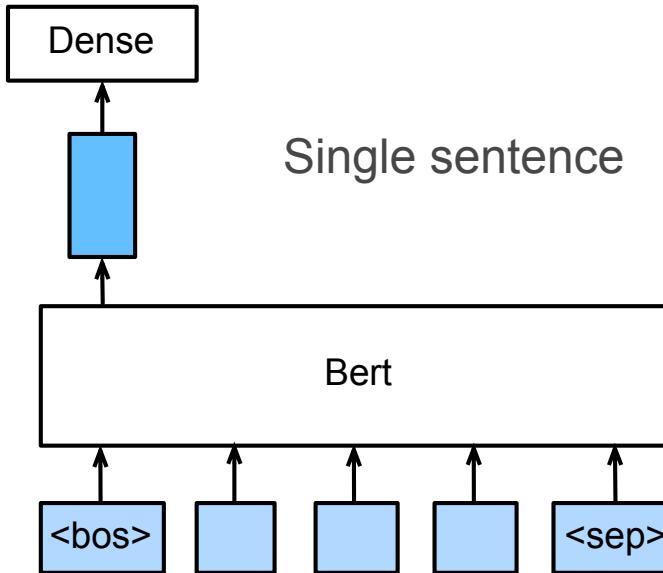
Bert for Fine Tuning

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

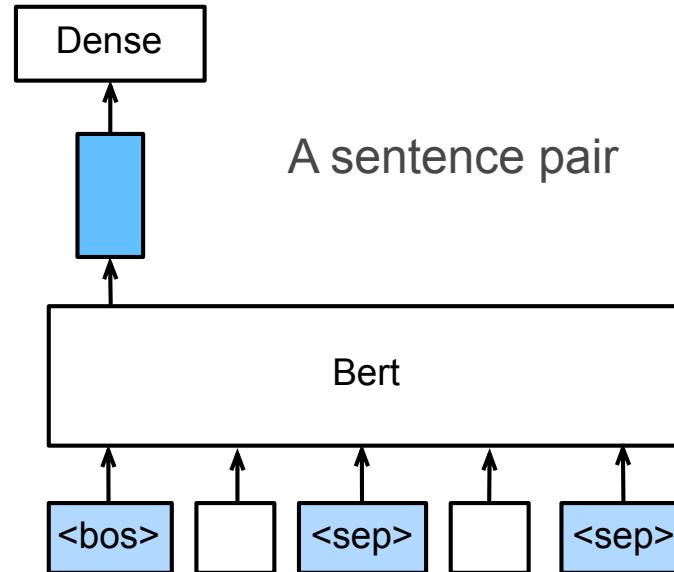


Sentences Classification

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



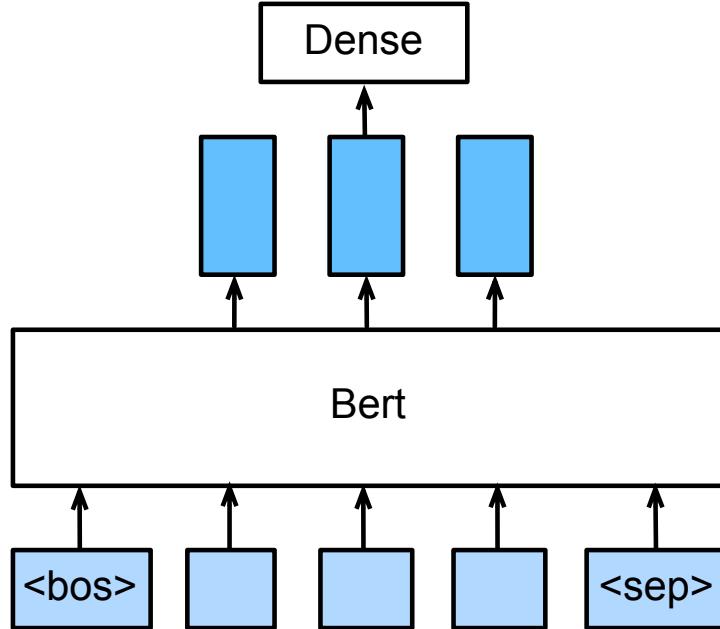
Single sentence



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

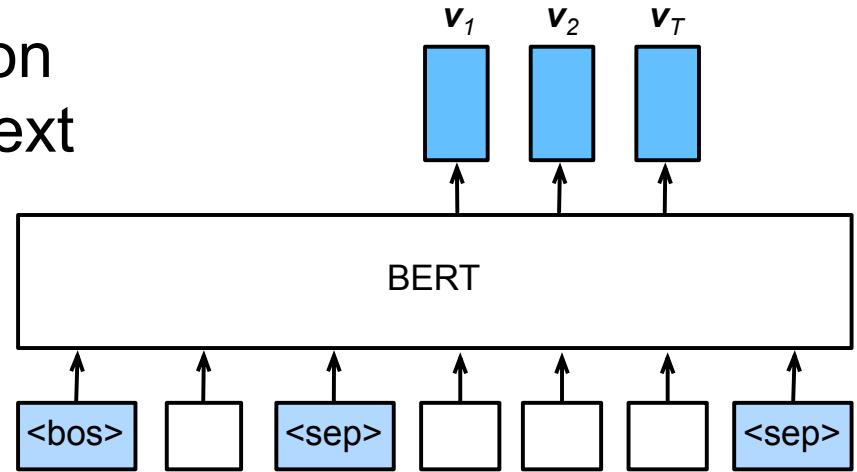


Question Answering

- 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...