

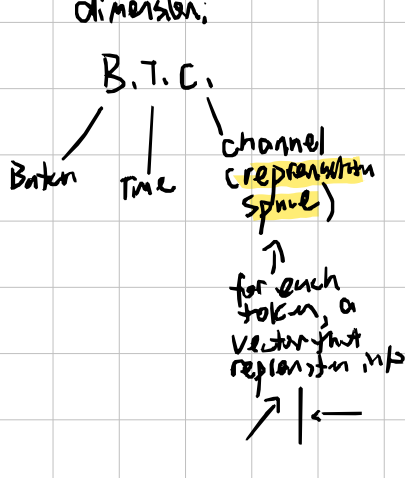
Mathematical Trick in self-attention

+ 8 tokens to talk to each other

1: should not future token
(only prev context to our time step)

↓ how?
average of preceding

i.e., average \rightarrow $M (M = M.E.)$
1st 2nd 3rd



• for each of $B \times T$ token, calculate average

" "
 $x_{\text{prev}} = \text{torch.zeros}(C, B, T, C)$

for b in range(B):

for t in range(T):

$x_{\text{prev}} = x[b, t+1]$

$x_{\text{bow}}[b, t] = \text{torch.mean}(x_{\text{prev}}, 0)$

$b, t, 0: t+1$

each token

↑
Oth dimension

EFFICIENCY

- MATH -

↓

Implementation, Transforms

Query: question
Key: meaning of myself \vec{v}
(what obj I mean)

+ Query: 'nn.Linear(C, head_size, bias=false)'

+ Key: 'nn.Linear(C, head_size, bias=false)'

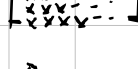
only matrix multiplication

$K = \text{key}(q)$
 $q = \text{query}(q)$ (loading input)

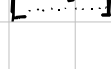
communication between queries and keys
(dot product)

$\rightarrow w_{ei} = q \otimes K \cdot \text{transpose}(-2, -1)$ (DOT PRODUCT), max dim remains fixed
 $\otimes (B, T, 16) \otimes (B, 16, T) \rightarrow (B, T, T)$

1. w_{ei} , masked_fill(tril=-D, float('-inf'))



2. $w_{ei} = F_{\text{softmax}}(w_{ei}, \text{dim}=-1)$



- 1. query purpose question: i.e., looking for vowels.
- 2. key purpose "answer": i.e., I am a vowel.
- 3. when similar: dot product gives high attentiveness score

$\rightarrow w_{ei}$ 4. upper-triangular masking to avoid future tokens communicating

5. normalize score i.e., -0.15 w/ softmax (how much information to propagate from information of the past)

* but we don't propagate tokens themselves, *
but the meaning of them; \vec{v}

softmax \rightarrow Meaning (or question) adjustment

Notes:

- There's no idea of space; thus we 'i' position-embedding
pos-emb = self-position-embedding-table(torch.arange(T, ...))
 $x = \text{token-emb} + \text{pos-emb}$
 $(B, T, C) \quad (T, C)$
- 'tril' (top-triangular mask) is not always the criterion, only for self-regression Transformers. Other e.g., sentiment prediction you don't 'tril', you use an "encoder" block, where no 'tril'
- we call this "decoder", since we're just decoding language to next token. you mask for future not just to prevent auto-regression. & it can be decoder-only or both etc.
- 'self-attention': we use x to generate all 'qkv'
- 'cross-attention': K, V come from other source & x . i.e., encoder block for external ctx (to read information from side) (different lang., different media)
- we use $\frac{1}{\sqrt{d_k}}$ for Attention(Q, K, V) = $\text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$
- as "scaled attention" to keep softmax diffused & saturated. (centrifugal)

Multi-head Attention:

- + Attend to information from different representation space
- + multiple relationships between tokens
- + multiple heads in parallel.

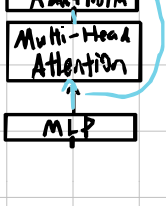
Multi-Layer Perceptrons

- + token-individual level
 - + process of updating meaning itself
 - + (factual info?)
- [nn.Linear(n_embd, n_embd)]
[nn.ReLU()]

Attention Blocks

- + Attention Heads \neq n-heads
- + MLP
- * Problem: Neural Network too deep \rightarrow optimization issue

1. Residual Connection



+ A skip connection with addition
(Attention forks off from residual pathing, while addition contributes grad equally)
 \rightarrow Gradient Hignany (supervisor to input) + Identity dropout (during training)

+ Projection \rightarrow Linear Transformation (After Head, Feed Forward)
'self.proj' = nn.Linear(num-heads * head-size, n_embd)

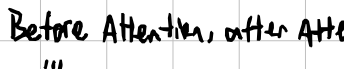
2. Inner-layer of FeedForward should be multiplied by 4

3. LayerNorm: Normalizes inputs \rightarrow 0 mean, 1 standard-deviation

\rightarrow LayerNorm BEFORE transformation

+ Per-token, w/ trainable parameters

- + conditioning the Addition properties of residual connection (we're adding things everytime)
- + decompose scale w/ direction (only direction)



1. Before Attention, after Attention

" "
def forward(self, x):
 $\hat{x} = x + \text{self.attention-heads}(\text{self.layer_norm}(x))$

" "
 $\hat{x} = \hat{x} + \text{self.feedforward}(\text{self.layer_norm}(\hat{x}))$

2. After Attention Blocks

4. Dropout:

- + when training, randomly reset some activation of neurons, but disable such behavior on inference
- + no co-adaption (don't depend on certain neurons)
- + scaling, against overfitting

1. after softmax (heads)
2. end of feedforward (back to residual pathing)
3. end of Multi-Head