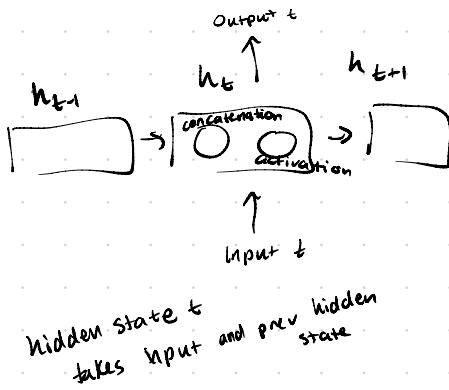


Efficient Mu  
transformers

# Part 1: Lecture 12

NLP: Discriminative Generative  
Pos/Neg

## RNN



Problem w/ RNN / LSTM

Y O(seq-length) for interactions b/wn tokens

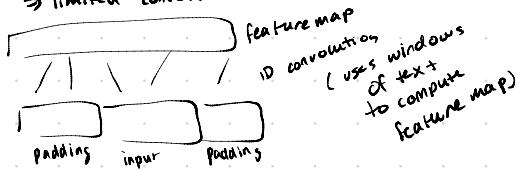
Y images have locality but languages don't

Y limited parallelism (has dependencies)

## Convolution Neural Network

→ no dependency b/wn tokens  
Y better scalability

→ limited context



## Transformer

### Tokenization

110 words  $\rightarrow$  162 tokens

~~~~~  
similar scale

### Embeddings

① one hot encoding each value is a #

101 00000

② word representation in vector

# Self Attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V$$

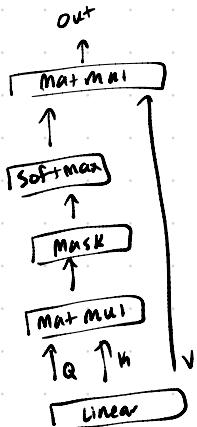
Softmax  
↳ logistic function  
 $\sigma(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$

⇒ multiply  $Q$  and  $K$  to get inner product, normalize by  $d_K$

⇒  $O(N^2)$  complexity

⇒ softmax ⇒ attention weights ⇒ so sum to 1 for each query ⇒ prob matrix

⇒ multiply attention weights by  $V$  (values for key)



$U \times N \times d$   
(one for each query)

$Q, K, V = (N \times d)$

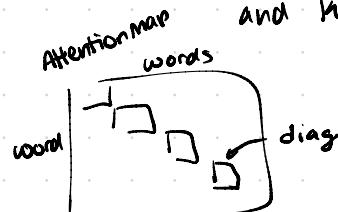
↓  
# of features

$$QV^T = N \times N$$

Attention = correspondence b/w words

$QKV$  doesn't rely on prev  
so done in parallel

↳ now each query and key are similar



↳ write self attention in pytorch

```

class SelfAttention(nn.Module):
    def __init__(self, n):
        super(SelfAttention, self).__init__()
        self.N = n
        self.Q = nn.Linear(d, d)
        self.K = nn.Linear(d, d)
        self.V = nn.Linear(d, d)
        self.softmax = nn.Softmax(dim=2)
  
```

```

def forward(self, x):
    queries = self.Q(x)
    keys = self.K(x)
    values = self.V(x)
    scores = torch.bmm(queries, keys.transpose(1, 2))
    attention = self.softmax(scores)
    weighted = torch.bmm(attention, values)
  
```

m batches  
n words  
d features per word

↳ linear( $d, d$ )  
↳ take input dim d  
and project to output of dim d  
↳ take input embeddings + put into  $Q, K, V$   
output = input  $\times$  weights $^T$  + bias

↳ softmax function applied on 3rd dim of vector so  $m \times n \times n$  applied here

↳ linear trained through back propagation + gradient descent

softmax(dim=2) → 3rd dim  
 $x = \text{torch.tensor}([$  so along each  
 $\begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix},$   $v_{\text{th}}$  should add  
 $\begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$  up to 1  
 $\dots,$   
 $\begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$   
 $\dots])$

Class Multihead Attention (nn. Module)  $\xrightarrow{\text{embedding dim}} = d_k \cdot h = d_v \cdot h$   $d_{\text{model}} = d_{\text{-embeddings}}$   
 $\text{def } \text{--init--}(\text{self}, d_{\text{-model}})$

`super(MultiHeadAttention, self). --init--(self, d_k=d_k, d_v=d_v, d_model=d_model)`

self.concat = nn.Linear( $d_{\text{out}} \times h$ ,  $d_k \times h$ )  $W^0$   
 already concatenated but apply linear transformation via  $W^0$

def forward(self, Q, K, V, mask) → can be  $x, x, x$  we don't need to apply 2 linear transformations  
 $\xrightarrow{x \rightarrow Q} \xrightarrow{Q \rightarrow W_Q Q}$   
batch-size, seq-length, d-model = .size()  
one linear transformation

queries = self.query(Q)  $\rightarrow Q \cdot N^Q$

queries = queries.reshape(batch\_size, seq\_length, h, d - k)

queries.transpose(1, 2) → (bs, h, seq\_length, d\_k)

```
Scores = torch.bmm(queries, key.transpose(-2,-1)) / d_k0.5
```

`scores = scores.masked_fill(mask=0, -1e-10)` seq.length x d\_k  $\rightarrow$  d\_k x seq.length

$$\text{Attention} = \text{nn.softmax}(\text{scores}, \text{dim}=-1) \Rightarrow (\text{bs}, \text{h}, \text{seq}, \text{seq})$$

add weighted = torch.bmm (attention, values)  $\rightarrow$  (bs, n, seq, d-k)  $\rightarrow$  bs, seq, n, dk  
 dropout weighted. transpose(1, 2)

weighted.reshape(b, seq, h x dk) concat(h<sub>1</sub> .. h<sub>n</sub>)

self. out (weighted)  $\xrightarrow{\text{d model}}$  Concat &  $W_0$

## Multi Head Attention

↳ each head captures dif semantics

$$\text{Multihead}(Q, K, V) = \text{concat}(\text{head}_1, \dots, \text{head}_n) W^O$$

each head produces attention map

$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$$W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_{K,V}}$$

## Attention Masking

↳ one token can only see tokens before it

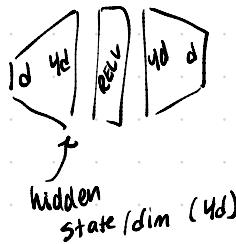
↳ causal self attention rather than global self attention  
(masking) (discriminate task)

## Feed Forward Layer

→ add non-linearity

→ two layer MLP w/ larger hidden state (Inverted bottleneck) + RELU/GELU activation

$$\text{FFN}(x) = \max(0, \text{ReLU}(W_1 + b_1) W_2 + b_2)$$



Inverted Bottleneck

- ↳ inverted residual blocks
- ↳ expand dim then reduce to learn richer feature representation



## Layer Norm / Residual Connection

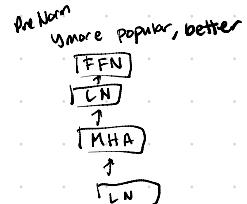
- performs normalization for each token
- then applies learnable affine transformation (linear mapping)
- for each token embedding, mean = 0, var = 1

→ not batching like CNN

$$y = \frac{x - E[x]}{\sqrt{Var(x) + \epsilon}} + \gamma + \beta$$

normalize

learnable  
affine transformation  
(linear)



## Positional Encoding

- far away dif meaning

$$\Rightarrow (i) p_t^i = f(t) = \begin{cases} \sin(w_k \cdot t) & \text{if } i=2k \\ \cos(w_k \cdot t) & \text{if } i=2k+1 \end{cases}$$

$w_k = \frac{1}{10000} \cdot 2\pi k$

$i$  (feature dim)  
 $t$  (token index)

unique encoding across  
each token position in sentence  
and across dif features

- Transformers  
better accuracy, less training cost

## Design Variants

Encoder + Decoder (T5), encoder-only (BERT), decoder-only (GPT)

(see everything)

Absolute Pos encoding → Relative Pos encoding

NV Cache Optimization

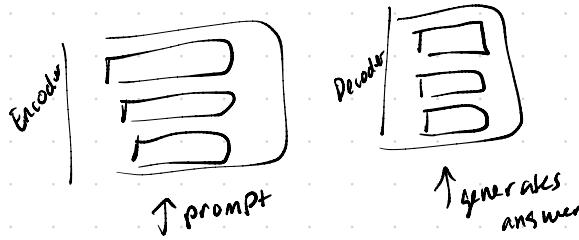
FFN → GLU (gated linear unit)

# Encoder - Decoder Architecture (T5 - Text to Text Transformer)

Text to Text Model for transfer learning

- Translation
- Acceptability
- Semantic Similarity (0-5)
- Summarization

↓  
pretrained  
for diverse  
set of tasks



## Encoder Only (BERT)

Two pretraining objectives

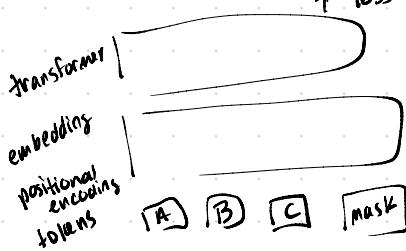
- ① Masked (more popular)  
15% masked at random  
train to predict words

pretraining stage  
→ train on unlabeled  
data to learn general language  
(special token)  
mask token! representations

- ② Next sentence Prediction

sentence B is sentence after A  
classification  
↑ loss

then pretrained model finetuned  
for downstream tasks



Full attention → can see all words, prev + future



## Decoder-Only Model (GPT)

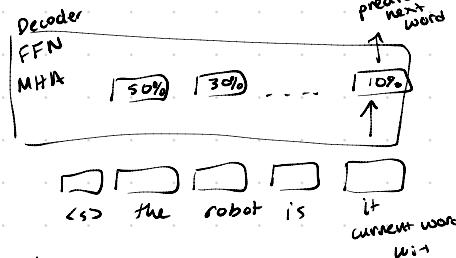
↳ generation pre-trained transformer

next word prediction

$$L_i(u) = \sum_i \log (P(u_i | u_{i-k} \dots u_{i-1}) \theta)$$

→ only see prev words

→ generating next word



Attention Mask



→ fine-tune for downstream

tasks OR zero shot/few shot  
(doesn't require finetuning)

## Absolute Positional Encoding

↳ in embedding thus ( $K, Q, V$ )

## Relative Positional Encoding

↳ impact attention score / bias to the attention map, or modifying query of key

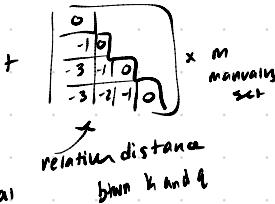
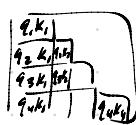
↳ can generalize to sequence length not seen in training

⇒ ALiBi

(offset to attention matrix

↳ train w/ short text

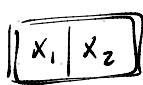
↳ test w/ long



⇒ Rotary Positional Embedding

⇒ split embedding  $d$  into  $d/2$  pairs  $\Rightarrow$  2D coord

⇒ apply rotation according to pos  $m$



$\Rightarrow m \theta_i$  rotation angle  
position  
 $1, 2, 3 \dots$

2D coord

$$\theta = \sum \theta_i = 1000^{-2(l-1)/d}, i \in [1, 2 \dots, d/2]$$

$$R \circ PE(q_j, m)$$

phase any of inner product of two complex vectors is phase difference

$$Rope(q_j, m), Rope(k_j, m)$$

$$= R \circ PE(q_j, k_j, m-n)$$

difference, relative

$$\theta = \sum \theta_i = 1000^{-2(l-1)/d}, i \in [1, 2 \dots, d/2]$$

# of dif tokens in sequence (large for long text)

$$f_{\theta_{i,13}}(x_m, m) = R_{\theta_{i,m}}^d W_{t_{i,13}} x_m$$

$$R_{\theta_{i,m}}^d = \begin{pmatrix} \cos(m\theta_i) & -\sin(m\theta_i) \\ \sin(m\theta_i) & \cos(m\theta_i) \end{pmatrix}$$

double length of sequence

yrotate by  $\theta'_i = \theta_i / 2$

bc differentiate by more tokens

$m \rightarrow 2m$  tokens in sequence

UV cache optimizations

GPT - KV Cache

$\Rightarrow$  store key values of all prev tokens

$\Rightarrow$  only need current query

Decoder

young attends to past tokens

Llama 2 7B, UV Cache Size

$$BS \times 32 \times 32 \times 128 \times N \times 2 \times 2 \text{ bytes} = 512 \text{ GB} \times BS \times N$$

minibatch   layers   heads   N   words   length   K/V  
 ↓              ↓              ↓              ↓              ↓              ↓  
 # of users    sequence    width    # of tokens

$$BS=1 \quad n\_seq = 512 \quad 1.25 \text{ GB of mem}$$

$$BS=1 \quad n\_seq = 4096 \quad 10 \text{ GB} (\sim \text{a paper})$$

$$BS=16 \quad n\_seq = 4096 \quad 160 \text{ GB} \quad (\text{two A100 GPU!})$$

↑              Heads for Query    Heads for KV     $\xrightarrow{\text{GPU} = 8 \text{ GB of memory}} \text{cost } \$20k$   
 # of words    N              N

Multi head Attention

Multi Query Attention

Grouped Query Attention  
(between multi head + multi query)

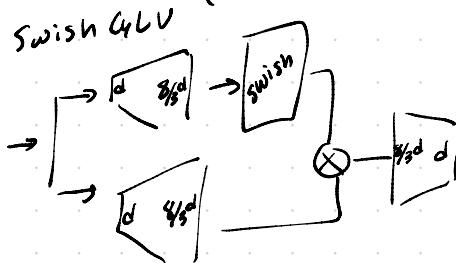
1     $\rightarrow$  avg val for all KV into 1

$G = N/8$      $\Rightarrow$  group every two heads, matches accuracy of multi head attention

FFN  $\Rightarrow$  GLU (Swish + Gated linear Unit)

$$\text{FFN}_{\text{SWiGLU}}(x, W, V, W_2) = (\text{Swish}(xW) \otimes xV)W_2$$

(3 layer)



GLU

$$y = \frac{x}{1-e^{-1.702x}}$$

swish

$$y = \frac{x}{1+e^{-x}}$$

similar shapes

Llama introduced Swish + GLU, RoPE  
relation positioning

Emergent: large models

modified arithmetic, word unscrambling  
few shot

---

Original Transformer

Attention is all you need

Encoder + Decoder  $\Rightarrow$  specific task

Dropout to output of sublayer before added to sublayer input + normalized (LN)  $P_{\text{drop}} = .1$

In Residual Connection:

$$\text{residual\_input} = x + \text{dropout}(\text{sublayer}(x))$$

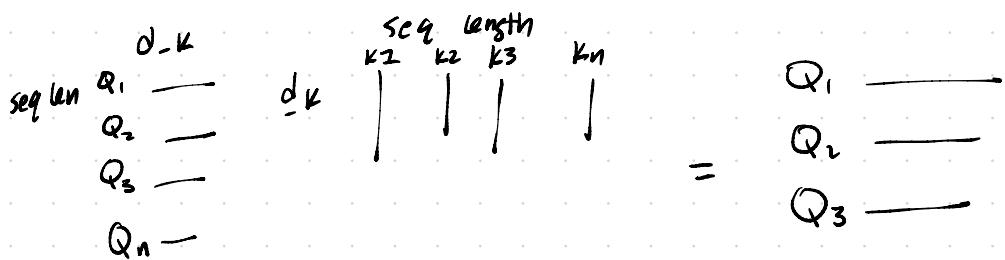
$$\text{output} = \text{LN}(\text{residual\_input})$$

In Sublayer:

- $\rightarrow$  after POS encoding
- $\rightarrow$  after RELU in FFN  $\rightarrow$  but before second so features not lost due to dropout
- $\rightarrow$  after softmax in Multi-head attention

y before multiplying times values

} after activation function



Attention scores  $\rightarrow$  query seq-length  $\times$  key sequence length

$\times$  value  $\rightarrow (q\text{-seq} \times \text{key-seq}) \approx (v\text{-seq} \times d_{\text{model}})$

Cross entropy loss

$\rightarrow \infty$  if predict 0 for actual

$$\sum_{x \in X} p(x) \log(q(x))$$

actual    pred

`super(class, self).__init__()`

## embeddings

nn.Embeddings (vocab-size, d-model)

## positional encoding

$$\text{even\_indices}(2i) : PE(\text{pos}, 2i) = \sin\left(\frac{\text{pos}}{1000^{2i/d\text{-model}}}\right)$$

$$\text{odd\_indices}(2i+1) : PE(\text{pos}, 2i+1) = \cos\left(\frac{\text{pos}}{1000^{2i/d\text{-model}}}\right)$$

y add to embeddings, absolute PE

PE pos (token pos / 0-seq-length) i (feature pos (0-d-model))

self.dropout = nn.Dropout (Pdrop)  $\rightarrow$  1

position = torch.arange (0, seq-length). unsqueeze(1)  $\rightarrow$  seq-length x 1

div = torch.exp(torch.arange (0, dmodel, 2).float() \* math.log(1000) / dmodel)  
i going from 0 to d-model, steps of 2  $e^{\frac{\log(1000) \cdot 2i}{d\text{-model}}}$

pe = torch.zeros (seq-length, d-model)

pe[:, 0:2] = torch.sin (pe / div-term)

$\downarrow$   
[:, 1:2]

$\uparrow$   
odd bc 1 out of 2

pe = pe.unsqueeze(0)  $\rightarrow$  (seq-length, d-model)  $\rightarrow$  (1, seq-length, d-model)

self.register\_buffer ('positional-encoder', pe)  $\rightarrow$  not a model param/not trainable

y self.positional-encoder but stored as part of the module

def forward (self, x)

x = x + self.pe[:, seq-length, :]. requires\_grad(False)

return self.dropout(x)

# LN

normalize token the n linear transformation

Layer Norm ( $x + \text{sublayer}(x)$ )

$$\alpha \left( \frac{x - E(x)}{\text{var}(x)} \right) + \text{self.bias}$$

+  $\epsilon$

.00001 to avoid divide  
by 0

residual output

use both  $x + \text{sublayer}(x)$   
to avoid the vanishing  
gradient problem

$$\alpha = \text{nn.Parameter}(\text{torch.ones}(1))$$

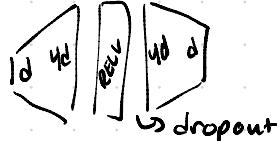
$\uparrow$  1 dim tensor

$$b = \text{nn.Parameter}(\text{torch.zeros}(1))$$

def forward(self):  
mean =  $x.\text{mean}(\text{dim}=-1, \text{keepdim=True})$  ✓ don't change # of input dim  
std =  $x.\text{std}(\text{dim}=-1, \text{keepdim=True})$

## FFN

$$(\max(0, xW_1 + b_1))W_2 + b_2$$



$$dff = d\_model \times 4$$

self.linear1(d\_model, dff)

self.dropout(Pdrop)

self.linear2(dff, d\_model)

forward(x)

self.linear2 / self.dropout(torch.relu(self.linear1(x)))

## Residual Connection

residual\_output =  $x + \text{dropout}(\text{sublayer}(x))$

forward(x, sublayer)

self.layer\_norm( $x + \text{self.dropout}(\text{sublayer})$ )  
 $\uparrow$   
residual

## Encoder

$N=6$  layers

↳ each layer is encoder block

## Class Encoder

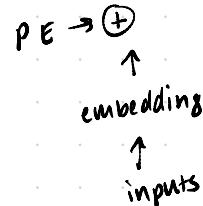
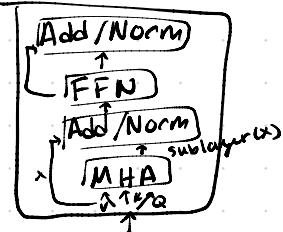
self.layers = nn.ModuleList(...)

EncoderBlock range( $N$ )

forward ( $x, mask$ )

for layer in layers  
layer( $x, mask$ )

## Encoder Block



## Class Encoder Block

self.MHA

self.FFN

self.residual\_connections = nn.ModuleList([ResidualConnectionL.range(2)])

forward ( $x, mask$ ) → don't need to compute QKV

attention ( $x, x, x$ )  $W_{QK} = Q$  separately

$q \leftarrow x \cdot W_q$   $k \leftarrow x \cdot W_k$   $v \leftarrow x \cdot W_v$  don't need two linear transformation  
→ residual  $\rightarrow$  FFN  $\rightarrow$  residual [?] (redundant)

src-mask  
(encoder input)

tgt-mask  
(decoder input)

## Decoder Block

output prob

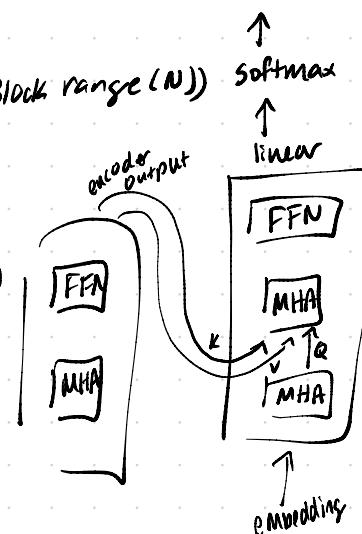
## Class Decoder

self.layers = nn.ModuleList(DecoderBlock range( $N$ ))

forward

for layer in layers

layer( $x, encoder\_output, tgt\_mask, src\_mask$ )



## Decoder Block

MHA1

MHA2

FFN

residual = nn.ModuleList([ResidualConnection, ... range(3)])

forward

$x = \text{MHA}_1(x, x, x, \text{tgt\_mask}) \rightarrow \text{residual}$

$x = \text{MHA}_2(x, \text{encoder\_output}, \text{encoder\_output}, \text{src\_mask}) \rightarrow \text{residual}$

$x = \text{FFN}(x) \rightarrow \text{residual}$

## Projection Layer

def \_\_init\_\_

self.proj = nn.Linear(d\_model, vocab\_size)

def forward(x)  $\rightarrow x.\text{shape} = [bs, \text{seq\_length}, d\_model]$

$\text{torch.logSoftmax}(\text{self.proj}(x), \text{dim}=-1)$

## Transformer

def \_\_init\_\_(self)

self.encoder  $\leftarrow$  nn.Module

so call

self.decoder self.encoder(x) will run forward

self.src\_emb

self.tgt\_emb

self.src\_pos

self.tgt\_pos

self.proj\_layer

def encode(self)

embed  $\rightarrow$  pos  $\rightarrow$  encoder

def decode(self)

embed  $\rightarrow$  pos  $\rightarrow$  decoder

def project

self.projection\_layer

## build transformer

```
for p in transformer.parameters():
    if p.dim() > 1 → only for tensors w/ dim > 1
        nn.init.xavier_uniform_(p) ← glorot initialization
                                         maintain variance through
                                         params to prevent vanishing/gradient
```

## tokenizer

check if tokenizer file otherwise tokenize

```
tokenizer = Tokenizer(WordLevel(unk_token='UNK'))
trainor = WordLevelTrainer(special_tokens=[“UNK”, “PAD”,
                                            “[SOS]”, “[EOS]”, min_frequency=2) → min # of times token appears in corpus
                                            to be in vocab
tokenizer.train_from_iterator(sentences, trainor=trainor)
tokenizer.save(filepath)
tokenizer.encoder(sentences)
```

## get Data

```
load dataset
DataLoader(data, batch_size, shuffle=True)
```

## Causal mask

```
mask = torch.triu(torch.ones(1, size, size), diagonal=1).type(torch.int)
```

```
model = build_transformer(vocab_size_src, vocab_size_tgt, seq_len_src, seq_len_tgt, d_model)
```

## Process Data (Dataset)

```
def __init__(self, tokenizer_src, tokenizer_tgt, seq_length):
    def __getitem__(i):
        pair
        y get padding
        src_tokens = tokenizer_src.encode(src).ids
        tgt_tokens
        src_padding = seq_len - 2
        tgt_padding = seq_len - 1
        src_tokens = src_tokens + [sos, eos]
        tgt_tokens = tgt_tokens + [sos, eos]
        sos = torch.tensor(tokenizer_tgt.token_to_id("[sos]"))
        if src_tokens[-1] == eos:
            src_tokens[-1] = sos
        if tgt_tokens[-1] == eos:
            tgt_tokens[-1] = sos
        check if too long ie padding < 0
        get encoder input
        torch.cat(sos tensor, EOS, padding
                  (tokens))
        get decoder input
        SOS tensor padding
        get label (correct decoder)
        tensor, EOS, padding
    return {
```

encoder input  $\rightarrow$  dim = seq-length  
decoder input  
encoder mask  $\rightarrow$  encoder\_input + padding token  
decoder mask  $\rightarrow$  not padding + causal-mask (seq-length)  
label  
src\_text, tgt\_text

## get\_next\_tokens()

```
sos_idx = tokenizer_tgt.token_to_id("[sos]")
y token for SOS
```

EOS\_idx

```
encoder_output = model.encode(source, source_mask)
```

```
decoder_input = torch.empty(1, 1).fill_(sos_idx).type_as(source).to(device)
```

while True

    decoder\_input.size(1) = max length  
    break

decoder\_mask = causal\_mask (decoder\_input.size(1))

out = model.decode (decoder\_input, encoder\_output, tgtmask, srcmask)

prob = model.projection (out[:, -1])  
    ↑ just  
    for last  
    decoder  
    input

max\_val, next\_word = torch.max (prob, dim=1)

decoder\_input = torch.cat ([decoder\_input, torch.empty(1, 1).fill\_(next\_word.item())  
    .item().tensor>python  
    .to(device)])

If next\_word = eos\_idx

    break

return decoder\_input.squeeze(0)

    y make a sequence of tokens

def run\_validation

model.eval()

with torch.no\_grad():

batch[encoder\_input]

batch[encoder\_mask]

result = get\_next\_tokens (model, -)

batch[src\_text]

batch[tgt\_text]

model\_out\_text = tokenizer.tgt.decode(result.detach().cpu().numpy())

train model

```
vocab_size = tokenizer.get_vocab_size()  
optimizer = torch.optim.Adam  
loss
```

label smoothing  $\epsilon_{IS} = .1$

Y hurts perplexity

Y increase accuracy (BLEU score)

↓  
Bilingual Evaluations  
Understudy

Y add prob to all classes, correct class less,  $\Rightarrow$  regularization

## Part 2: Lecture 13

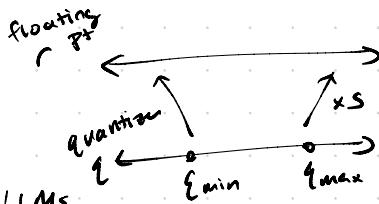
### Quantization

14 bit range  $\rightarrow$  integer values

$\Rightarrow$  important for CNN, destroy performance for LLMs

$\Rightarrow$  LLMs have outliers in activations so can't quantize  
(0-70)

b) weights easy to quantize, range 0-1



### Smooth Quant / W8A8

$\Rightarrow$  activation outlier mitigation

Accuracy remains the same

$$Y = XW \xrightarrow{10} \begin{matrix} \text{enlarge weight,} \\ \text{make activation} \\ \text{smaller} \end{matrix}$$

\* used in post training

(not pretraining)  
hard to learn quantization levels  
(in pretraining  $\rightarrow$  degradation)

weight  $\rightarrow$  harder to quantize

activation  $\rightarrow$  easier to quantize

$$X \left[ \begin{array}{|c|c|c|} \hline 1 & -6 & 2 & 6 \\ \hline -2 & 8 & 1 & 9 \\ \hline \end{array} \right] \quad W \left[ \begin{array}{|c|c|c|} \hline 0 & -1 & -2 \\ \hline 1 & -1 & -1 \\ \hline 2 & -1 & -2 \\ \hline 1 & 1 & -1 \\ \hline \end{array} \right] \quad \xrightarrow{\quad} \quad S = \left[ \begin{array}{|c|c|c|} \hline 1 & 4 & 1 & 3 \\ \hline \end{array} \right]$$

$$\hat{X} = X \text{diag}(S) \quad \hat{W} = \text{diag}(S) W \quad \xrightarrow{\quad} \begin{matrix} \text{no extra overhead} \\ \text{runtime} \end{matrix}$$

$\Rightarrow$  computed during compile time

$$Y = (\underbrace{X \text{diag}(S)}_{\text{fold into activation function}})^{-1} (\underbrace{\text{diag}(S) W}_{\text{compute offline}}) = X W = \hat{X} \hat{W}$$

$$s_j = \max(|x_j|) / \max(|w_j|)^{1-\alpha} \quad j = 1, 2, \dots, c \quad \# \quad \text{scaling factor is constant}$$

$\Rightarrow$  Reduces Memory By  $\frac{1}{2}$

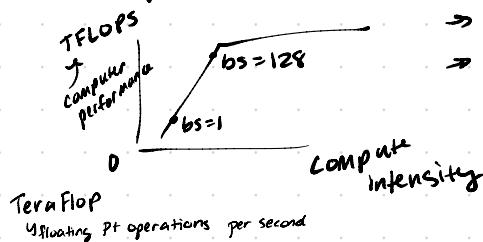
$\Rightarrow$  lower latency

$\Rightarrow$  use quantization in all compute intensive operations, linear s, BMMs

lower perplexity  $\Rightarrow$  better  
Y - log likelihood  
Y perplexity = X  $\Rightarrow$  choosing  
between equally likely  
options  
Y measurement for  
accuracy of predicting  
the next word

## Single batch Serving

W8A8 quantization good for batch serving



→ W8A8 good for batch serving

→ small batch  $\Rightarrow$  memory bounded  $\Rightarrow$  need low bit weight  
(time to compute restricted by memory)  
only quantization

Why weight only quantization:

- Main bottleneck is storing weights (memory)
- Smooth Quant (W8A8) good for batch not  $bs=1$
- focuses on quantizing memory intensive high precision format (FP16)

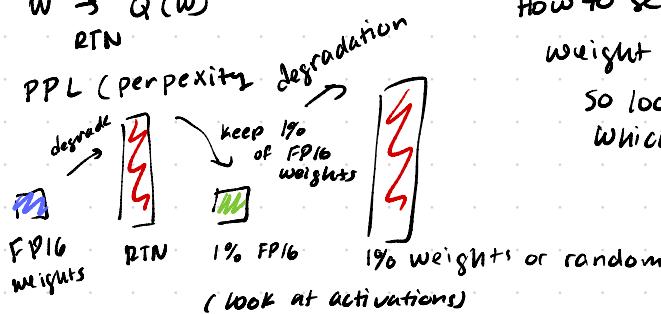
Low bit weight Only Quantization

↳ Leaves activations in high precision format (weights bigger matrix)  
 $x \cdot w = y$

- ↳ compute cheap, memory expensive
- ↳ RTN (round to nearest)

$$W \xrightarrow{\text{perplexity degradation}} Q(W)$$

RTN



How to select important weights?

weight is mult times activations

so look at activations to choose  
which to save

keep big outliers

A WQ for low-bit weight Quantization

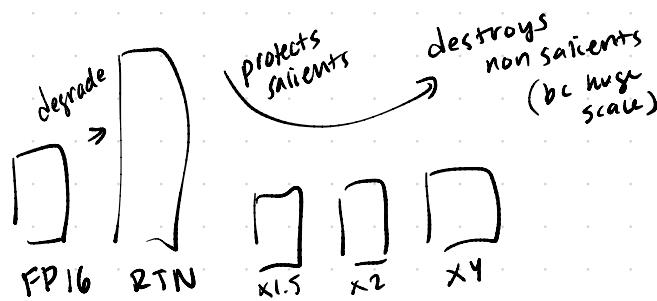
protecting salient weights by scaling

$$w \mathbf{x} \rightarrow Q(w_s)(s^{-1} \mathbf{x})$$

multiply salient channels  $w_s$ ,  $s > 1$  reduces quantization error

has significant impact on output  
(identify via activation)

better than prev bc  
before 1% are unquantized  
→ hard to write kernels



A WQ Match:

Consider linear layer:  $y = w \mathbf{x} \rightsquigarrow y = Q(w) \mathbf{x}$

$$Q(w) = \Delta \text{Round}(w/\Delta) \quad \Delta = \frac{\max(|w|)}{2^{N-1}}$$

scaled version:  $Q(w \times s) / (x/s)$

A WQ

Y also works for multimodal LLMs

## Tiny Chat

- light weight chatbot for LLM on the edge
- 4 bit AWD model
- ↑- quantize

# Pruning and Sparsity

## Pruning (Wanda)

→ consider activation when pruning weights

→  $|weight| \propto \|activation\|$  as pruning criteria

↳ not just weights

## Attention Sparsity (spAttention)

↳ don't need all tokens

↳ small cumulative  $\Rightarrow$  pruned tokens away  
attention score

## Cascade Pruning

↳ prune tokens + heads

If  $QK^T$  is small don't fetch V and can prune away

→ progressive quantization  $\Rightarrow$  gradually reduce our time

↳ lower precision

→ low precision first if not confident, no one that is significant  
then fetch the high precision

## Attention Sparsity

→ only keep local tokens and heavy hitter ( $H_2$ ) tokens in KV cache  
↳ if attention score is small  $\Rightarrow$  prune away

## Contextual Sparsity

↳ context matters

## DejaVu (Input dependent sparsity)

→ static sparsity (locations of 0 weights are predefined pretraining)

→ hurts accuracy w/ medium/high sparsity  $\Rightarrow$  replace w/ dynamic sparsity

## Dynamic Sparsity:

→ contextual sparsity (small, input dependent sets of redundant heads)

→ rate  $\Rightarrow$  inference without hurting model quality

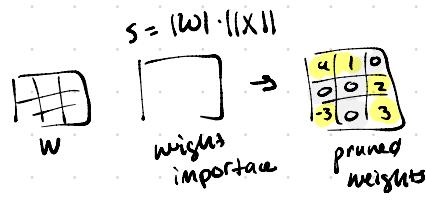
## Mixture of Experts

### MoE

→ for each token don't activate entire network

→ some total params, reduce inference cost per token

More experts  $\rightarrow$  larger total Model size  $\rightarrow$  better perplexity/lower loss



## Auto regressive

→ automatically predict next tokens  
from previous inputs to system

Capacity factor  $C \Rightarrow \frac{\text{tokens per batch}}{\text{number of experts}} \times \text{capacity factor}$

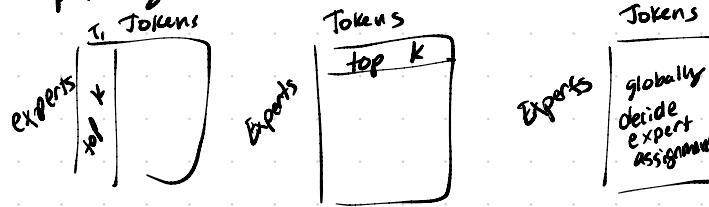
6 tokens, 3 experts:  
at most

$C=1$ , every expert can process  $\frac{6}{3} \sim 1 = 2$  tokens

$C=1.5 \quad \frac{6}{3} \cdot 1.5 = 3$  tokens

↳ have extra slack

## Sparingly Activated Token Routing Mechanisms



### VLLM and Paged Attention

→ runtime allocation of memory is slow

→ preallocate → but wasteful

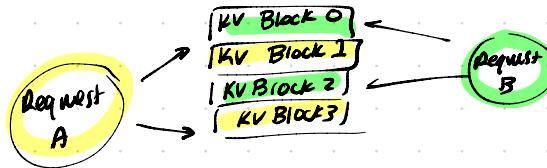
waste

→ internal fragmentation (over allocation due to unknown output length)

→ reservation (not used in current step but future step)

→ external fragmentation (diff. seq. lengths)

→ like OS system → virtual memory + paging

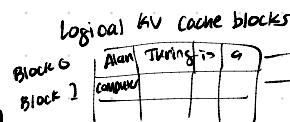


### Paged Attention

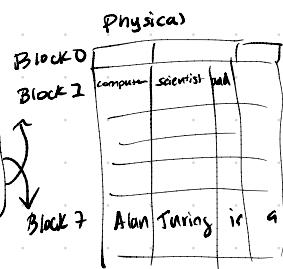
$y$  tokens per block

Maximum waste is  $3^y$

↳ no longer worry  
about fragmentation!



| Physical | Block Table |             |
|----------|-------------|-------------|
|          | Physical    | # of splits |
| Block 0  | 7           | 4           |
| Block 1  | 1           | 3           |



Paged Attention is good for  
Several suggestion, parallel sampling

Streaming LLM

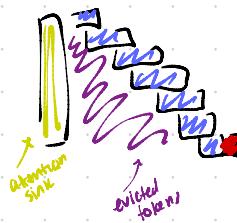
- Memory grow exponentially
- multi round dialog
- Window Attention → memory constant → perplexity grows  $\Rightarrow O(T^2)$  PPL $\downarrow$ 
  - 4 tokens ahead
- Dense Attention → everything → bad efficiency + performance  $\Rightarrow O(TL)$  PPL $\downarrow$
- sliding window w/ recompilation of KV cache (bad perplexity)
  - for each incoming token  $\Rightarrow O(TL^2)$  PPL $\uparrow$

Attention Sink (found in SpAttention paper)

- heavy attention to first token
- due to softmax function → have to 1 for all contextual tokens
- initial tokens become sinks due to high visibility
- if not quite related, jump attention on first token

Solution → Streaming LLM

- always keep KV + sliding windows (w/ out recompilation) of attention sinks KV to stabilize



How many attention sink need

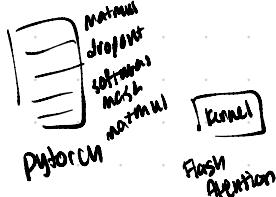
- 4-8 → plateau so use 4
- 2 attention sink → same performance

## Flash Attention

$N \times N$  attention matrix

↳ slow high bandwidth memory

↳ kernel fusion  $\Rightarrow$  merge into one kernel



## Speculative Decoding

Small model  $\Rightarrow$  auto regression, one by one

Decoding  $\Rightarrow$  token by token  $\Rightarrow$  memory bounded

Large model  $\Rightarrow$  parallel, verification of words  $\Rightarrow$  many tokens at once, no memory bottleneck  
↳ so decoding is parallel  $\Rightarrow$  accept, inject correct

Draft Model  $\Rightarrow$  small

Target model  $\Rightarrow$  big (try to accelerate, slow w/ lots of params)

## LoRA

low rank adaptation of LLM

instead of full rank  $\Rightarrow$  learn low rank  $\Rightarrow$  lower number of training params

## Q LoRA

↳ LoRA w/ quantized base model weights

↳ double quantization, quantize scaling factors

↳ new data-type for LLM weight quantizer NF4



## LoRA Adapter

$\Rightarrow$  train to new scenarios w/ out changing original params

## Prompt Tuning

↳ add prompt to task  
↳ train prompt for diff tasks  $\Rightarrow$  in training

↳ prompt engineering

↳ model gets larger  $\Rightarrow$  same accuracy

↳ mix learned prompts into single batch

↳ same model to handle diff tasks

↳ trained w/ task prompt

Tensor RT Public Release

↳ all of these methods + tensor parallelism + pipeline parallelism

