

Transformers

ChineJune

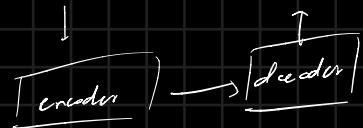


The Illustrated Transformer

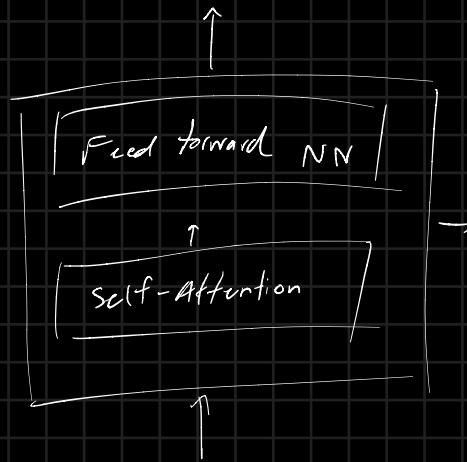
Lends itself to parallelization

As a machine translation application

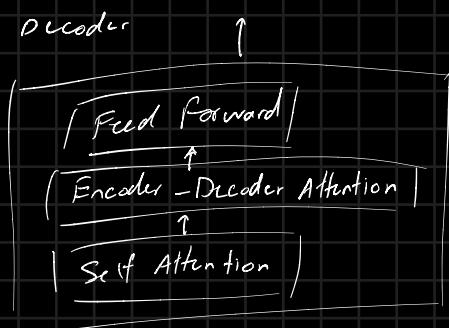
Je suis étudiant I am a student



Encoder: do not share weights



Self-Attention: Helps encoder look at other words in the input sentence



Turn input word into embedding vector

All encoders receive a vector of size 512.

Data flows up.

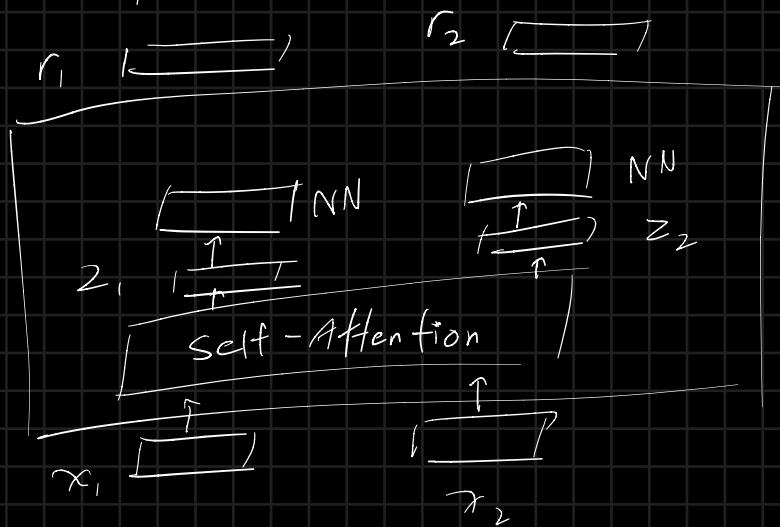
Size of SIZ is a hyperparameter \Rightarrow max sentence size.

words in each position flows through its own path.

\Rightarrow various paths can be executed in parallel.

Encoder : Self-Affection \rightarrow feed forward neural networks

\rightarrow Next encoder



What is Self-Attention?

Associate words with other words

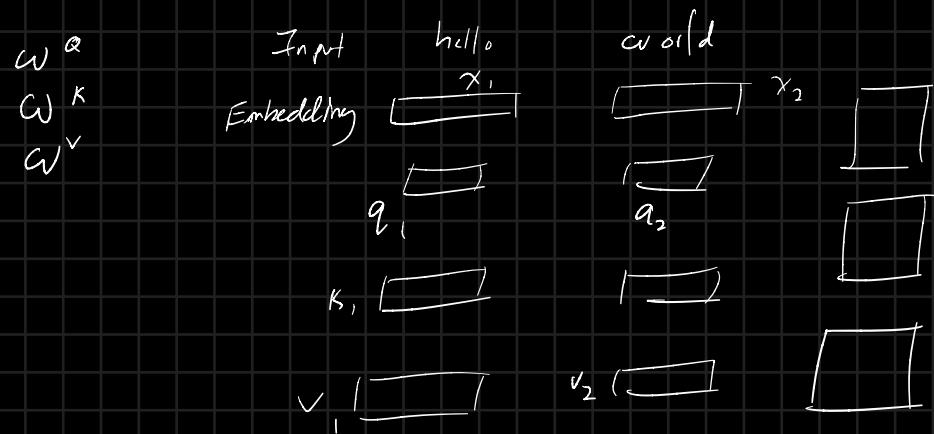
↳ When processing an embedding, it looks at other positions to lead to a better encoding of the word.

As vectors first:

Created by:

Multiply embedding by
3 trained matrices

- ① Create 3 vectors
- ② Query Vector
- ③ Key Vector
- ④ Value Vector



Step 2: Calculate a score

Score all other inputs against a word (how much activation contributed by other parts of the input sentence i.e. attention)

Dot product of query vector with the key vector of the respective word to score.

e.g. q_1 dot k_1 ,
 q_1 dot k_2 ...

Divide Scores by $\sqrt{8}$ (Sqrt of dim of Key vectors)

\Rightarrow Softmax operation. $(\overline{T_{d_k}})$

Scores how much each word will be expressed at this position

Step 5. Multiply each vector value by the softmax score.

\Rightarrow Keep values of the words we want to focus on, and make irrelevant word activations go to zero.

Step 6. Sum up the weighted Value vectors

\Rightarrow Send to the neural networks

$\leftarrow \text{---} \text{---}$

As a matrix operation...

Calculate Q, K, V matrices

$$n \begin{bmatrix} s_1 \\ x \end{bmatrix} \otimes \begin{bmatrix} w^k \\ w^v \end{bmatrix} = \begin{bmatrix} Q \\ K \\ V \end{bmatrix}$$

All rows correspond
to word in input
sentence (x)

$$\begin{bmatrix} x \end{bmatrix} \begin{bmatrix} w^k \end{bmatrix} = \begin{bmatrix} K \end{bmatrix}$$

$$\begin{bmatrix} x \end{bmatrix} \begin{bmatrix} w^v \end{bmatrix} = \begin{bmatrix} V \end{bmatrix}$$

Step 2-6:

$$\text{Softmax} \left(\frac{\begin{pmatrix} Q \\ K^T \end{pmatrix}}{\sqrt{d_k}} \right) \odot [V] = [Z]$$

Multi-headed attention:

$$\underbrace{\begin{pmatrix} x \end{pmatrix}}_{\text{Attention head } \#0} \quad \underbrace{\begin{matrix} \text{Attention head } \#1 \\ \dots \\ \text{Attention head } \#7 \end{matrix}}_{\text{Attention head } \#1} \quad \dots \quad \underbrace{\begin{matrix} \text{Attention head } \#7 \end{matrix}}_{\text{Attention head } \#7}$$
$$\begin{matrix} (Q_0) & [w_0^Q] \\ (K_0) & [w_0^K] \\ (V_0) & [w_0^V] \end{matrix} \quad \begin{matrix} (Q_1) & [w_1^Q] \\ (K_1) & [w_1^K] \\ (V_1) & [w_1^V] \end{matrix} \quad \dots \quad \begin{matrix} (Q_7) & [w_7^Q] \\ (K_7) & [w_7^K] \\ (V_7) & [w_7^V] \end{matrix}$$

Do same thing 8 times

$$A \# 0 \quad A \# 1, \dots, A \# 7$$

$$[Z_0] \quad [Z_1] \quad [Z_8]$$

$(z_0) \quad (z_1), \dots (z_8)$

$\underbrace{\quad}_{(z_7)}$

(z_7)

concatenate

$$[z_7] \begin{bmatrix} w_o \end{bmatrix}$$

Summary

Input Encoder
hello world $[x]$

Encoding only for
first layer

otherwise, we take
in R .

Single attention layer:

$$\begin{bmatrix} w_o^Q & w_o^K & w_o^V \end{bmatrix} [z_1] \quad (z_1) \quad \left(\begin{array}{c} u_o \\ v_o \end{array} \right) = R$$

$$\begin{bmatrix} w_1^Q & w_1^K & w_1^V \end{bmatrix} [z_2] \quad (z_2) \quad \left(\begin{array}{c} u_1 \\ v_1 \end{array} \right)$$

$$\begin{bmatrix} w_7^Q & w_7^K & w_7^V \end{bmatrix} [z_7] \quad (z_7) \quad \left(\begin{array}{c} u_7 \\ v_7 \end{array} \right)$$

Representing the order of the sequence using positional encoding:

To account for word order:

Transformer adds a vector to each input embedding

This added vector is trained.

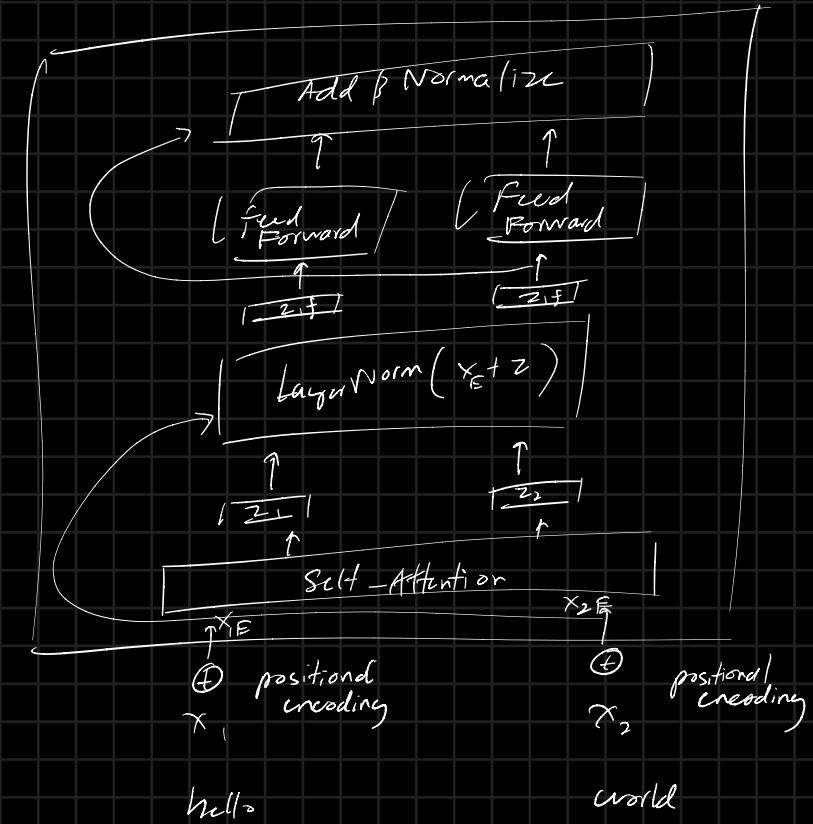
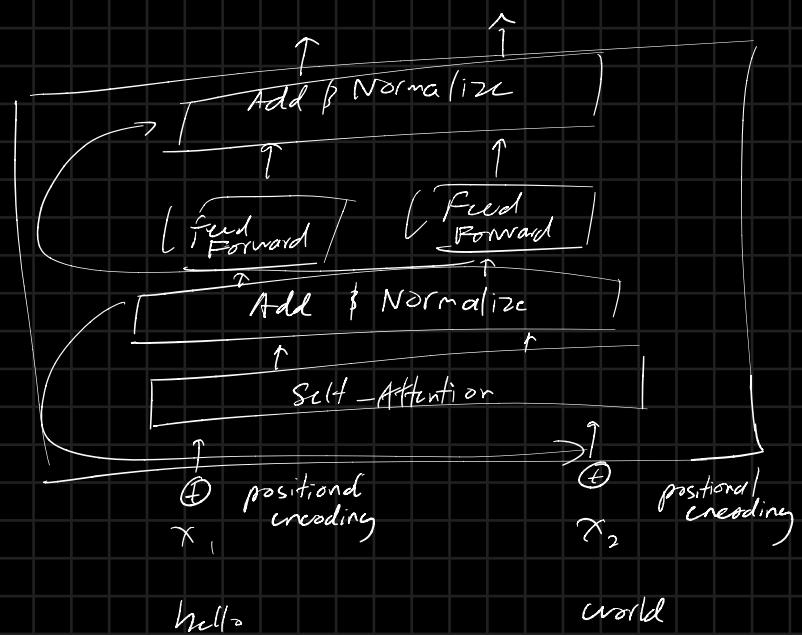
Represents "distance" between words

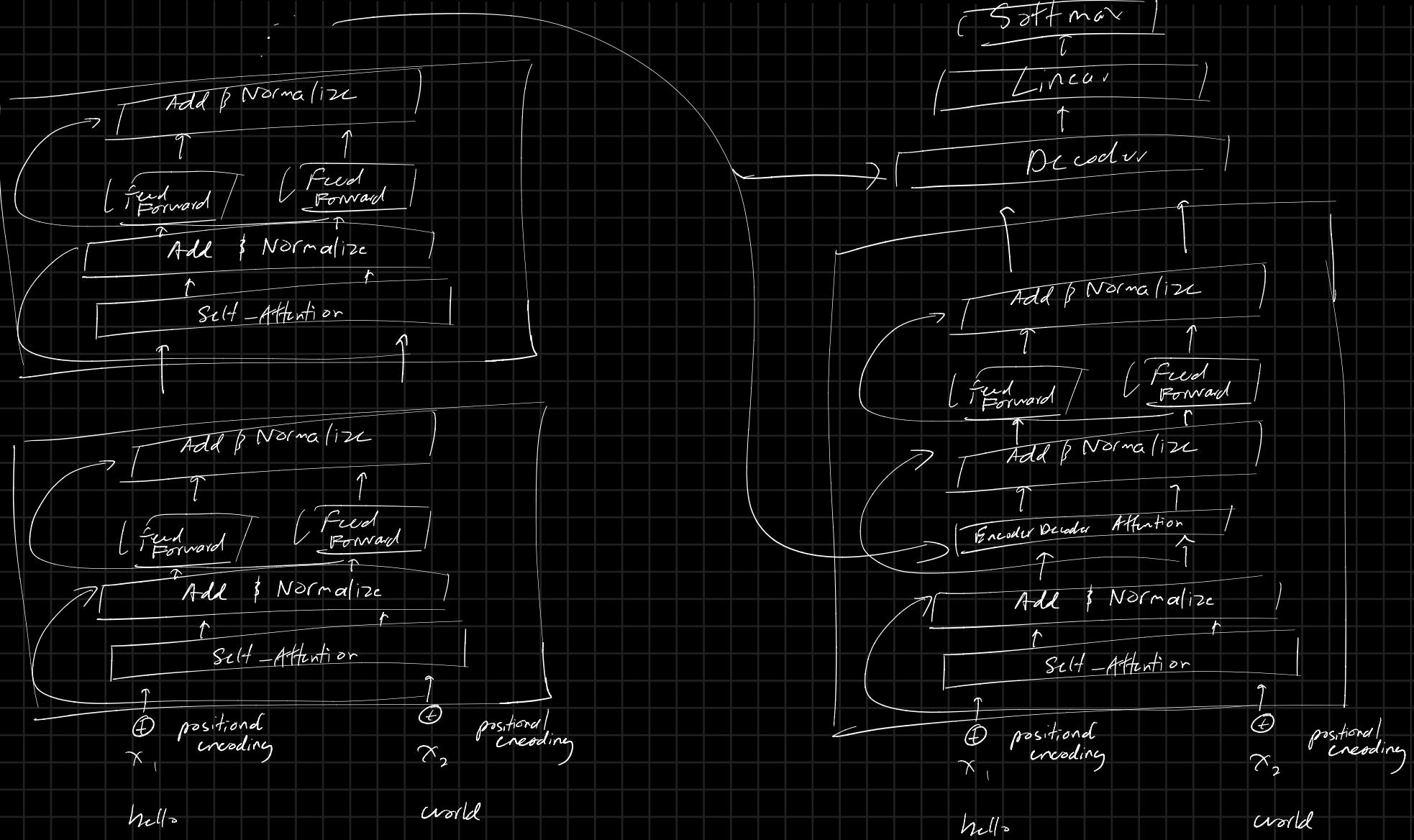
They are added element-wise

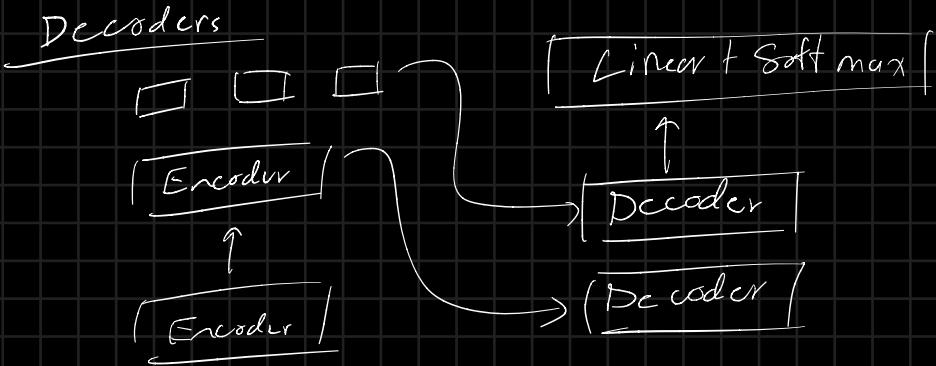
$$x_i + t_i = x_{E,i}$$

e.g. one single word

In the paper: residual connection + layer normalization







Repeat until special end token is outputted

Output of each step is feed into bottom decoder at next time step

In decoder, self-attention keys is only allowed to attend to earlier positions in the output seq
 Set future positions to -inf before softmax.

Encoder - Decoder Attention is exactly what we did, but creates queries from layer below and false keys / values matrices from output of encoder stacks.

Linear layer projects vector output from decoders to a logits vector (vocab size long)

Cell of highest probability is chosen

Logits \rightarrow Log probs

Define softmax

one hot encode vectors

maximize correct, minimize incorrect activations of

final log probabilities

Beam Search: Sample multiple times, take the one with least error

Visualizing neural machine translation models

Sequence to sequence models

Encoder + decoder (item by item)

Context ↗

Both are RNNs

context vector size (hyperparameter)

words are first embedded as vectors

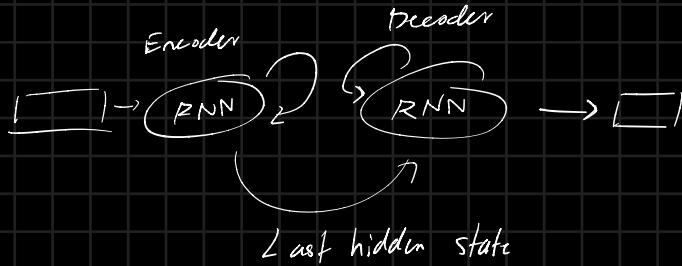
(200 - 300) size

$$h_1 \xrightarrow{\text{RNN}} h_2 \dots$$

↑
o.

x₁

Encoder



Decoder is an attention decoder RNN

All hidden states passed to decoder

Decoder attention also will:

Give each hidden state a score, softmax the scores
and then multiply hidden states by softmaxed
scores, to create the context vector

Done at each time step in decoder

Attention decoder RNN takes in <END> token embedding,
and initial decoder hidden state.
Both are used to produce new hidden state.

use encoder hidden states and new hidden state to
calculate a context vector for this time step.

Concat context and hidden state to a feed forward
NN, that indicates output word of time step.

Reproduce GPT 2:

124 mil parameter model

Mini Series

=> Biggest model => GPT 2

Down stream metrics

124M, 12 layers, 768 dim in Transformer

plotting wpc. weight

Every row is a position (0, 1024)
positional learns sinusoids

Response more or less (activation) based on
the position

Jagged (not fully trained)
Expected to see smooth curves

Positional embeddings are fixed to sinusoids + cosines at diff
frequencies

Decoder only transformer
Cross attention missing

GELU (always a change)

Dataset: tiny shakespeare

3 to 1 compression: 3 words, 1 token

wc input.txt

lines, words, bytes

view stacks up
sequences for batching

load 1 additional token

Reasonable starting point for the loss

Every vocab gets uniform prob

$$-\ln\left(\frac{1}{\text{vocab-size}}\right)$$

share same weights between embedding layer
and pre-softmax-linear layer

(weights are tied together)

want 2 matrices to behave similar

Intuition: Similar tokens are nearby in
the token embedding space (similar embeddings)
=> Get better performance this way

↳ 2 dependencies

Initialization:

stddev = 0.02

nvidia-smi

import code

code.interact(local=locals())

By default, all parameters are float 32

Don't need too much precision for training

FLOAT ↓ 7FLOPs ↑

7FLOPs

Sparsity

int8 used for inference, not training

Floating point operations (Trillions)

Memory bandwidth (less bits for datatype, more quick training)
↳ often the bottleneck

↓

Memory speed to access memory

Tensor Cores (instructions in A100 architecture)

(4,4) @ (4,4) + (4,4)

Varying precision FP32

Calculations get broken up like this (fastest way to multiply matrices)

Most work happens in the form of nn. Linear

bits

S 28 m23 FP32) 8x faster
S 28 m10 FP32 Can tolerate this empirically

Cost-reduced precision (approximate)

torch.cuda.synchronize()

waits for GPU to finish work

Metric: tokens per second

Drop to BFLOAT16 => less precision and more truncation

FP16 → reduced int range

Automatic mixed precision

torch.autocast()

Some converted tensor, some are left alone

LSTM cell

Layer Norm

No gradient Scalar

`torch.compile()`

compiler for neural networks

GCC for neural nets

Use always if not debugging

Reduces Python overhead to GPU

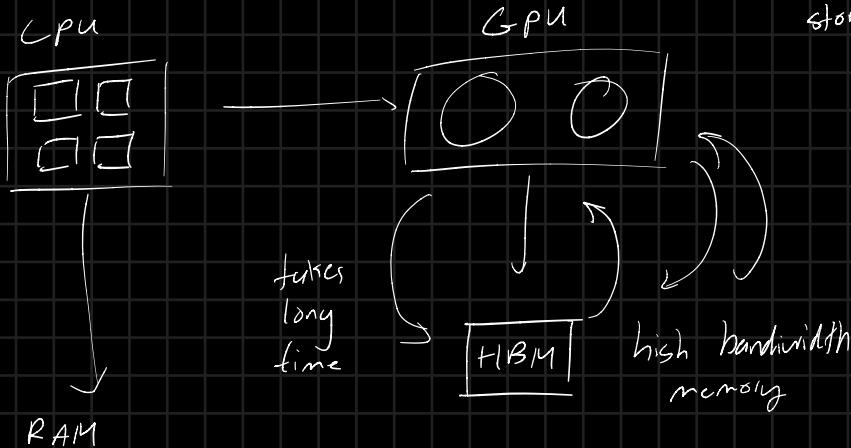
Compiles NN as an object



`torch.compile()` analyzes operations and optimizes it

Interpreter runs in = eager mode"

All tensors stored on GPU

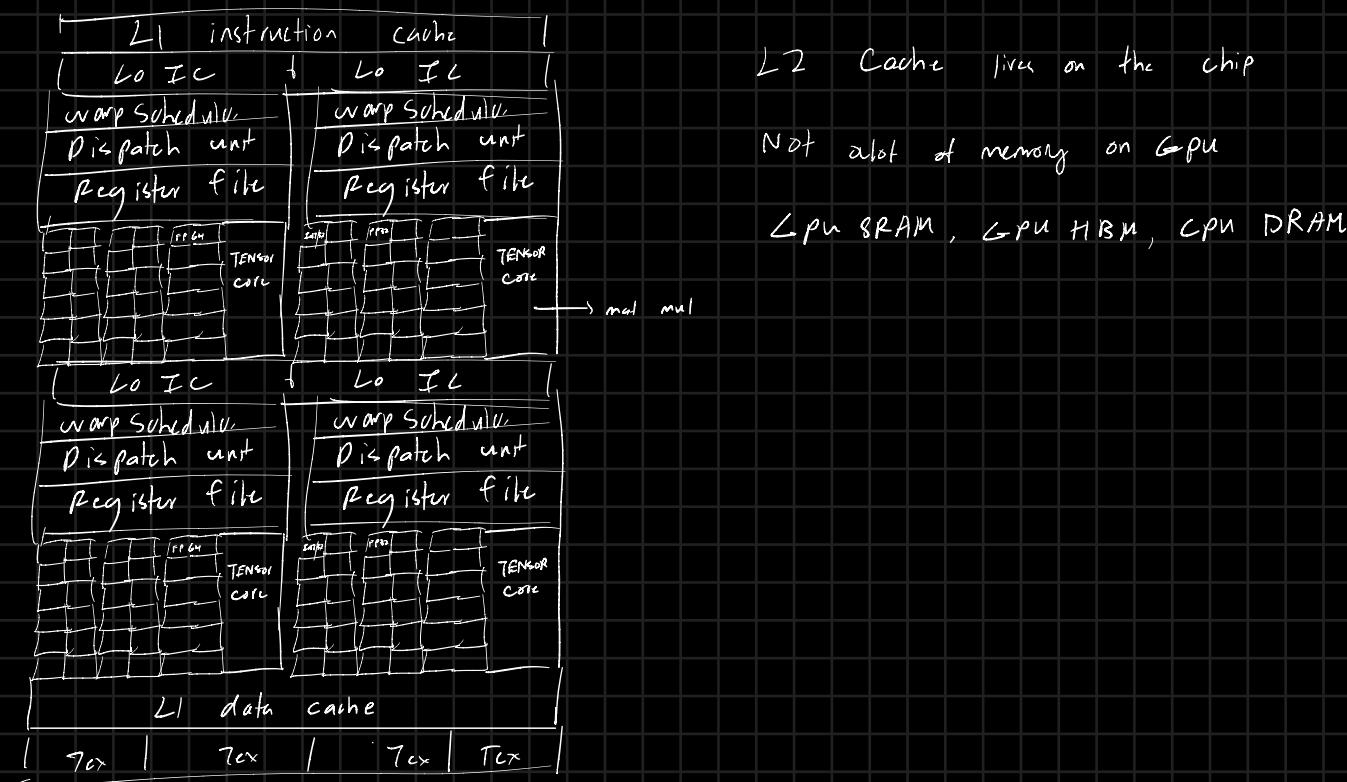


Kernel fusion !!! `torch.compile()` optimizes read/writes

Calculations occur in GPU
memory in HBM (2 chips)

Streaming multiprocessor (SM)

SM



L2 Cache lives on the chip

Not a lot of memory on GPU

L1 SRAM, GPU HBM, CPU DRAM

Floch attention: Kernel fusion

online normalizer for Softmax

mat mul, dropout, softmax, mask, mat mul

ugly numbers vs. nice numbers

Everything in CUDA written in pow of 2

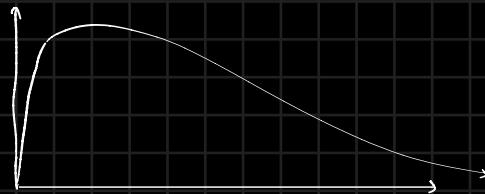
$$P_1 = 0.9$$

$$\beta_2 = 0.95$$

length of vector
↑

Clip gradients norm

Cosine decay learning rate (10%) with warmup



update batch size as well, increase over time

Learn which tokens actually show up first

Gradients are very similar

Gradients after will start to differentiate

Data is sampled without replacement

weight decay 0.1, to implement regularization

Fused Adam $\text{CW} \Rightarrow$ kernels are fused during calculations

Batch size of 0.5 million \Rightarrow gradient accumulation

Remember to do accumulate in the loss. `backward()` / normalize

Data Distributed Parallel

Launch x processes, assigned a GPU

Average all gradients

web text dataset

Scraped all webpages + text on Reddit

Common Crawl dataset

WebText 2

Books 1

Books 2

wikipedia

Good datasets

Red pajama

Slim pajama

Fine web

Sample -10 BT

Download dataset

Tokenize dataset using GPT tokens

Start tokens with <|cot> end at text token

Tokenize in uint16

Tokenize as np.array()

Shards → numpy files

Train + val Shards

100 shards

Advance shard and loop, get tokens, adjust position

(calculate total) warmup tokens

$$\text{Steps} \times \frac{\text{sec}}{\text{step}} \times \frac{1 \text{ min}}{60 \text{ sec}} \times \frac{1 \text{ hour}}{60 \text{ min}}$$

torch.Generator(device=device)

()

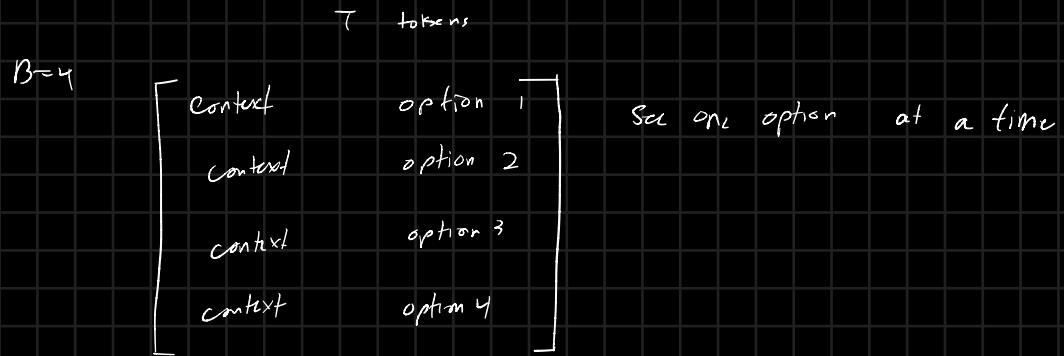
✓
Sampling from diff seed, for each GPU
all diff from original

falla swap

Multiple choice constructs, one is natural, rest are not

Shared context

Token completion



Fine tune into chat format

\Rightarrow Swap out dataset into user assistant structure

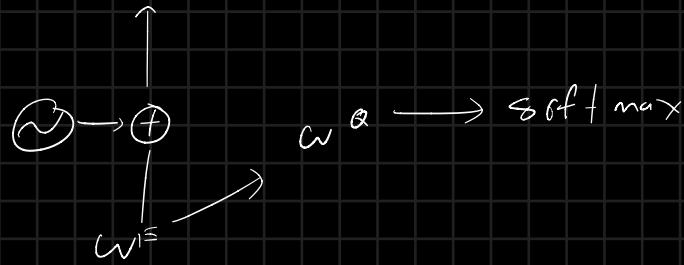
Fine tune on this

Fill in user, sample or assistant

One hot encodings vs. token embeddings

$S \times \text{dim}$

\Rightarrow positional embedding



parallelization

Stabilize gradients

Output of self-attention: context enriched word embedding

Multi-head Self-Attention

Layer Norm: subtracts mean, divides by std deviation, across all vectors

Batch Norm normalizes across all features independently

Layer Norm normalizes across each of the inputs independently across all batches

BN - CV tasks

LN - NLP tasks \rightarrow varying sentence length