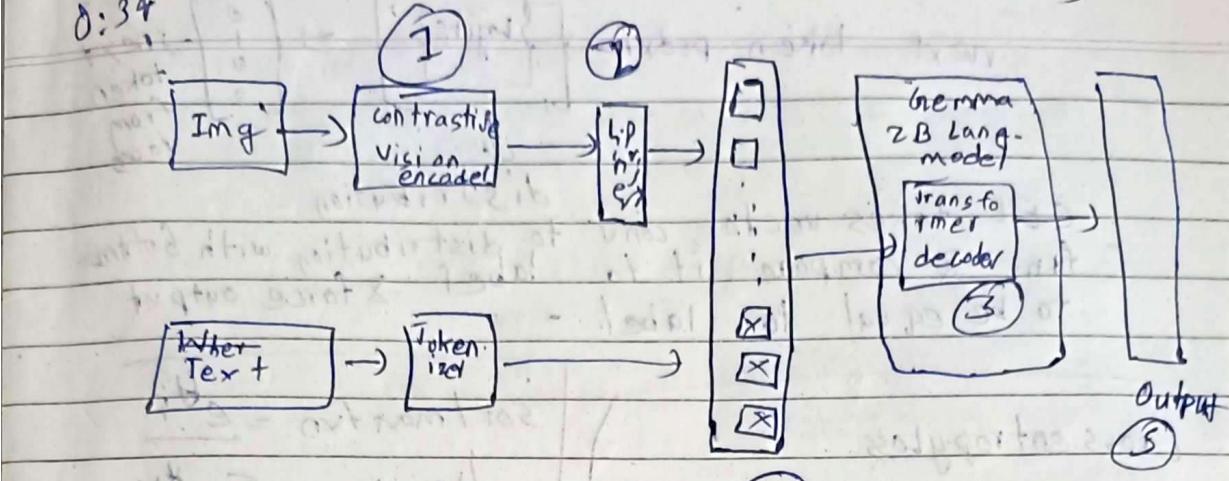


VLMs from scratch in Pytorch (Paligemma)

0:34



VLMs can extract info from Img!

Topics

- Vision Transformer
- Contrastive learning (CLIP, Sig Lip)
- Language Model (Gemma)
- KV-Cache

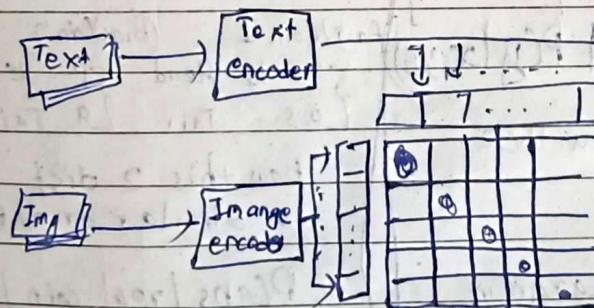
Rotary Positional encoding

- Normalization (Batch, Layer, & RMS)

Part 1

Q) What is contrastive learning

1) Contrastive pre-training

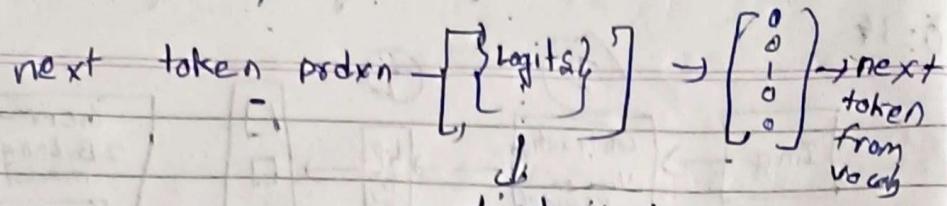


{ we want when
similar image & text
things are dot product
it shall give high value.
if img & text are
diff the value is low

We will first find everything in 3

high &
each box & then we will implement
a loss fxn that forces the
TL-BD diagonal to be high & others to
be low. [Cross entropy loss]

Why cross entropy loss (CEL),



CEL takes vector conv
fan & compare it to
to be equal to label distribution
to distribution with softmax
label & force output

Cross entropy loss

$$H(P^* | P) = -\sum_i P^*(i) \log P(i)$$

Probabilistic theory argues
"most natural output will be
entirely prob dist."

$$\text{softmax fn} = \frac{e^{u_i}}{\sum e^{u_k}}$$

→ Real values $\sum e^{u_k}$
probab.

→ Only used in output
Higher probab
as actual output.

Param (θ)

$y = \text{class}$

x_i Predicted: True

$$P(y|x_i; \theta), P^*(y|x_i)$$

An intuitive loss fn for APP should be:

we should focus on
minimizing KL divergence btwn
model dist & true dist

$$D_{KL}(P^*(y|x) || P(y|x_i; \theta))$$

which is further
processed

* KL divergence
(fallback library)

dist measure btwn
distributions

$$\text{discrete: } D_{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

$$\text{continuous: } D_{KL}(P||Q) = \int P(x) \log \frac{P(x)}{Q(x)} dx$$

[Trying to measure
dist btwn prob dist]

Fair Coin I (Bias Coin 2)
Goes to Head {P Head
L.S Tail } Q Tail

How this 2 dist
is similar or diff

$P(\text{obs} | \text{real coin})$
 $P(\text{obs} | \text{coin 2})$

A natural measurement
of dist btwn prob
dist. motivated by looking
at how likely 2nd dist
will be able to generate
samples from 1st distri.
w.r.t. 2nd part

$$\underset{\theta}{\operatorname{argmin}} D_{KL}(P^* || P) = \underset{\theta}{\operatorname{argmin}} - \sum_y P^*(y|x_i) \log \frac{P^*(y|x_i)}{P(y|x_i; \theta)}$$

$$\underset{\theta}{\operatorname{argmin}} D_{KL}(P^* || P) = \underset{\theta}{\operatorname{argmin}} H(P^*, P)$$

For this to work our dist
should be valid dist i.e.
such as sum is 1 & 0 only
& this is maintained by
softmax fn & get 2nd part

Problem with clip [Computationally expensive]
we are using clip

softmax has a problem: it is numerically unstable
as exp fxn can grow fast & may not fit
in 32 bit FP number

Soln: do not make exp grow to ∞

$$s_i = \frac{c \cdot e^{a_i}}{c \cdot \sum_{k=1}^N e^{a_k}} = \frac{e^{\log(c)} e^{a_i}}{e^{\log(c)} \cdot \sum_{k=1}^N e^{a_k}} = \frac{e^{a_i + \log(c)}}{\sum_{k=1}^N e^{a_k + \log(c)}}$$

Normally choose $c = \max(a)$

This will make the arg of exp towards negative
& the exp itself towards zero.

Too many component & not parallelizing efficiently
too

assymetry in that table we need
every to calc softmax for each row & every column

Siglip ~~clip~~ Improvement on Clip

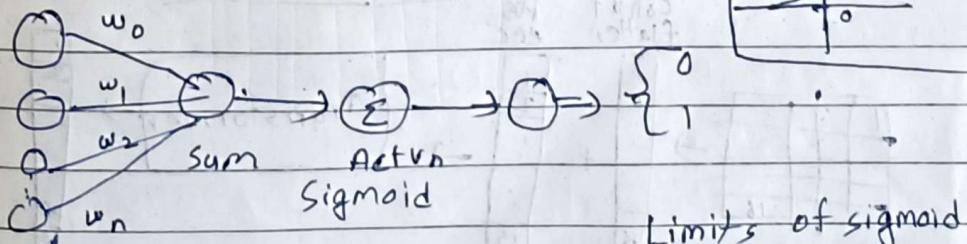
To use sigmoid loss instead of cross entropy loss.

Sigmoid fn

plays imp role in logistic regression

Logistic reg: technique to predict outcome of binary classif. problem.

$$f(x) = \frac{1}{1+e^{-x}}$$



Limits of sigmoid

Making predictions

$[0, 1]$

every predm $\geq 0.5 \approx 1$ predm
 $\leq 0.5 \approx 0$

Each box in table is treated independently as binary classif. problem. But in CLIP we opted towards depending on boxes itself.

And this is done via sigmoid fn.

So, by this way the principal diagonal becomes 1 and other becomes 0

coz of sigmoid & helps in parallelizing & no dependency.

Part 2

Why contrastive vision encoder

why not normal encoder for capture meaning of any img?

Because we need to integrate it with text & this app. we discussed just does that.

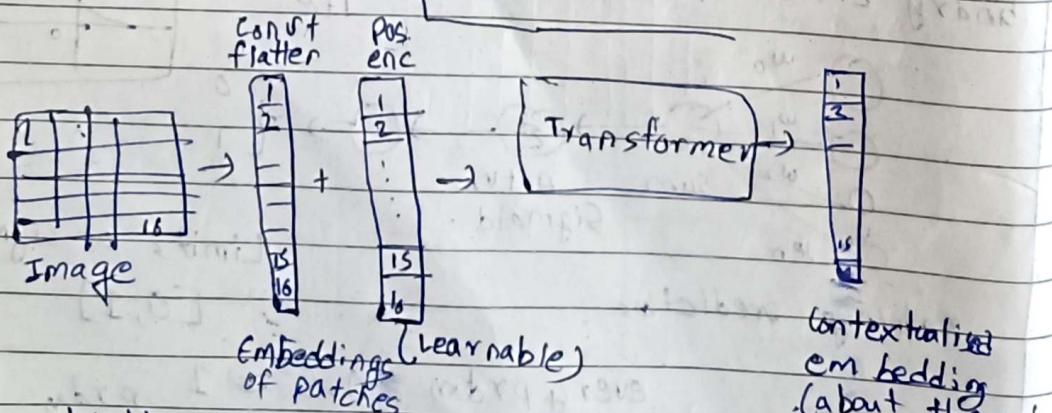
Cheaper to train bcz each img has its description in web.

Noisy dataset but it is trained on trillions so, eventually it becomes good.

Part 3 Vision Transformer

Explained in paper

"An Image is worth 16x16 words"



With LM we do causal mask + all others but with ~~EM~~ ViT we don't care about auto regressive nature like I output is depending on its previ part only.

Coding Siglip & input is image [patch embedding] flattened to things but outputs are also 16 things not 1 single thing

- make config & input \rightarrow hidden_size = embedding_size

- structure [Siglip Vision Model]

-- init -- config, call transformers?

forward (self, pixel_values) \rightarrow tuple:

[Batch_size, channels, Height, Width] \rightarrow

[Batch_size, num_patches, embed_dim]

- Siglip

Transformer

config

-- init --

config, embed_dim, embeddings,

encoder, post_layernorm,

forward (self, pixel_values, tensor) \rightarrow torch.Tensor

hidden_states = self.embedding

encoder \rightarrow last_hidden_state, (last hidden-state)

which is basically return last_hidden-state

transformer's encoder

in vit_norm is done before FFN &

MHA

- Siglip Vision Embedding

init ~~init~~ -- init --

config, siglip

config, embeddings, encoder,

post_layernorm

conv2D

num_patch

num_positions

pos_embeddings

forward forward (pixel_values)

patch_embed

embedding = patch_embed

flatten

embeddings, transpose(1, 2)

emb = embedding + post_embedding

Conv2D

inputdim = (H, W, C)

K = Filter

F = kernel size

S = stride (jumping)

P = padding (padding)

filter is for each

channel diff so, if

if K=2, then

it has basically 6 diff

filter.

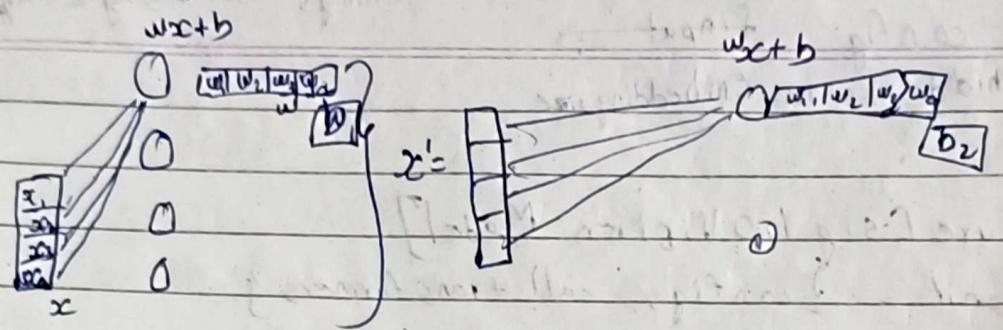
element wise mult.

can be implemented

[Unfold & Matrix mult

Resize]

Normalization



Problem: covariate shift

When input vector changes from 1 batch to another that changes in magn., the output also changes in magn., & it simultaneously affect loss, gradient & stuff, back propagation & affects model.

Network Learns Slowly

Soln:

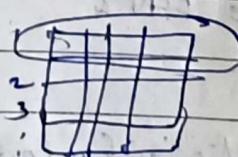
- ① Batch Normalization
 - dist. datapoints in mean ≈ 0 & $\sigma \approx 1$
 - we do it along l_{in}
 - problem with BN is that each statistic depends on what other items are in batch, to get good result we must use big batches

② Layer Norm

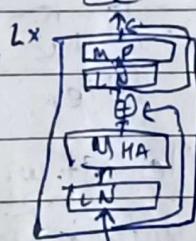
we do it along row

same mean = μ , $\sigma^2 = 1$

& independent of each other data points



Contextualised embedding



Patch + Post-norm embedding.

Siglip Encoder
Decoding Siglip Encoder

Siglip Encoder Layer

init. ()

embed super

embed_dim

self_attn = siglip Attention ()

layer_norm1

m/p

layer_norm2

forward

residual = x

hidden = layer_norm2(x)

$h_s = \text{self_attn}(x)$

$h_s = \text{residual}, h_s$

residual = h_s

$h_s = \text{layer_norm2}(h_s)$

$h_s = \text{m/p}(h_s) \rightarrow$

$h_s = \text{residual} + h_s$

why m/p?

non-linearity
& understand
complexity

Siglip MLP

init.

super

config, selfcl, fc1, fc2

forward

$h_s = \text{fc1}$

$h_s = \text{nn.gelu ("tanh")}$

$h_s = \text{fc2}$

Understanding MHA

Way of contextualizing stuff

Goal: transform initial seq. of uncontextualized embeddings to a seq. of contextualised embedding

Vision Transformer

$$x = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{bmatrix} \Rightarrow x = \begin{bmatrix} p_1, p_2, p_3, p_4 \\ p_1, p_2, p_3, p_4 \\ p_1, p_2, p_3, p_4 \\ p_1, p_2, p_3, p_4 \end{bmatrix}$$

$$x = \begin{bmatrix} C \\ I \\ cat \\ pizza \end{bmatrix} \Rightarrow x = \begin{bmatrix} C \\ I \\ I \\ I \end{bmatrix} \begin{matrix} I \\ eat \\ eat \\ pizza \end{matrix}$$

It is done in $1/C$ part it is
powerful coz of

Self Attention

$$\begin{aligned} &\text{-- init --} && \text{from } x \text{ to } Q, K, V \\ &\text{-- super --} && Q = X \times W_Q = (4, 1024) \times (1024, 128) \\ &\text{config, embed-dim} && K = X \times W_K = (4, 8, 128) \\ &\text{num-heads, head-dim} && V = X \times W_V = (4, 8, 128) \\ &\text{forward} && \end{aligned}$$

Coding

SI: from x to Q, K, V

$(4, 1024)$ \downarrow $(4, 8, 128)$
 \downarrow \downarrow \downarrow
 seq embed-dim seq head dim
 n_head

I have another note for other thing
 I will just binge watch this section

Each head will just watch a section of embedding like a head watches 1st 128 dim & so on...

because it is more eff

$$(4, 8, 128) \rightarrow (8, 4, 128)$$

In this way, ($4, 128$) to ($8, 128$) we can understand that each section can be computed independently as those 8 sections carries very imp meaning to see context in diff. way

- why?
- we want to parallelise comp
- each head should learn to relate tokens (or patches) diff

s-3: calc. attn for each head in, 11
($4, 128$)

$$Q_{\text{head}} = \begin{bmatrix} C & I \\ E & J \\ C & J \end{bmatrix}$$

$$K^T_{\text{head}} = \underbrace{\begin{bmatrix} I & J & J & I & J & I \end{bmatrix}}_{K}$$

$$\frac{Q \times K^T}{\sqrt{d_{\text{head}}} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} Q}$$

This repr relationship btwn tokens bigger value is good relationship.

Attn masking



because we tell we don't want current token to know about future token

$$\text{softmax} \left(\frac{Q \times K^T}{\sqrt{d_{\text{head}}}} \right)$$

$$\text{softmax} \left(\frac{Q \times K^T + \text{mask}}{\sqrt{d_{\text{head}}}} \right) = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Attention weight

Step 4: Multiply by V seq

each row rep weighted sum of

$$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot \end{bmatrix}_{(4, 128)} \times \begin{bmatrix} C & I \\ C & J \\ C & J \\ C & J \end{bmatrix}_{\text{eat}}^{I} = \begin{bmatrix} C & I \\ C & J \\ C & J \\ C & J \end{bmatrix}_{\text{eat}}^{I}$$

$$\begin{bmatrix} C & I \\ C & J \\ C & J \\ C & J \end{bmatrix}_{\text{pep}}^{J} = \begin{bmatrix} C & J \\ C & J \\ C & J \\ C & J \end{bmatrix}_{\text{pep}}^{J}$$

$$\begin{bmatrix} C & I \\ C & J \\ C & J \\ C & J \end{bmatrix}_{\text{pizza}}^{J} = \begin{bmatrix} C & J \\ C & J \\ C & J \\ C & J \end{bmatrix}_{\text{pizza}}^{J}$$

$$\begin{bmatrix} C & I \\ C & J \\ C & J \\ C & J \end{bmatrix}_{\text{I eat pep - pizza}}^{I} = \begin{bmatrix} C & I \\ C & J \\ C & J \\ C & J \end{bmatrix}_{(4, 128)}^{I}$$

Step 5: Transpose back

$(8, 4, 128)$ to $(4, 8, 128)$ to $(4, 1024)$

Step 6: concatenate all head

Given each head in computing contextualized
 $(4, 8, 128)$ to $(4, 1024)$

Step 7: Multiply by W_0

Parameter

$$\begin{matrix} (4, 1024) & \xrightarrow{W_0} & (4, 1024) \\ \left[\begin{matrix} C & P \\ C & P \\ C & P \\ C & P \end{matrix} \right] \times & \left[\begin{matrix} : & : \\ : & : \\ : & : \\ : & : \end{matrix} \right] = & \left[\begin{matrix} C & P \\ C & P \\ C & P \\ C & P \end{matrix} \right] \end{matrix}$$

Why multiply by W_0 ?

not much diff - To combine info from attn head

To restore correct dim.

We understand attention

that and attention rep. things in multiple dimension not only $2D$ for particular input token but in img we must have $(n \times 1)$ dim for any particular image which is the main input for contrastive learning.

① This is solved via by only taking 1st output part of contextualize embedding from transformer as a representative for whole img, force model to put all info in 1st contextualised embedding

②

② Just take avg of all output embedding to generate 1 single embedding

We are done with understanding how to build contrastive vision encoder

part - A

Processing Paligemma

We need to make tokeniser in this paper. Paligemma we are making the use of gemma tokeniser but our approach to understand image also, so some modifications are maintained to achieve what we want to achieve.

Paligemma further can be used for img segmentation or obj detection as well.

But we are doing paligemma for conditional modeling and other aspect is out of our scope.

Image normalisation mean & standard deviation

Same meaning as norm. that we do in neural network, so that we want always have same dist.

$$\text{Imagenet_standard_mean} = [0.5, 0.5, 0.5]$$

$$\text{Imagenet_standard_std} = [0.5, 0.5, 0.5]$$

Preprocess img:

S1: take img.

S2: resize

S3: change to np array

S4: rescale

S5 = normalize with std_mean & std_std

S6 = transpose to Cchannel, Height, width

Passed this S6 to make tensor ^{single} of tensor further to torch tensor.

And add ~~img~~ text aspect as well

Coding Gremma

Palindrome for Conditional Generation

Init -

config

vision-tower

multimodal-projector

vocab-size

language-model

Pad-token-id

tie-weight

weight tying

technique

of using param
of 1 layer to
another

forward

input-ids

pixel-values

attn-mask

kv-cache

input-embeds

selected-image-feature

image-feature \rightarrow Linear Projector

merge input-id-with-input-features

outputs = langmodel()

Gremma config

6

=
=

Gremma Config

group query-attn

when q \neq k, v has diff

no. of heads

so, we just make place to fit
every thing with final-embedding

then we make,
text-mask
image-mask
pad-mask

we make it applicable
to make it compatible with given
dim

text-mask-expanded
pad-mask-expanded
image-mask-expanded

Now add text embed

final_embedding [add text mask-exp]
if true with some
same for where fn.

final

Now add image things

with scatter. Can't use
masked_scatter [where]

final_embedding

Now for padding just
insert zeros.

K-V cache

Problem - k-v cache is solving:

Causal mask - helps in understanding causal mask of
present token & all the previous tokens.
if we don't use causal mask it understands
everything in a way [vision encoders kind of follows
it]

Some explanation of transformer on
how it works during inferencing

During inference we are generating
embeddings that we aren't using.

If it has 1000s token attn mech. will
gen 1000×1000 matrix from that only a small
portion will be used

And KV cache exactly helps with that helping us to exactly do that generate embeddings that matters

$$\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \times \begin{bmatrix} C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

(3,3)

(3, 128)

Attn weight

previously we used to multiply a1)

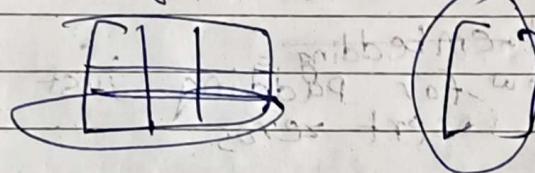
Attn wt with v to create

$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

Only this is imp. bcoz with

but this can be gen by just

using last row of Attn weight & v



Encoder

I now only

Q, K, V

$$K = \begin{bmatrix} I \\ I \\ I \end{bmatrix}$$

$$V = \begin{bmatrix} I \\ I \\ I \end{bmatrix}$$

[::]

$$\begin{bmatrix} Q \\ K \\ V \end{bmatrix} \rightarrow \begin{bmatrix} I \\ I \\ I \end{bmatrix}$$

[::] Logits

[::] Softmax

[::] Softmax

[::] Calc Atn

KV cache basically allows us to during inferencing to store token gen in to avoid generating all embeddings of all input seq but only last contextualised embedding

Prefilling : we can do it for multiple as user can give multiple tokens & generate all that can be // maintained

Only during prefilling phase we allow gen. of multiple output embeddings & discard that we don't need.

- because it will make it fast
- & output can be effectively utilised

2 step

i) Prefilling

ii) Token generation

In paligemma we don't mask out anything later like in GPT-2 transformer for img token we need to understand each thing with another

why prompt isn't causal masking as well to tokens of the prompt.

coz textual prompt itself is very short

& this prompt actually represents what we want our model to perform we want all the tokens to understand each other

Image, prefix, Suffix / Target

We only mask out what we expect Model to output

It's what paligemma chose.

The output gen aren't that long but they can be fine-tuned to be so as well.

During training we will have causal mask for suffix but during inference we don't have any causal mask working with KV cache.

Rotary positional encoding (RoPE)

Absolute posn / embedding
Relative posn / embedding

- Learned from data
- Posn vect for 1-512
- Max length is bounded
- Sinusoidal fn

It kinda solves the distance problem in words

like 1 & 3 will have same effectiveness as 501 & 503

This approach is slow

Empirically: Both approach have similar performance

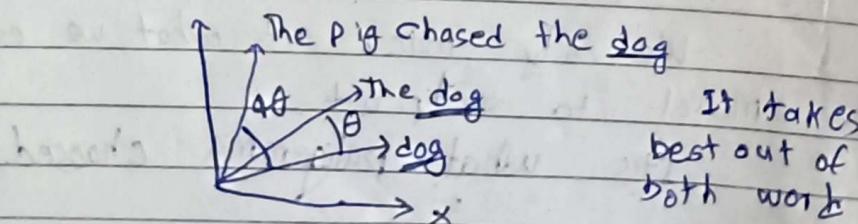
- Extra step in self-attn layer
→ Difficult for KV cache

Another problem is every single posn embedding is independent of each other

i.e. Pos. 1 & 2 is not diff from 1 & 500 which is intuitively wrong

RoPE:

Instead of adding posn in vector to encode posn of word in a sentence, they propose to apply rotation to vector.



Gremmafor Causal LM

Transformer Model + lang. modeling Head

GremmaModel

init
config

tie_weights

forward

(attn-mask
pos-ids
inputs-embeds
kv-cache)

output = attn-mask, pos-ids, inputs-embeds,
kv-cache

hidden = outputs

logits = ~~set~~ lm-heads (hidden)

logits = logits.float()

GremmaModel

init

config

padding_idx

Vocab-size

embed-tokens

layers = ModuleList [Gremma Decoder Layer]

norm = Gremma RMSNorm

get-input-embeds

forward (attn mask, pos-id, input-embeds, kv-cache)

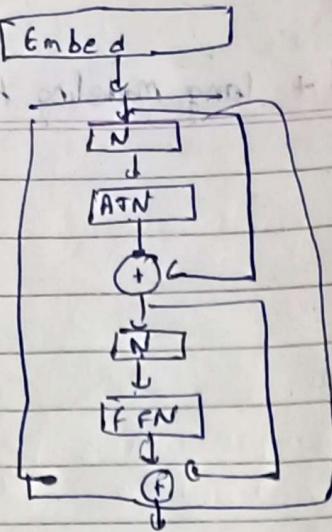
hid-states = input-embeds

normaliser

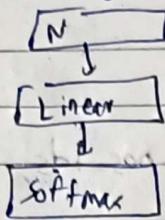
hidden-state_j = hid-states * normaliser

for loop

all transformer blocks



$$N = \text{RMS norm}$$



RMS norm

Claim: success of LayerNorm is art. bcz of re-centring invariance but coz of re-scaling invariance

what the μ mean is the mean need not be around 1. It could be around any no. given that it is normally dist. around that number

we only need to compute RMS instead of σ , μ like in LN

$$\bar{a}_i = \frac{a_i}{\sqrt{\sum a_i^2}}, \text{RMS}(a) = \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2}$$

faster than LN

```
class GammaRMSNorm
    init
        eps
        weight
        norm
```

forward

EPS is req. if no is close to 0
division by 0
EPS is added

class Gremma DecoderLayer
init
Gremma Attn
Gremma MLP
RMS1
RMS2

forward

like in fig

class GremmaMLP {get+}

init

config

hidden_size

intermediate_size

gate_proj

up_proj

down_proj

forward

gelu actvn fxn

Gelu

motivation

- combines both

relu & dropout

- deterministic

non-linearity

weighted sum - x

$$g(x) = \text{GELU}(x) = x \cdot \text{erf}(\frac{x}{\sqrt{2}})$$

m = cumulative dist.

fxn (Normal dist +

we have gate-proj now here which
is basically upping & downing no. of neurons

class GremmaAttention {config, layer_idx} Gremma is a decoder only
init

key_value_heads

key_value_group

rope_theta

is_causal

q_proj (hidden_size, num_heads * head_dim, attention_bias)

k_proj

v_proj

o_proj

It is diff from
normal attn.

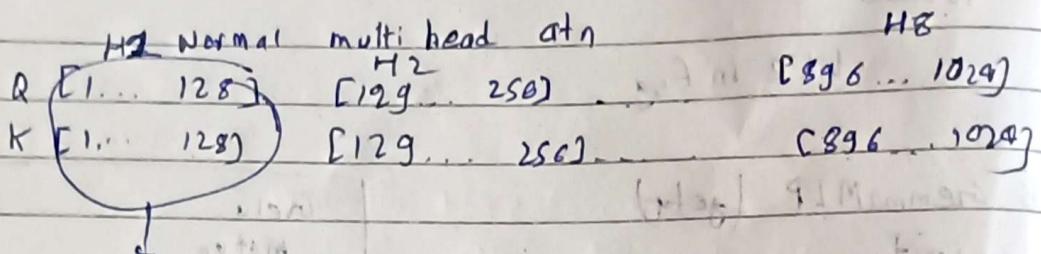
Wq size [1024, 1024] # Wv size = [1024, 1024]

Wk size = [1024, 128]

Wo size [1024, 128]

We have less head in key & values then query which is Grouped query, attn concept.

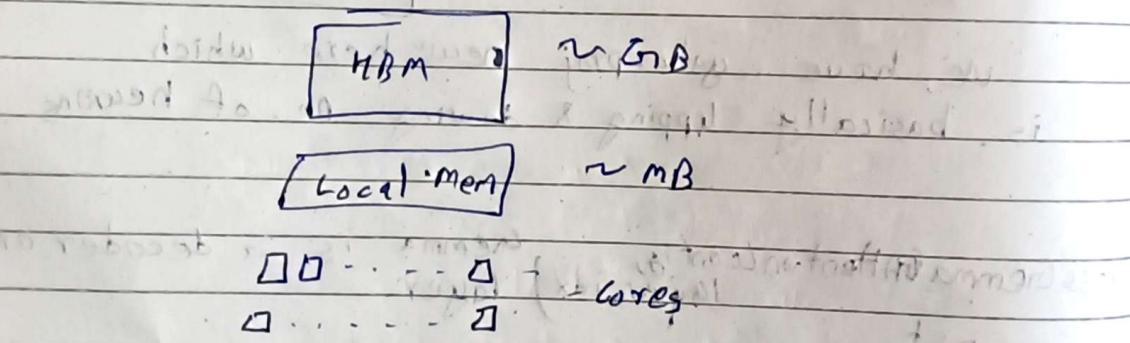
Grouped query attr



dot product. similarly all other in #

Problem with MHA:

claims: problem isn't with no. of computations that we are doing, But the no. of data transfer happening in GPU



matrix opn can be done in 11

Mem copying is very slow.

Similarly Concept in Flash after we redo computation instead of saving coz saving means copying which costs time

So, Our approach is to reduce data transfer
1) Use less head for keys

A group of ~~query~~ query will use 1 key

In a way it reduces some quality but not much which we can negotiate based on our usecase

Also reduces size of KV cache

KV cache is also have been bottleneck and it also helps in that

rotary-embeds = Embeddings /

head_dim, max_pos_emb,

base = rope_theta

forward

seen as usual + rotary embeds

KV-cache thing

class KVcache{}

init

key_cache

value_cache

num_items

if == 0

return 0

else

return key_cache[0].shape[2]

update

if never done

else(done)

update

repeat_kv [basically repeats for the keys & values of the query that are missing for heads of the query]

Custom cuda kernel we are doing this coz we don't have

attn weight
softmax
dropout
attn output
return

New positional encoding used in decoder layer.

in vanilla transformer we don't use absolute posn embedding.

we use absolute posn embedding i.e. sinusoidal things

But now adapts we use RoPE (Relative posn embedding)

We don't add directly into tokens embedding of each token

but they modify the mech itself so that attn mechanism takes into consideration the posn of tokens to relate them differently based on the posn

paper explores qk

can we find encoding of embedding vectors such that when we do

dot product that encodes posn / embed or, $q \cdot k$'s dot product only depends on embedding of 1st token, 2nd token & relative dist btwn them

Huggingface implementation of RoPE is explained which is a bit diff from the approach we have in paper but the result is same