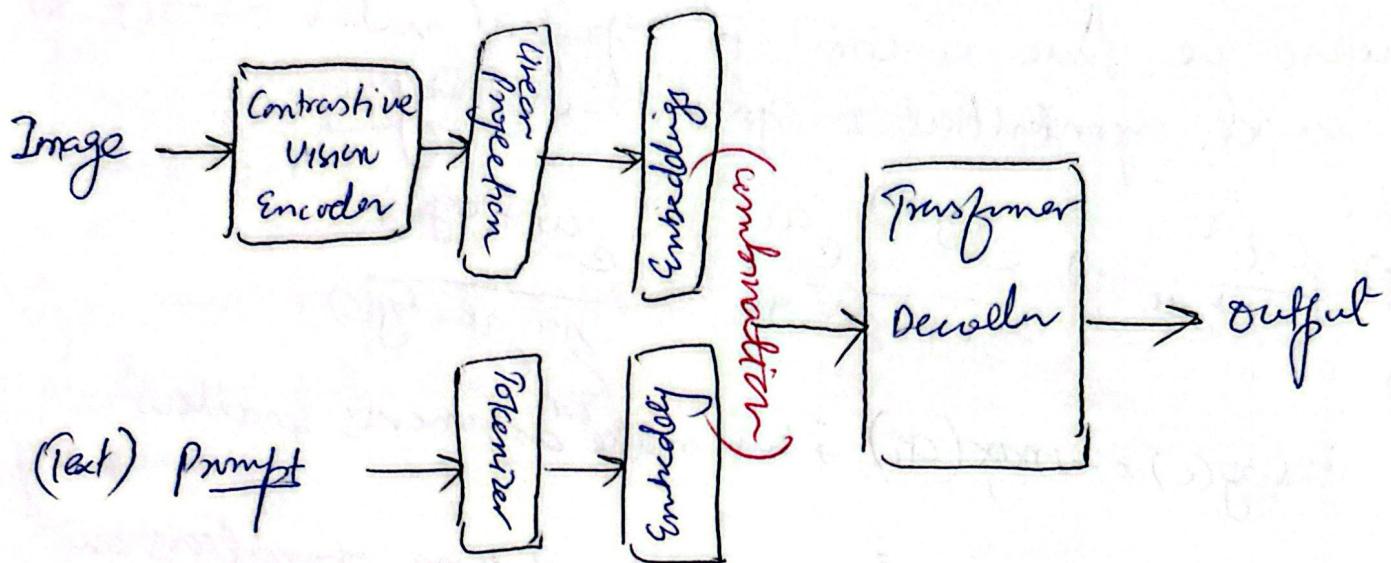


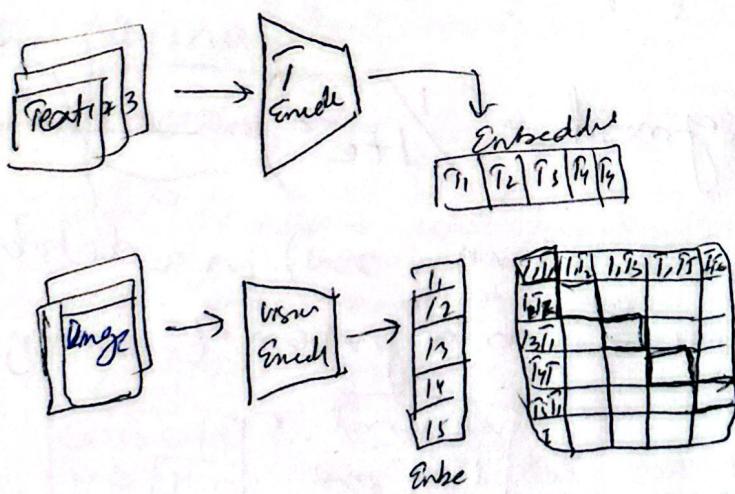
Multimodal Vision Language Model (PaLigenna)

→ Vision NLP → can extract info from image



→ Contrastive learning ($CIP \rightarrow \text{Sig}^2(p)$)

1) Contrastive Pre-training.



we multiply each T embedding value with each I embedding to get a Matrix. and we want to train these encoders to produce a high dot product value for image and their corresponding text & low vice versa.

Problem with CIP's

$$\text{softmax}_2 S_2 = \frac{e^{a_i}}{\sum_{j=1}^B e^{a_j}}$$

we find a loss function to make correctly pair high by CrossEntropy loss. (we compute loss for both and then average.)

Numerically Stable softmax

The exponentials in softmax cause issue because they grow very fast and if big number comes it's a problem because we have a limit to represent until 32-bit. So we avoid exponential to grow to infinity.

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

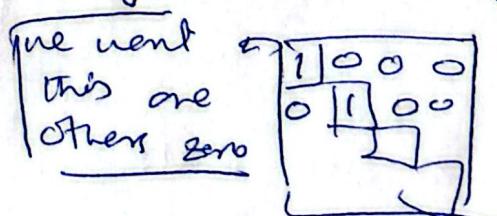
$\therefore \log(c) = -\max_i(a_i)$ → this make arguments smaller.

So, in step we have to calculate all those operations on each row and column, which make it computationally heavy due to a lot of operations.

So, In SIGLIP → It help parallelize & high batch size.

Replace softmax with sigmoid = $(\frac{1}{1+e^{-x}})$ or $\sigma(x)$

So, instead of treating (each column or row) as a distribution we independently treat them as a problem of binary classification ('1 vs all')



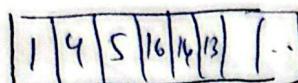
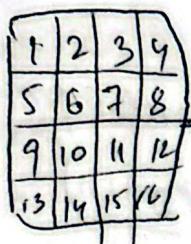
why contrastive encoder.
we want the embeddings that representation of image and have properties that can be later used with text embedding

Hence, contrastive + caption information of image.

Vision-Transformer (seq-seq Model)

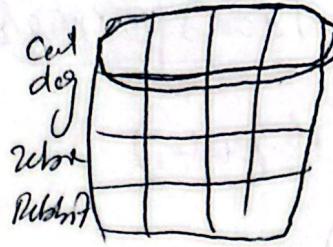
(divide to 16 patches)

P.E (learned)



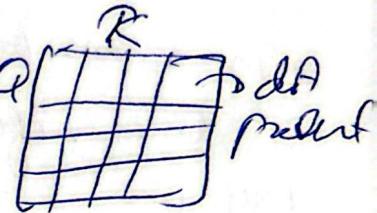
layer - Normalisation. (Item-dimension)

- In this each item is treated separately → Global training
- It's only depends on cat dimension



$Q, K, V \Rightarrow (1024, 1024)$ or $(4, 8, 1024) \xrightarrow{\text{seq. length reads.}} (8, 4, 1024)$ for parallel process
so, each head will compute independently. dimensions

$$Q \times K^T = \underbrace{(4, 128)}_{\sqrt{d_{\text{head}}}} \times (128 \times 4) \rightarrow (4, 4)$$



$$\times V = (4 \times 128) \rightarrow (4 \times 128) \xrightarrow{\text{weighted sum.}}$$

- concatenate (heads) W^0 . $\Rightarrow \text{concat}(128, 128 \dots) \xrightarrow{4} 1024$
- $W^0 = (1024 \times 1024)$, we do to mix result of different heads.

Linear

In layer norm, they treat each independent feature from batch, making it $\mu = 0, \sigma = 1$, so it is coming out from distribution (gaussian)

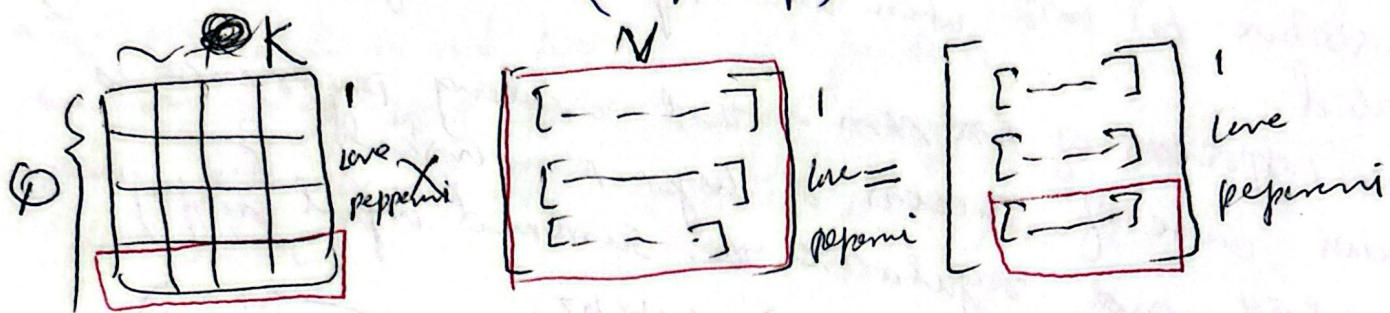
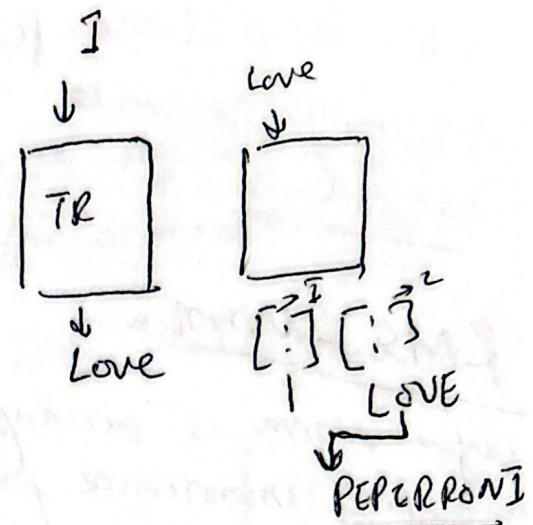
[Some formulae but we calculate]
statistics by row instead of cols]

- * Gamma and Beta are learnable parameters

KV-Cache

problem - generating alert of embedding, but using only one to project the next token (that have info about all previous)

attention scores matrix = $(\text{seq} \times \text{seq})$



\rightarrow pepperoni is result of last $xn \times nk$ all N (columns)
so, we do.

instead of passing all previous tokens, we just passed last single token (Query) and we append it to previously created KV-cache buffer which we used for calculation of self-attention, and we'll get output embedding (2-tokens) contextualized

prefix

\rightarrow adding many embeddings of token in KV-cache in single forward pass

Causal Mask

\rightarrow we can't mask image tokens, and are also not masking prompt mask if short q also because it define the task to do.

but, it's a Chriso PaliGamma, as we are not generic
we only generate Target/Output which we made
causal.

RMS - Normal

layer Norm is successful, because its centering and
rescaling invariance property. due to which model is more
sensitive to noise when both input and weights are randomly
scaled.

The hypothesis of RMS_norm is that rescaling invariance is
main cause of success of layer Norm instead of
centering, and regularize the summed input simply
according to (root-mean-square) statistic.

$$a_i = \frac{a_i}{\text{RMS}(a)} g_i, \text{ where } \text{RMS}(a) = \sqrt{\frac{1}{n} \left(\sum_{i=1}^n a_i^2 \right)}$$

when means of RMS sum = 0 (RMS = layer Norm), although
RMS doesn't \propto center.

→ The goal is values around their mean, which can be
any value {not 0} as in layer Norm?.

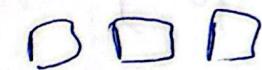
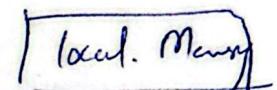
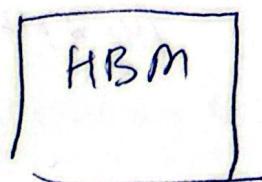
Multi-Query Attention

The difference is we have less heads in key and value which results in smaller projection for embedding of each token.

$$\{1 \dots 1024\} / 8 \text{ Heads} \rightarrow [1 \dots 128] [129 \dots 256] \dots [1024]$$

Multi-query paper

→ The bottleneck is not how many dot products are we doing but the memory transfer from high-bandwidth to local memory (time taken)

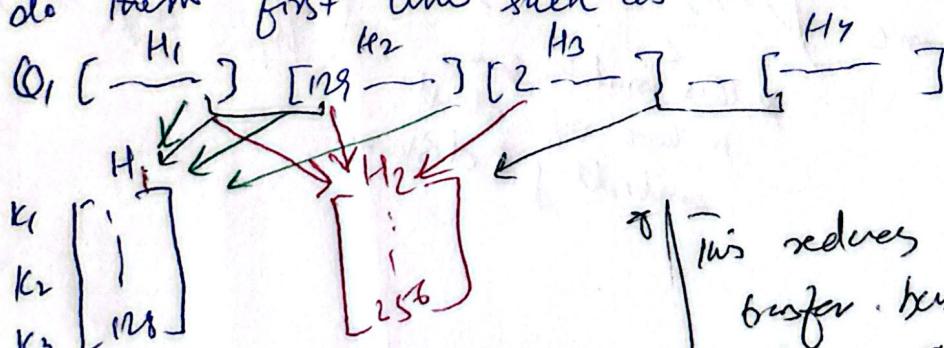


→ So, compute is less important than data transfer. Flash attention exploit this difference b/w data transfer and compute. It's twice faster to reduce computation than copying different stuff from GPU. This is also why we do gradient checkpointing. During backpropagation we do redo computations instead of copying them from HBM \rightarrow LM.

To Reduce data transfer while we do multi head attention

→ One way is to use less heads in keys.

→ so instead of copying each head of key for queries, we do them first time such as



* This reduces data transfer based on some cost on quality of model

Multi-query attention = $\text{f}_1 \text{v}$ (heads = 1)

Grouped Multi-Query attention = kv (head = 2 or more)

also reduces the size of kv-cache

$$(4, 8, 128) \begin{bmatrix} \text{head } 1 \\ \text{head } 2 \\ \vdots \\ \text{head } n \end{bmatrix} \xrightarrow{\text{grouped}} \begin{bmatrix} \text{kv} \\ \text{kv} \\ \vdots \\ \text{kv} \end{bmatrix} \quad \text{for parallel}$$

query needs to attend all past key and values
 $\Rightarrow (\text{expand} = \text{repeat})$

Rotary Positional Encoding

They are not directly added to input, but they modifies attention mechanism in a way that it take in consideration the position of token.

$$f_q(x_m, m), f_k(x_n, n) = g(x_m, x_n, m-n)$$

The dot product depends on embedding + relative distance b/w them

2D-Cart

$$\text{RotMat} = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}, \text{so, } f(q, k)(x_m, m) = R \begin{pmatrix} w_{q1} & w_{q2} \\ w_{k1} & w_{k2} \end{pmatrix} \begin{pmatrix} x_m^1 \\ x_m^2 \end{pmatrix}$$

p this make sure to have info about embedding + distance

General form

$$R_{\text{Rot}}^d =$$

$$\begin{pmatrix} \cos \theta & -\sin \theta & 0 & 0 \\ \sin \theta & \cos \theta & 0 & 0 \\ 0 & 0 & \cos \theta/2 & -\sin \theta/2 \\ 0 & 0 & \sin \theta/2 & \cos \theta/2 \end{pmatrix}$$

So, better way: as we already know where are zeros (0)

→ It Sparse,
mean most of
elements is 0,
and we will be
doing lot of
computation.

Computationally efficient realization of Rotory Matrix Multipl

using element-wise multiplication.

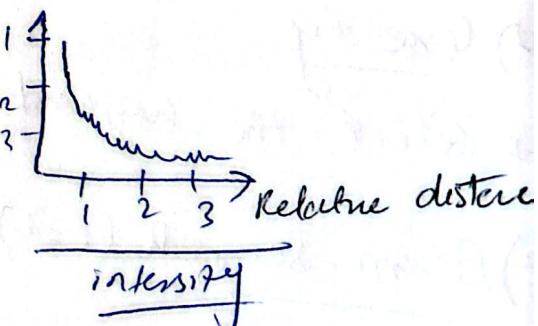
$$R_{\text{Rot}}^d = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \cos \theta_1 & \cos \theta_2 & \cdots & \cos \theta_{d/2} \\ \sin \theta_1 & \sin \theta_2 & \cdots & \sin \theta_{d/2} \end{pmatrix} + \begin{pmatrix} -x_2 \\ x_1 \\ -x_4 \\ x_3 \\ \vdots \\ -x_d \\ x_{d-1} \end{pmatrix} \otimes \begin{pmatrix} \sin \theta_1 & \sin \theta_2 & \cdots & \sin \theta_{d/2} \\ \cos \theta_1 & \cos \theta_2 & \cdots & \cos \theta_{d/2} \end{pmatrix}$$

positions to encode

embeddings dimensions rotated dimension embedding

$\therefore \theta_i = 10000$
 $i = [1, 2, \dots, d/2]$

Relative UB



decaying effect

similar = high

less = decay