

Deepseek from scratch

zero-shot

→ No examples are given.

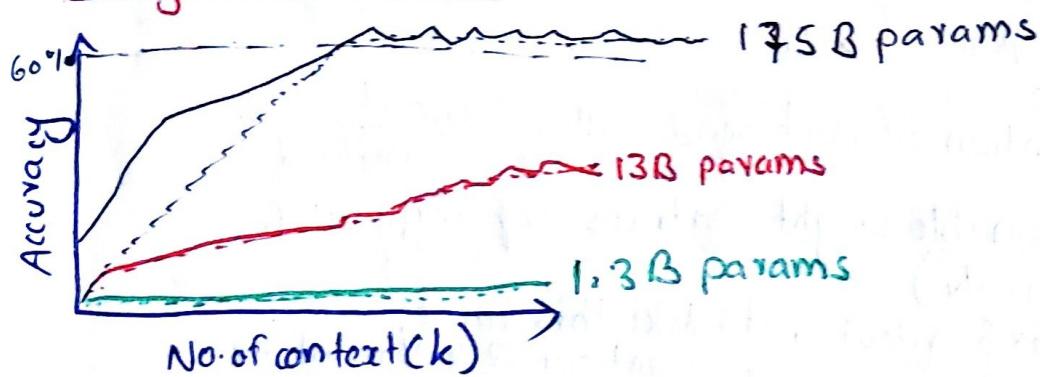
few-shot

→ few numbers of example is shown to model to predict the answer.

oneshot

→ One example is given.

* As we increase the no. of parameter, we human notice some emergent behaviour.



Our learning journey.

Phase 1

Innovative Architecture

- Multi-head (MLA)
- Latent Attention
- Mixture of experts (MoE)
- Multi-token prediction (MTP)
- Quantization
- Rotary Positional Encodings (ROPE)

Phase 2

Training Methodology

- Rise of Reinforcement learning (RL)
- Not only relying in human labeled data, use large scale RL to teach complex reasoning to the model.
- Rule based reward system.

Phase 3

CPU optimization tricks

- NVIDIA Parallel Thread Execution (PTX) not Cuda.

Phase 4

Model ecosystem

- Model distillation into smaller models (even down to 1.5B parameters).

Multi-head attention

* Example

$$b, \text{num_token}, d_{\text{in}} = (1, 3, 6).$$

$$x = \text{torch.tensor}([[[1.0, 2.0, 3.0, 4.0, 5.0, 6.0], [6.0, 5.0, 4.0, 3.0, 2.0, 1.0], [1.0, 1.0, 1.0, 1.0, 1.0, 1.0]]])$$

step2 : Decide d_{out} , num_head

↳ suppose $\text{d}_{\text{out}} = 2$.

So, then dimension of each head will be $\frac{d_{\text{out}}}{\text{num_head}} = 3$.

step3 :- Initialize trainable weight matrices for key, query, value (W_k, W_q, W_v)

Suppose $6 \times 6 \rightarrow d_{\text{out}}$, its like this cause we want our Q, k, V must be as input dim.

Eg:- If we have

$$\begin{matrix} 3 \times 6 \\ \text{num_token} \quad \text{din} \end{matrix} \cdot \begin{matrix} 6 \times 6 \\ \text{weight} \end{matrix} \rightarrow \begin{matrix} 3 \times 6 \\ \text{num_token} \quad \text{d_out} \end{matrix} \quad \text{by the way same dim.}$$

so, when

$(1, 3, 6)$ then we get,
keys $(1, 3, 6)$, Queries $(1 \times 3 \times 6)$ & values $(1 \times 3, 6)$

batch no. of tokens d-out

Now, as query, key, value is splitted as no.of heads.

Here we have

$$\begin{matrix} (1, 3, 6) \\ \text{batch} \quad \text{no. of tokens} \quad \text{dout} \end{matrix} \xrightarrow{\text{we unroll it as we ke}} \begin{matrix} (1, 3, 2, 3) \\ \text{made d_out as num_head, head_dim} \end{matrix}$$

$$T_1 \{ [1, 2, 3], [3, 5, 6] \}^{h_1}$$

$$T_2 \{ [7, 8, 9], [1, 2, 9] \}^{h_2}$$

$$T_3 \{ [1, 2, 8], [1, 2, 9] \}$$

heads, $T = \text{tokens}$.

In our case

$$\begin{matrix} (1, 3, 2, 3) \\ \text{head_dim.} \end{matrix} \xrightarrow{\text{we reshape the place.}} \begin{matrix} [[1, 2, 3], [7, 8, 9], [1, 2, 8], [1, 2, 9]] \\ T_1 \\ T_2 \\ T_3 \end{matrix} \xrightarrow{\text{make } (1, 2, 3, 3) \text{ for easy calculation.}} \begin{matrix} H_1 \\ H_2 \end{matrix}$$

& similarly, we will have for Queries, keys, values.

Now,
for attention $Q \cdot k^T$

$$\Rightarrow (1, 2, 3, 3) \cdot (1, 2, 3, 3)$$

$$= (1, 2, 3, 3) \cdot \text{Keys} \cdot \text{transpose}(2, 3)$$

$$\Rightarrow \begin{bmatrix} Q_1 \cdot k_1^T \\ Q_2 \cdot k_2^T \end{bmatrix} = (1, 2, 3, 3) \cdot (1, 2, 3, 3)$$

batch no.of head
 no.of token
 no.of tokens

d_{out}
 no.of tokens

batch no.of head
 no.of tokens

no.of tokens

As. $Q \cdot k^T$ gives attention matrix

score
 batch no.of heads

1, 2, 3, 3
 batch no.of heads

Now, we have to find attention weights

$$\text{Attention} = \text{softmax}\left(\frac{Q \cdot k^T}{\sqrt{d_k}}\right) \cdot V$$

So,
Now, we Mask the attention scores to implement causal Attention.

$$\text{tensor}\left(\begin{bmatrix} [-8.2, -\inf, -\inf], \\ [-1.3, -7.07, -\inf], \\ [-5.89, -2.72, -1.016] \end{bmatrix}, \begin{bmatrix} [4.65, -\inf, -\inf], \\ [-1.31, 1.39, -\inf] \\ [2.38, 2.72, 0.848] \end{bmatrix}\right)$$

$$\text{Divide by } \sqrt{\text{head-dim}} = \sqrt{6/2} = \sqrt{3}$$

& apply softmax.

$$\hookrightarrow \begin{bmatrix} [[1, 0, 0], \\ [0.96, 0.03, 0]] \\ [[1, 0, 0], \\ [0.17, 0.82, 0], \\ [0.38, 0.46, 0.15]] \end{bmatrix}$$

This gives attention weights matrix.

Now, last step:

$T = \text{Token}$ | $H = \text{head}$

Context vector = Attention weights * Values.

$$= (b, \text{num-head}, \text{num-token}, \text{num-token}) * (b, \text{num-head}, \text{num-token}, \text{head-dim})$$

$$\Rightarrow (b, \text{num-head}, \text{num-token}, \text{head-dim})$$

$$\Rightarrow (1, 2, 3, 3)$$

Then we
concatenate

$$\Rightarrow \text{tensor}([[[[1, 2, 3],$$

$$T_1 \rightarrow [2, 3, 6], H_1 \\ T_2 \rightarrow [6, 3, 3]]]$$

$$T_3 \rightarrow [1, 2, 3, 6], H_2$$

$$[2, 1, 1, 0], H_2 \\ [1, 1, 0, 0]]])$$

Now, we merge
such that,
we get,

$$(b, \text{num-token}, d_out)$$

$$\Rightarrow (b, \text{num-token}, \text{num-head}, \text{head-dim})$$

$$\Rightarrow (1, 3, 2, 3)$$

$$\Rightarrow \text{tensor}([[[[1, 2, 3], H_1, \text{we will merge them} \\ T_1 \rightarrow [3, 2, 1], H_2]$$

$$T_2 \rightarrow [3, 2, 1], H_2$$

$$T_3 \rightarrow [1, 5, 6]]]$$

$$[3, 2, 2], H_1 \\ [5, 5, 0]]]$$

$$[1, 1, 1]$$

\Rightarrow Final context vector :

$$\Rightarrow \text{tensor}([[[[-3.6, 2.09, 1.38, 1.59, -0.9, -0.34]$$

$$T_1 \rightarrow [-1.2, 2.08, 1.3, 1.07, -0.1, -0.96]$$

$$T_2 \rightarrow [-1.1, 2.08, 1.3, 1.00, -0.9, 1]]])$$

$$T_3 \rightarrow [1, 1, 1]$$

Note

- i) we only need context vector for the last token in input sequence to predict next token.
- ii) we observe, many computations seem to be repeated during the inference (as we are finding the context vector of all previous token each time).

K-V cache

store the keys & values of a previous token, so to minimize the unnecessary re-computation.

* without caching the computational difficulty is $O(n^2)$
with caching its $O(n)$.

* size of KV cache: \rightarrow context length. \rightarrow number of bytes per floating point!

$$l * b * n * h * s * 2 * 2 \rightarrow \text{caches } (k, v)$$

↓ ↓ ↓ ↓ ↓ ↓
layers batch no. of heads head dimension

In deepseek R1/V3

$$L = 61$$

$$b = 1$$

$$n = 128$$

$$h = 128$$

$$s = 100000$$

$$\text{KV cache size} = 400\text{GiB}$$

~~different query, same keys, value weight matrix (w_k, w_v)~~

* we use,

- ↳ Multi-query attention
- ↳ group -query attention

 to solve KV cache problem.

Example keys weight matrix.

•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
wk ₁ (8,2)	wk ₂ (8,2)	wk ₃ (8,2)					

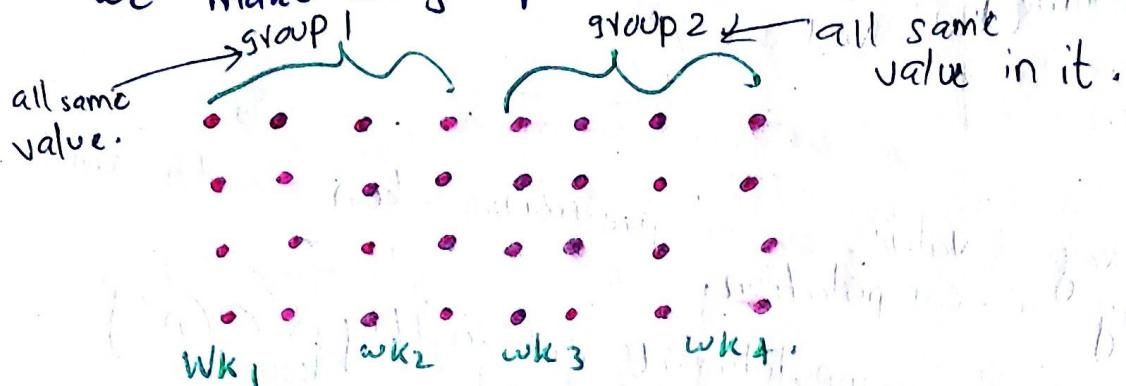
values weight matrix.

•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
wv ₁	wv ₂	wv ₃					

what we have done, we just make number of heads (n) = 1.
so, it reduced size by $\left(\frac{1}{n}\right)$ factor.

* Group-query & attention

↳ it tells instead of making all value same in w_k, w_v ,
we make as. group of heads as $\left(\frac{1}{n}\right)$ si same value.



Here, GQA is a medium level approach for having low cache size & medium language model performance.

* To achieve both less cache & high model performance

MLA is introduced.

↳ multi-head

latent attention.

* they ask a simple question: do we have to cache keys & value separately?

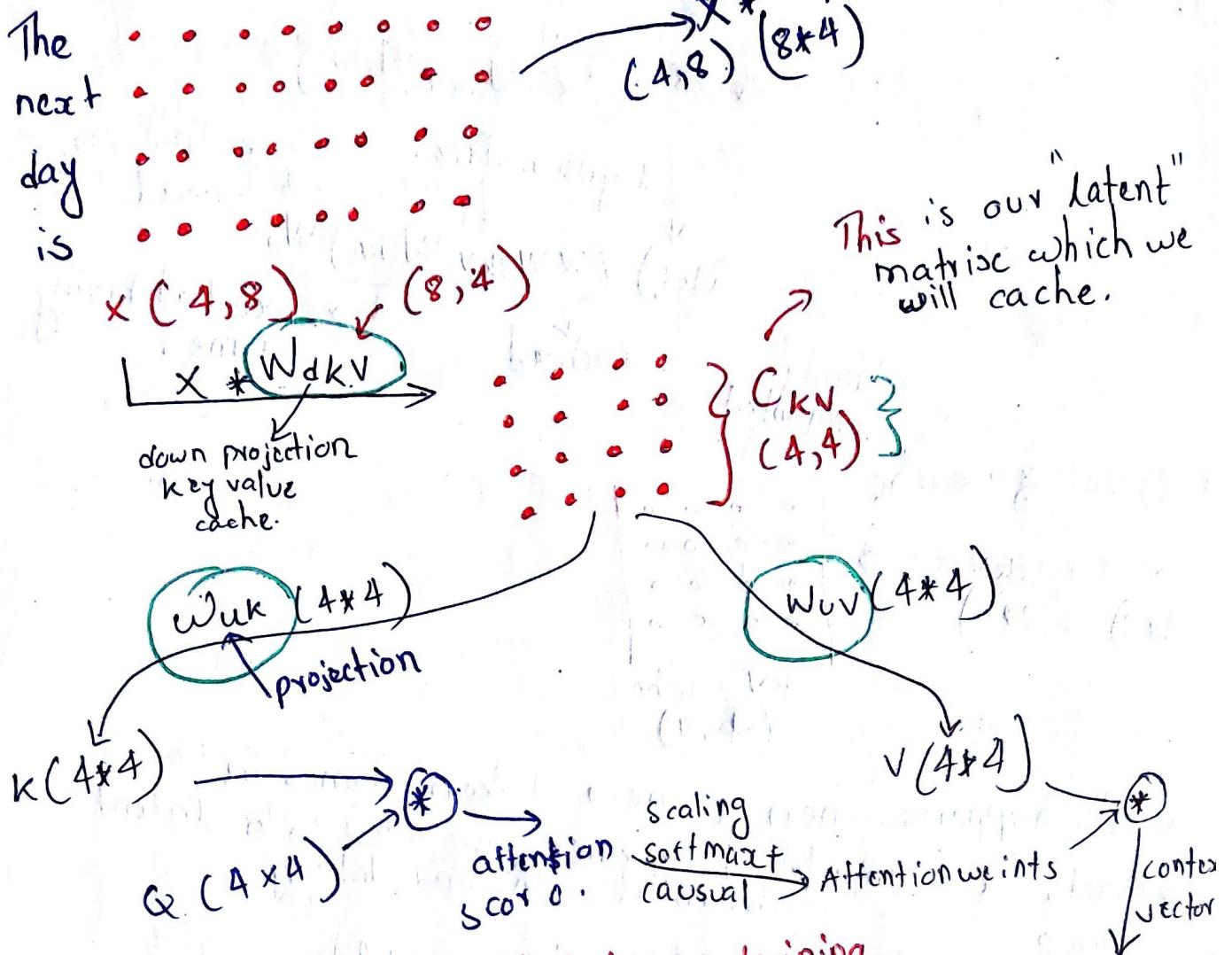
① What if we don't have to cache only one matrix.

what if we cache only one matrix.

what if this matrix has less dimensions than " $n * h$ "?

↳ so, to get that matrix, we start by projecting the input embedding matrix into a latent space.

Let say we have



Note: Wakv, Wuk, Wuw are trained during pre training.

* How does adding this latent matrix help???

$$1) Q = X * W_q$$

$$\text{ii) } C_{kV} = X^* W_d k_V$$

$$\text{iii) } R = C_{KV} * W_{UK} = \frac{X * W_{dKV} * W_{UK}}{X * W_{dKV} * W_{UV}}$$

$$\text{iii) } K = C_{WV} * W_{WK} = X * W_{DKV} * W_{WV}$$

$$\text{iv) } Q = C_{WV} * W_{WV} = \text{"action trick"}$$

$\nabla = \text{C}_{\text{UV}} * \text{W}_{\text{UV}} -$
 "Absorption trick"
 $\hookrightarrow \underline{\text{Attention scores}} = Q * k^T$

$$\hat{w}_{dkv} = \hat{w}_{dkv}^T * \left(x * w_{dkv} * w_{dkv} \right)^T$$

$$\Rightarrow x^T w_q + (x^T w_{dkv}^T + w_{uk})$$

$$\Rightarrow x^T w_q + (w_{uk}^T + w_{dkv}^T + x^T)$$

$$\Rightarrow x^T (w_q * w_{dk})^T (x w_{dk})$$

Absorbed
Query

fixed at
training

This need to be
cached.

it Ckv
latent cache.

④ Context Vector Matrix: Attention weights \mathbf{V}

$$(\mathbf{Q}^T) \times (\mathbf{W}_{d\mathbf{K}} \mathbf{V})$$

logits matrix

$$(\mathbf{Q}^T) \times (\mathbf{W}_{d\mathbf{K}} \mathbf{V}) \mathbf{W}_0$$

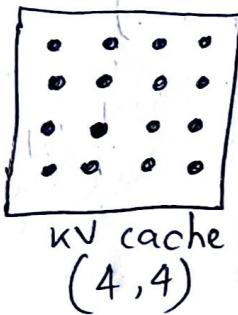
output projection head
fixed at training time.

already computed

cached

* latent KV cache

$$\mathbf{x} * \mathbf{W}_{d\mathbf{K}} \mathbf{V} \rightarrow (4 \times 8) \times (8 \times 4)$$



* what happens when a new token comes in?

↳ First, we compute the queries projected into latent space

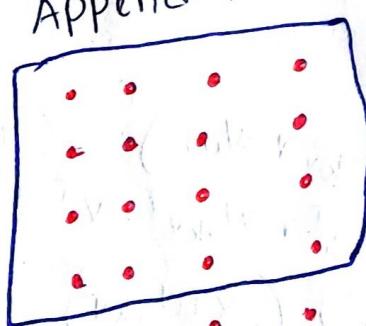
$$\textcircled{1} \quad Q: \mathbf{x} \left(\mathbf{W}_q \mathbf{W}_{u\mathbf{k}}^T \right) \xrightarrow{\substack{\text{query vector} \\ \text{for bright}}} (1 \times 4)$$

suppose
 x_{bright} (1, 8) (8, 4) (4, 4)

(II) Compute KV vector.

$$\mathbf{x}_{\text{bright}} * \mathbf{W}_{d\mathbf{K}} \mathbf{V} \rightarrow \dots (1, 4)$$

Append to latent KV cache
previous KV cache.

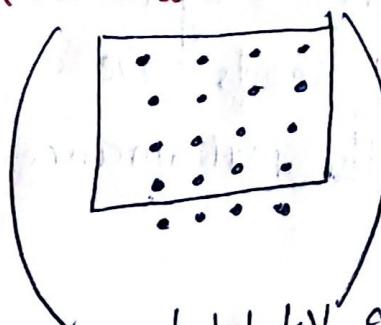


updated KV cache
(5, 4)

③ Multiply Q with updated KV cache.

$$Q: X (W_q W_{kv}^T) *$$

$$\begin{matrix} \dots & \dots & \dots \\ \text{x}_{\text{bright}} & \downarrow & \downarrow \\ (1,8) & (8,4) & (4,8) \end{matrix}$$

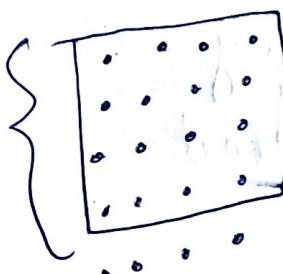


updated KV cache
(8,4)

Attention scores for "bright"
(1,5)

Attention "bright" weights for.
(1,5)

context vector for
bright (1,4)



updated KV cache
(8,4)

* W_{kv}
(4,4)

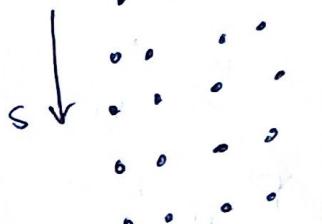
This is later absorbed into the logits calculation.

values
matrix.

Now we don't need two separate K-cache & V-cache, we just need one latent cache (C_{kv}).

Note, we store, we
we store.

C_{kv} only.



cache chosen by ourself.

now size

$2 * l * b * s * dk$ & also dimension of KV cache

one "2" is gone no "any k & v" nah!!!.

Here, W_{Qk} & W_{Vk} have weights different for each attention head. Thus, all heads have different k, V values.

→ This solves the performance issue with M QA.

multi-query attention.

* Binary positional encoding

I saw a dog. The dog was black in color.

suppose,

64: 0 1 0 0 0 0 0 0
 65: 0 1 0 0 0 0 0 1
 66: 0 1 0 0 0 0 1 0
 67: 0 1 0 0 0 0 1 1
 68: 0 1 0 0 0 1 0 0

; so on.

Here, lower indices oscillate fast between positions.

Higher indices oscillate slow between positions.

* Main problem:

- * Since, the values of the integer positional encoding are discrete, they lead discontinuities in the resulting vectors.
- * During pre-trained, it becomes difficult for the LLM optimization routine to deal with those "jumps" or discontinuities.

* best practice

→ graph with smooth curve no jumps.

* Sinusoidal Positional Encoding:

$$PE(pos, 2^i) = \sin\left(\frac{pos}{10000^{2^i/d_{model}}}\right)$$

↑ index of dimension

$$PE(pos, 2^i+1) = \cos\left(\frac{pos}{10000^{2^i/d_{model}}}\right)$$

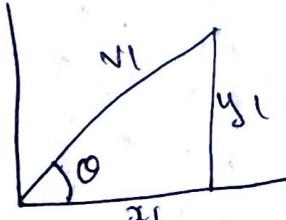
suppose,

$$i=0, PE(pos, 2) = \sin\left(\frac{pos}{10000^{2/1/d}}\right)$$

$$\& PE(pos, 3) = \cos\left(\frac{pos}{10000^{2/1/d}}\right)$$

suppose $\frac{pos}{10000^{2/1/d}} = 0$.

} for index 2 & 3

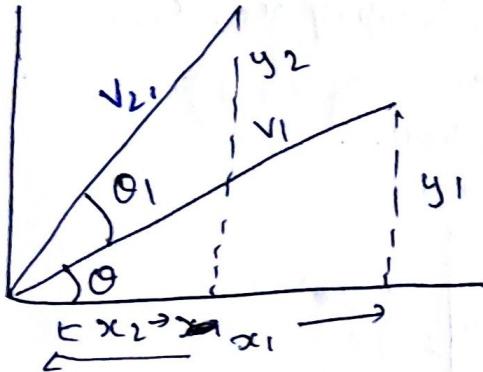


Here $\sin \theta = y_1$ & $\cos \theta = \frac{x_1}{d}$ positional encoding value of 3rd position.
 positional encoding
 value of 2nd position
Now, to find position encoding for

p+k position

~~$tk = \sin(\omega \cdot p)$~~ → pos $p+k+1$

~~$\frac{1}{10000} \sin(2\pi \frac{p+k+1}{10000})$~~ P.E (pos^{tk}, 2^(p+k+1)) = $\sin \left(\frac{\text{pos} + k}{10000^{2/d}} \right)$
 ~~$\frac{1}{10000} \cos(2\pi \frac{p+k+1}{10000})$~~ & P.E (pos^{tk}, 2^(p+k+1)) + 1 = $\cos \left(\frac{\text{pos} + k}{10000^{2/d}} \right) \rightarrow \sin \left(\frac{\text{pos} + k}{10000^{2/d}} \right)$



Here $\theta + \theta_1 = \omega(p+k)$.
 Here ~~sin~~.

$\rightarrow \cos \left(\frac{p+k}{10000^{2/d}} \right)$

\therefore Here, P.E for $(\text{pos} + k)$, we just need to rotate the vector v_1 by an angle.

here also conclude.
 $\theta = \omega p$ & $\theta_1 = \omega k$.

④ v_1 & v_2 are rotations of each other.
 * What this mean is : Relative positional encodings are just rotations of each other.

* Positional embeddings are directly added to token embeddings.
 This pollutes the semantic information carried by token embeddings.
 ↳ solution:
 ↳ can we preserve original vector magnitude by merely rotating the vector.

\therefore So, RoPe introduced,
 i) In rotary embedding, we are going to look at augmenting the queries & the keys vectors. we are not going to touch my token embedding. vectors at all.
 ii) we are just rotating.

Rope Rotary Positional Encoding

suppose

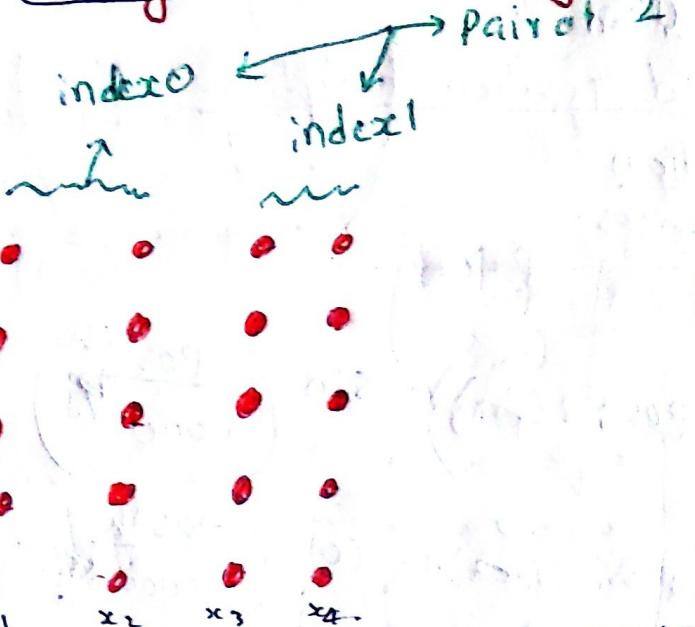
Pos 1 → The

Pos 2 → dog

Pos 3 → chased

Pos 4 → another

Pos 5 → dog.

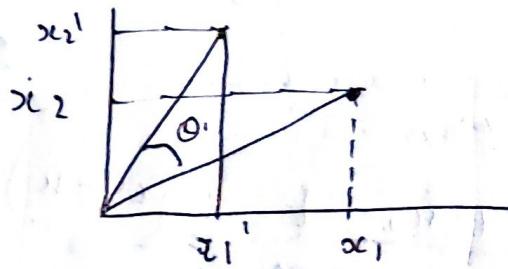


Here for Pos 1, we take the first two suppose x_1, x_2 .

If we apply $x_1', x_2' = \sin, \cos$ of w.p.

where, $\omega = \frac{1}{10000} 2^\circ / \text{pd}$ & p = position.

& again same of x_3 & x_4 .



we apply this in query.

so, our original query-key

$x_1 x_2 x_3 x_4$

& our positionally encoded query-key.

$x_1' x_2' x_3' x_4'$

$$\text{w.p. } \frac{p}{10000} 2^\circ / \text{pd}$$

* This also conclude

low index, high oscillation

↳ means model captures these small shifts, example,

"I just told her the truth" vs "I told just her the truth".

↳ The low index fast oscillations capture the change brought about by varying the position of the word "told".

Here, "told" change position like from 2 to 3, this is captured by low index, high oscillation.

- * Higher index components ensure that even with large position differences, the relationship is preserved.
- ↳ It captures the long range context dependencies.

Decoupled Rotary Position Embedding

using MLA with Rope.

$$\Rightarrow \text{Attention scores} = Q \cdot k^T$$

$$\Rightarrow X \cdot W_q \cdot (C_{kv} \cdot W_{uk})^T$$

$$= X \cdot W_q \cdot W_{uk}$$

(absorbed)

but,

with rope.

$$\text{Attention scores} = \underset{\text{Pos}}{\text{Rope}} (X \cdot W_q) \cdot \cancel{\underset{\text{Pos}}{\text{Rope}}} R_{\text{pos}}$$

due to R_{pos} ,

now, $X \cdot W_q \cdot W_{uk}$
cannot be absorbed.

↳ due to this reason, we may need to recompute the keys for all tokens during inference.

* To avoid this, Deepseek proposed Decoupled RoPE

What is it says??

$$\Rightarrow \text{Attention scores} = Q \cdot k^T$$

$$= [Q_c : R_R] \cdot [k_c \cdot k_R]^T$$

ROPE normal:

$$= Q_c \cdot k_c^T + Q_R \cdot k_R^T \quad \left\{ \begin{array}{l} \text{like the tensor is in} \\ \text{concat so we can do to} \\ \text{half of it as soon nah!!} \end{array} \right\}$$

* we do some extra computation, but it's fine as we insure our absorb technique of MLA.

How de coupled Rope works ???

δ = embedding dimensions.

This is without
Rope used

The
next
day
is

$\text{W}_{\text{QKV}}(8, 4)$

$$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}$$

This is
only
cached
latent

$\text{W}_{\text{UQ}}(4, 8)$

$\text{C}_{\text{KV}}(4, 4)$
 $\text{W}_{\text{UQ}}(4, 8)$

token
taken

$\text{W}_{\text{UQ}}(4, 8)$

$$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}$$

token
 d'
 $\text{C}_{\text{Q}}(4, 4)$

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

k_{Q}

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

d'

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$\text{Q}_{\text{C}}(4, 8)$

Note: we do

$\text{W}_{\text{DQ}} \rightarrow \text{C}_{\text{Q}} \rightarrow \text{W}_{\text{UQ}} \rightarrow \text{Q}_{\text{C}}$ (downprojection, upprojection).

↳ According to the paper, it saves activation memory during the training type.

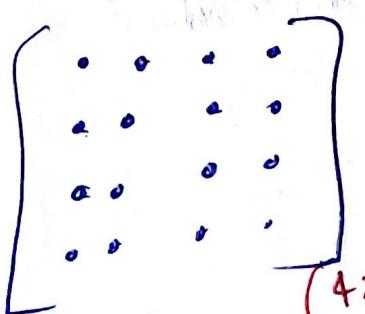
meaning: down-projecting the query matrix into a smaller matrix reduces the activation memory or memory requirements during training time & it's also seems to lead to a better performance.

* With rope

In here, we did not find W_k notice & also, we use same value across different heads

The
next
day
is

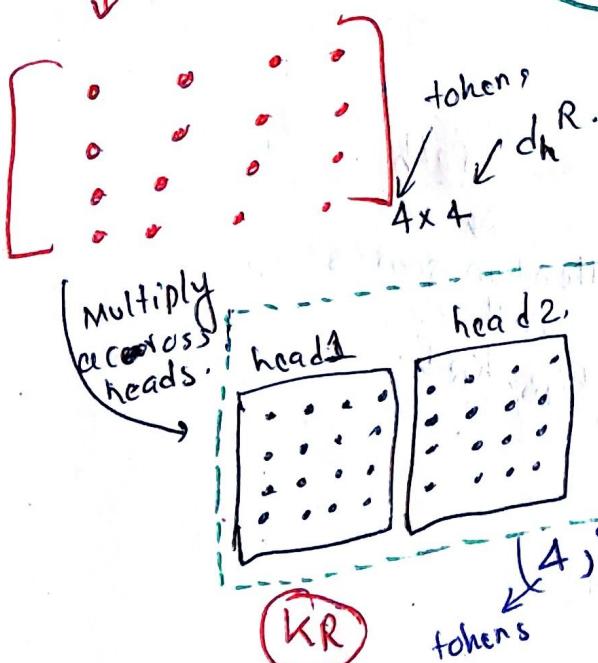
WKR $(8, 4)$ dR we are going to use some weights across diff. multi head.



dimension of head in rotary.

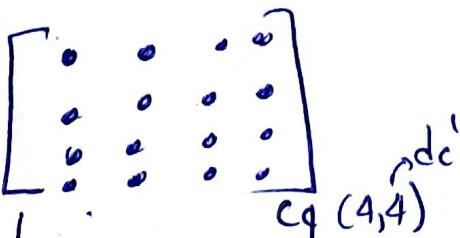
here, key are same in multi head but not query in Rope variant

ROPE

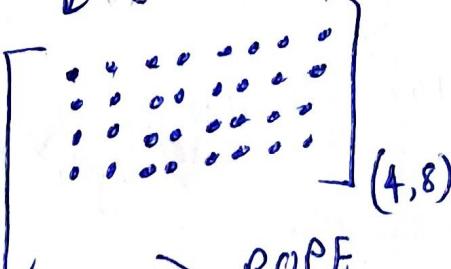


* WDQ $(8, 4)$

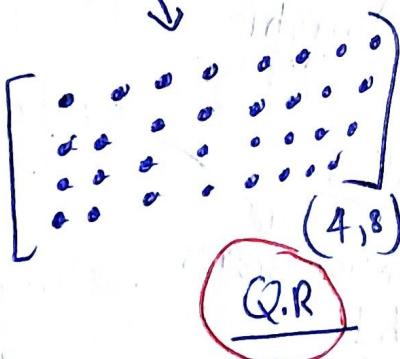
we are going to use some weights across different heads.



* WQP.R $(4, 8)$ $n_h * d_h$.



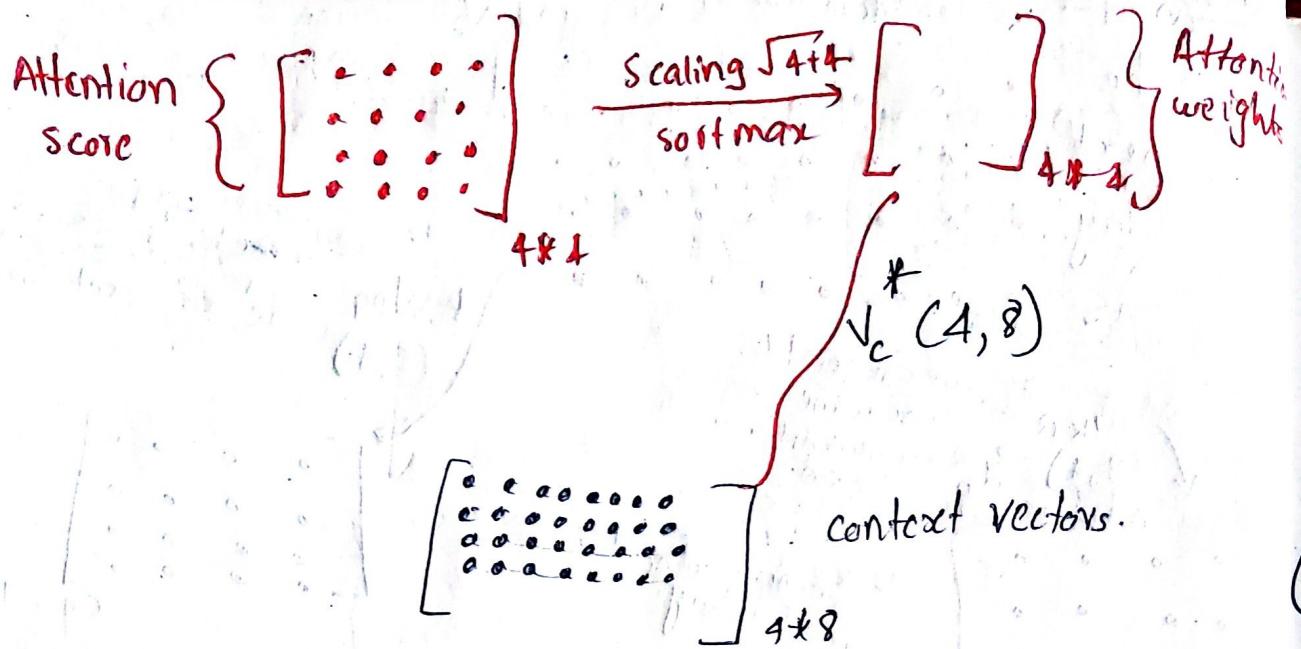
ROPE



Now, attention scores = $Q \cdot K^T$
we concatenate both.

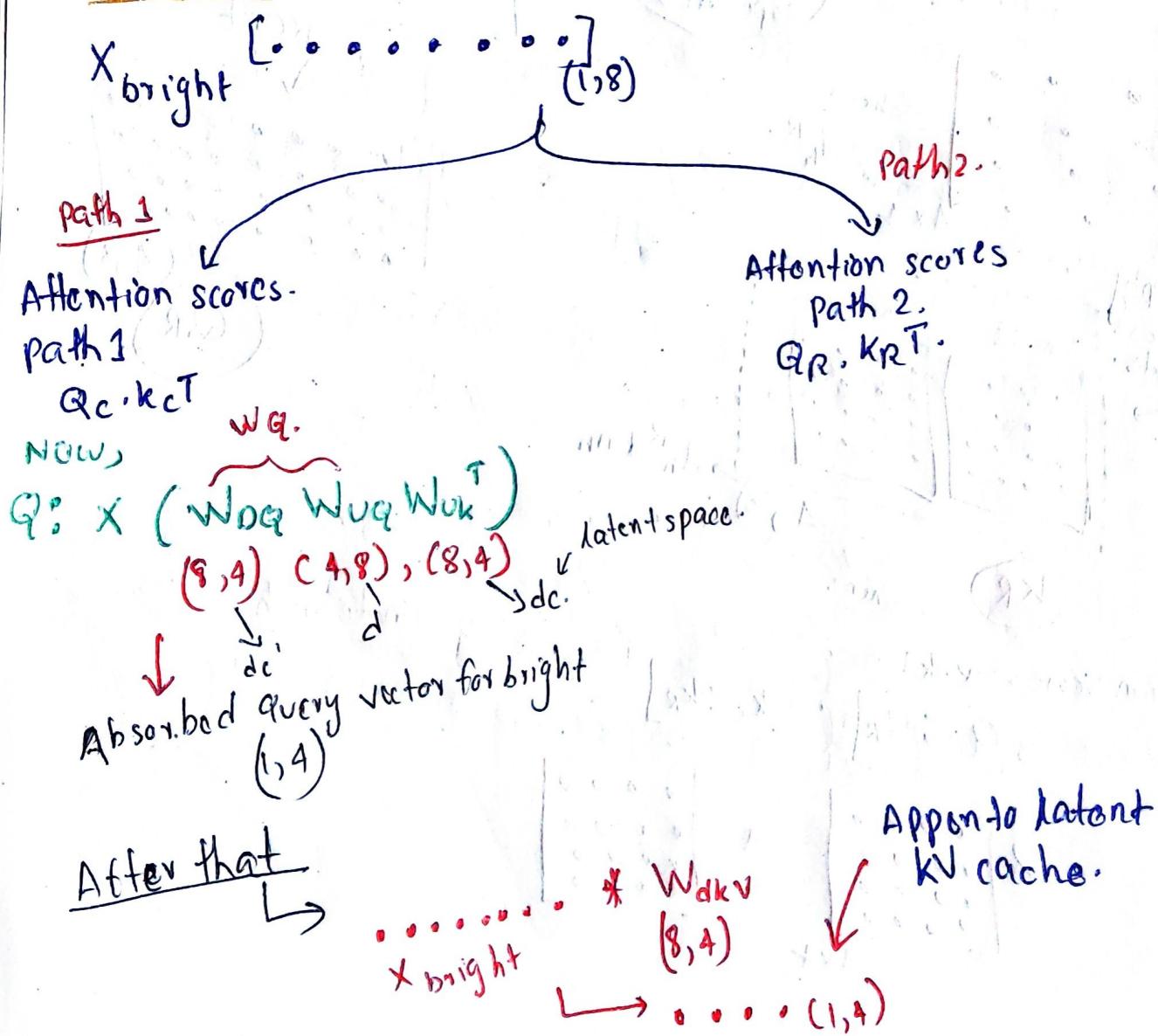
$$\Rightarrow [Q_c : Q_R] \cdot [K_c : K_R]^T$$

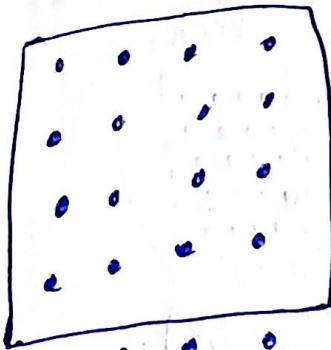
$$\Rightarrow \begin{bmatrix} \dots \\ \dots \\ \dots \\ \dots \end{bmatrix}_{4 \times 4} + \begin{bmatrix} \dots \\ \dots \\ \dots \\ \dots \end{bmatrix}_{4 \times 4}$$



* we only cache C_{kv} & K_R

So what exactly happens when a new token comes in?





original

C_{KV}

why absorbed query help?
because of it we didn't need to cache whole KV values, we just cache the ~~entire~~ C_{KV} latent matrix.

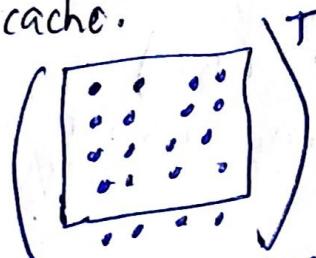
updated KV cache. $(5, 4)$

③ Multiply Q_c with updated KV cache:

$$Q_c : X \left(W_{DQ} W_{UQ} W_{UK}^T \right)^*$$

$(1, 8)$

$(1, 4)$



updated KV cache
 $(5, 4)$

...
Attention scores for "bright"

The next day is bright

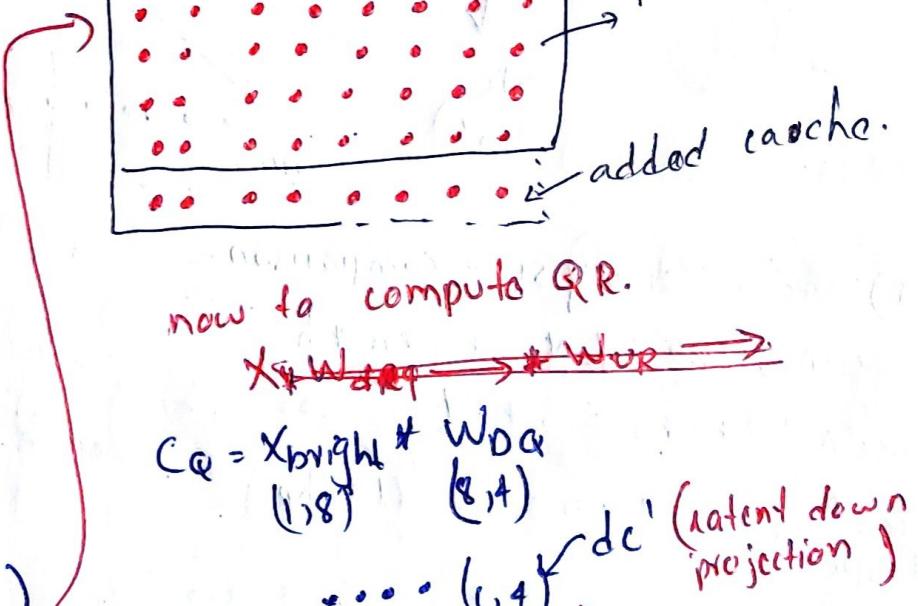
Path 2

$Q_R \cdot K_R^T$

$$x_{bright} * W_{KR} \rightarrow d_h^R$$

$(1, 8) \quad (8, 4)$
...
Rope

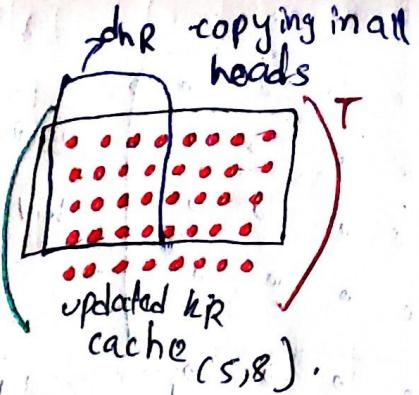
expand
across heads.
we have
to cache it
 K_R



As we have both updated ICR & QR

we find
Attention scores
"bright"

Attention scores for = (1,8)



... (1,5)

After scores for "bright"

Now)

$$\text{Total affection score} = \frac{\text{Attention score from path 1} + \text{Attention score from path 2}}{\text{Path score}} \quad (15)$$

So, total caching dimension = $dc + dh^R$

The diagram illustrates the computation of attention scores and context vectors. It starts with an input vector $v_c(s, 8)$ at the top right. An arrow labeled "Resulting attention score" points to a box containing the formula $\text{Attention}(v_k, v_q) = \frac{\exp(v_k^T v_q)}{\sum_i \exp(v_k^T v_i)}$. Another arrow points from this box to a vector $v_k(s, 4)$. A third arrow points from $v_k(s, 4)$ down to a vector $v_c(s, 8)$. To the left of $v_k(s, 4)$, the text "update k V cache" is written above $(s, 4)$. Above $v_k(s, 4)$, the text "WU9" is written above $(4, 8)$.

4) Cache Memory size comparison.

MHA ; 2*2 * L * b * s * n * h

MLA: $2 * l * b + s * dh$, $\rightarrow^{dh} h_{1/2}$ {in original paper,
here $h = \text{attention dimension}$

$$MLA : 2 * l * b * s * (dL + dHR)$$

MoE

Top-k \rightarrow a token routes to how many experts.

Mixture of Experts

* what is the need to add multiple experts?

- ↳ Adding multiple experts allows models to be pre-trained with far less compared to a dense model (without experts).
- ↳ Allows much faster inference as compared to a dense model (without experts).

* In a dense model, every input token passes through all the parameters (i.e., all layers & neurons). In contrast, a MoE model has multiple "experts" (think of them as sub-networks or feedforward layers), but only a small subset. (e.g., 2 out of 64 are activated for any given input).

Example

Suppose we have 64 experts each 128 dim.

& 2 active experts then

~~2*128~~ 128 parameters

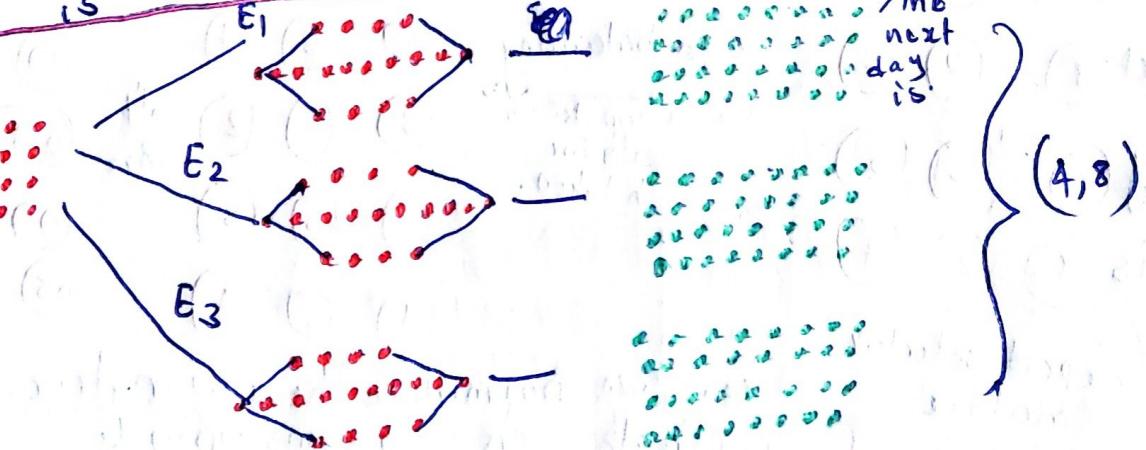
which is 3 times less than a dense layer of 768 dim.

The whole idea is sparsity,

Example :

Suppose input matrix of size (4,8)

The
next
day
is

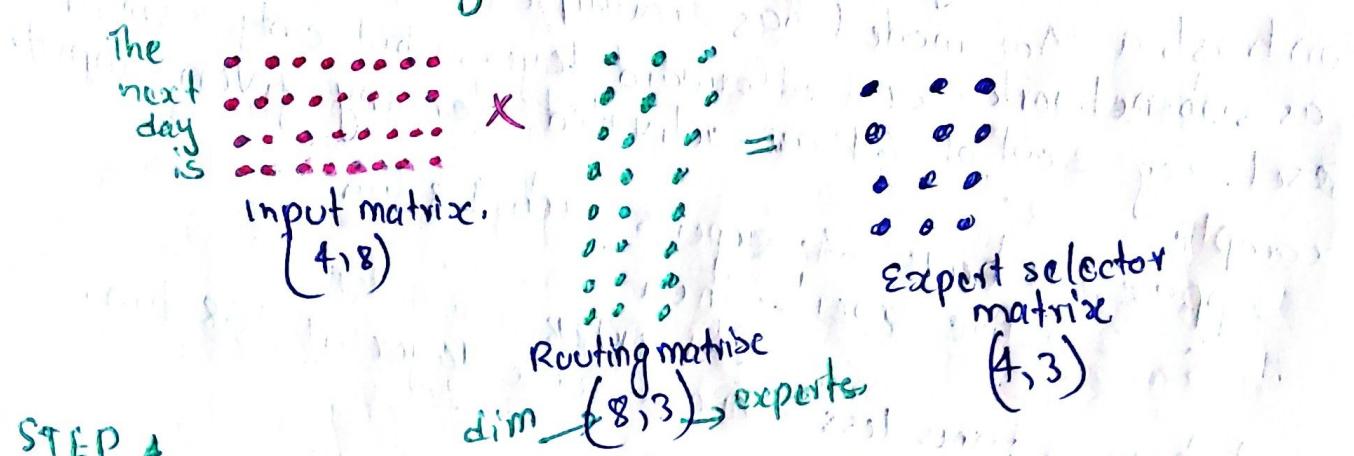


Here, comes the problem, we want to route to only 2 experts the output must be merged but how to do??

Step 2: To merge, we use a method called "load balancing". Ideally, since there are 3 experts; every token will be routed to the 3 experts.

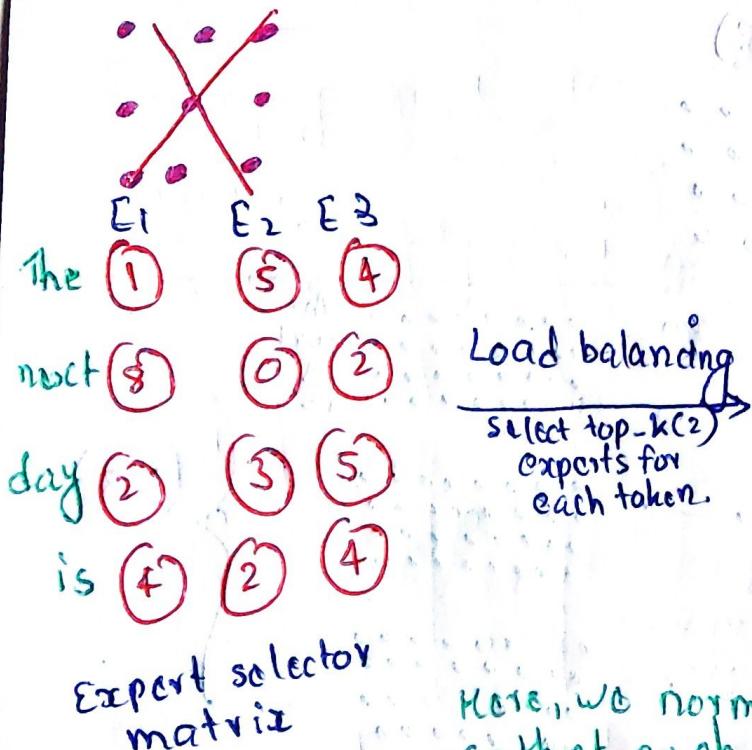
* Every token will only be sent to a selected no. of experts. This is called as load balancing!

Step 3: Once we have decided how many experts will be assigned to each token, the next question is: "How many weightage needs to be given to each expert?" This is decided by the routing mechanism.



STEP 4

Suppose



E1	E2	E3	E1	E2	E3
0	5	4	0	0.6	0.4
8	0	2	0.9	0	0.1
0	3	5	0	0.4	0.6
4	0	4	0.5	0	0.5

Here, we normalize the raw matrix, so that each row sums upto 1. We do softmax.

weight matrix

Step 5
we multiply weight \times output from expert \rightarrow final output.

Example

$E_1 \text{ output} = 0.6 \times \dots + 0.4 \times \dots$
 $E_2 \text{ output}$
 $E_3 \text{ output}$

The \dots $= \dots$ The
Moe layer one token prediction

so, now Auxiliary loss :

\hookrightarrow we need to ensure there is balance, while pre-training,
we also ensure all the experts are essentially utilized.

Here, $E_1 \quad E_2 \quad E_3$ This indicate there is no balance in moe learning.
 $0.6 \quad 0.4$
 $0 \quad 0 \quad 0.1$
 $0.9 \quad 0 \quad 0.6$
 $0 \quad 0.4 \quad 0.6$
 $0 \quad 0 \quad 0.5$
 $0.5 \quad 0 \quad 1.6$
 $1.4 \quad 1 \quad 1$

So, we introduce a loss.
 $\text{Loss} = 1 (C_V)^2$. $C_V = \text{coeff. of variation}$
 $= \frac{\text{standard deviation}}{\text{mean}}$

why we balance?

\hookrightarrow The rich get richer problem :

At the start of training, all experts are initialized randomly. By pure chance, Expert A might be slightly better at a specific token

than Expert B.
• The router sees this & sends more tokens to Expert A.
• The router sees this & becomes even smarter.
• Expert A gets more "practice" (training) & learns nothing. We end up

This happens as of 62 experts sit idle and learn nothing. We end up with a huge model that has the "brain" of a tiny model.
To prevent this, we add a "load Balancing" or "Auxiliary loss" to the main training objective.

Why we need balance?

A) Hardware Efficiency (The Bottleneck).

- ↳ In large-scale training different exports are stored on different GPUs (Expert Parallelism).
 - if one expert is "popular", that specific GPU will be 100% busy.
 - The other 7 GPUs will be 0% busy, waiting for the first one to finish.
 - Impact: our training speed drops to a crawl because we aren't using our hardware in parallel.

B) Waste of Parameters.

- ↳ The whole point of MOE is to have a "Large Brain"
 - If the router only uses 2 experts out of 64, you have wasted 62 experts worth of a GPU memory (VRAM). You are paying for a 50B parameter model but only getting the performance of a 2B model.

* Load Balancing

↳ Assigning equal importance to tokens \neq uniform token routing.

Now??

Suppose

E ₁
1
0
0
0

Export importance of E₁ = 1

This is imbalance, but it's ideal for an example

E₂

0.25
0.25
0.25
0.25

→ Export importance of E₂ = 1

Here although export importance is same, we clearly see not uniformly distributed tokens.

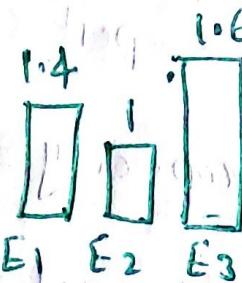
If a token sent to each expert is non-uniform, it can lead to a high memory & a reduced MoE performance.

To balance it we use load balancing.

first, we calculate the probability that the router chose the given expert. (P_i)

suppose, we have 4 tokens & 3 Experts.

so, $P_i = \frac{1}{4} \times 1.4$ Expert importance
probability of one token



Next step is to calculate the fraction of tokens dispatched to each experts. (f_i)

How is it calculated ??

Suppose

	E1	E2	
T1	0.6	0.4	
T2	0.4	0.6	
T3	0.6	0.4	
T4	0.4	0.6	

$2 = \text{Expert importance} = 2.$

$$P_i = \frac{1}{4} \times 2 \text{ for } E1$$
$$f_i^o = \frac{1}{4} * \frac{1}{sum\ of\ count\ of\ token\ to\ that\ expert}$$

$$\text{for } E2 \Rightarrow P_i^o = \frac{2}{4}$$

$f_i^o = \frac{1}{4} \times 2.$
as only T1 & T3
will be routed
to E1 as it is
high probability.

$$\text{So, Loss} = \sum f_i^o \times p_i^o$$

we will minimize this

* Capacity factor, Export capacity

↳ suppose all token got routed to only one expert
this cause the previous problems.

* So,
Export capacity = the maximum number of tokens that a single expert can process in a given batch.

$$\text{Export capacity} = \frac{\text{Tokens per batch}}{\text{Number of experts}}$$

* Capacity Factor.
if capacity factor = 1 & 1000 token in batch & experts.

Then, 250 max token is handle by each expert.

Tokens per batch = batch size * sequence length + top-k.
top-k = no. of experts chosen for each token (load balancing)

* if capacity factor = 1, each experts have equal tokens.
if capacity factor > 1, some have more token & some has less tokens.

if it's too huge it make heavily imbalance.

so, generally it is in range of 1.125 - 2

if capacity factor < 1, some tokens may get dropped if all experts are full.
its not prefer to drop, as deepseek V3, they did not drop any token.

How deepseek Rewrote MoE

Deepseek Innovation

Auxiliary loss free load balancing

we know,

$$\text{load balancing loss} = \text{Scaling factor} * \text{No. of experts} \sum_{i=1}^{\text{num experts}} f_i p_i$$

The use of this loss helps maintain load balance in expert models, but it also acts as a regularization term that interferes with language modeling.

If the scaling factor is too small or absent, it leads to poor balance and routing collapse, reducing model efficiency by limiting expert utilization.

A large scaling factor ensures load balance but degrades overall model performance.

This tradeoff presents a fundamental problem: balancing experts effectively without compromising training quality.

$$\text{model-loss} = \text{Training loss} + \lambda \sum_{i=1}^{\text{num experts}} f_i p_i$$

expert set selector weight matrix.
(A, B).

Then How?

Load Balance
without loss

	E ₁	E ₂	E ₃
T ₁	0	0.6	0.4
T ₂	0.3	0	0.1
T ₃	0	0.4	0.6
\uparrow	0.5	0	0.5
	2	2	4

Token routed

Total token routed = 8

$$\text{Avg. load per expert} = 8/3 = 2.67.$$

Based on the average load per expert, we can know whether a given expert is overloaded or underloaded.

$$E_1 = 2.67 - 2 = 0.67 \quad \text{underloaded}$$

$$E_2 = 2.67 - 2 = 0.67 \quad \text{underloaded}$$

→ Load violation.

$$E_3 = 2.67 - 4 = -1.33$$

overloaded.

Load violation

= Avg load per Expert - Load for a given expert.

Now, next step, we introduce bias.

For each experts we introduce a bias term b_i^0 .

The bias is initialized as 0 for each expert.

Here is how the bias term is updated for each expert.

$$b_i^0 = b_i^0 + \alpha * \text{sign}(\text{load violation error})$$

α = pre-defined constants.

Here, for, expert E_1 , E_2 both are underloaded, so, we increase the bias & for E_3 is over loaded, we decrease the bias.

$$b_i^0 = b_i^0 + \alpha * (+)$$

$$= b_i^0 + \alpha.$$

	E_1	E_2	E_3
$0+b$	$0.6+b$	$0.4+b$	$0.4+b$
$0.9+b$	$0+b$	$0.1+b$	$0+b$
$0+b$	$0.4+b$	$0.6+b$	$0.6+b$
$0.5+b$	$0+b$	$0.5+b$	$0.5+b$

Adding bias term will make sure that the values of expert 1 & expert 2 increases. This will increase the probability of this experts.

Note: The bias value is not a "trainable parameter" in the pytorch / tensorflow sense (it ~~has~~ has require_grad = False). It is a running statistic updated by a simple if-else rule.

* The bias is only added during the Selection phase. Once the exports are chosen, the model often uses the original "unbiased" scores to determine how to combine.

shared experts

In this approach the experts are divided in two groups.

- **shared Experts**: Experts that process every token, regardless of routing
- **Routed Experts**: Experts that handle selectively, based on the usual routing strategy.

Traditional MoE face:

i) knowledge Hybridity

↳ limited experts, so,

each experts will have diverse knowledge.

ii) knowledge Redundancy

↳ multiple experts may converge in acquiring shared knowledge in their respective parameters, leading to redundancy in expert parameters.

- * Common & general information are handled by shared experts.
Note, all token needs to be pass through shared experts.

Fine-Grained Expert Segmentation

↳ The core idea is to divide each expert into multiple smaller experts while maintaining the overall model capacity & computational cost.

↳ In fine-grained expert segmentation, each large expert FFN is split into m small parts, reducing the hidden dimension by a factor of $1/m$.

To compensate for the reduced capacity per expert, the model activates m times more experts per token while keeping the same total computational cost.

* MBPP evaluation
↳ most basic python programming

Multi-Token Prediction

* Deepseek renovate this technique also.

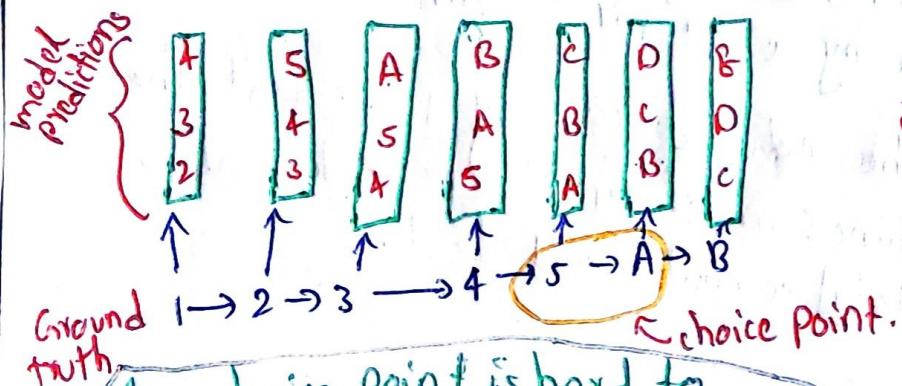
* Instead of prediction only one token we predict the three next token & do back propagation.

why is Multi-Token Prediction useful?



MTP-trained models achieved better results on standard benchmarks like HumanEval, MBPP with the same amount of training data, solving about 15% more code problems on average.

suppose:



since token A is a part of multi predictions (appears in predictions starting from tokens 3, 4, & 5), errors related to predicting A appears repeatedly in the loss calculation. Hence, the training process implicitly prioritizes improving predictions of such consequential tokens, focusing the model's capacity on more critical decision rather than inconsequential ones.

It assigns high weight to consequential choice point tokens.

The choice point is hard to transition in one token prediction, we need something which must see some unknown transition will occur in near future.

what is speculative Decoding?

↳ use two model to predict fast.

↳ Draft model (SLM):

- quickly predicts the next k tokens

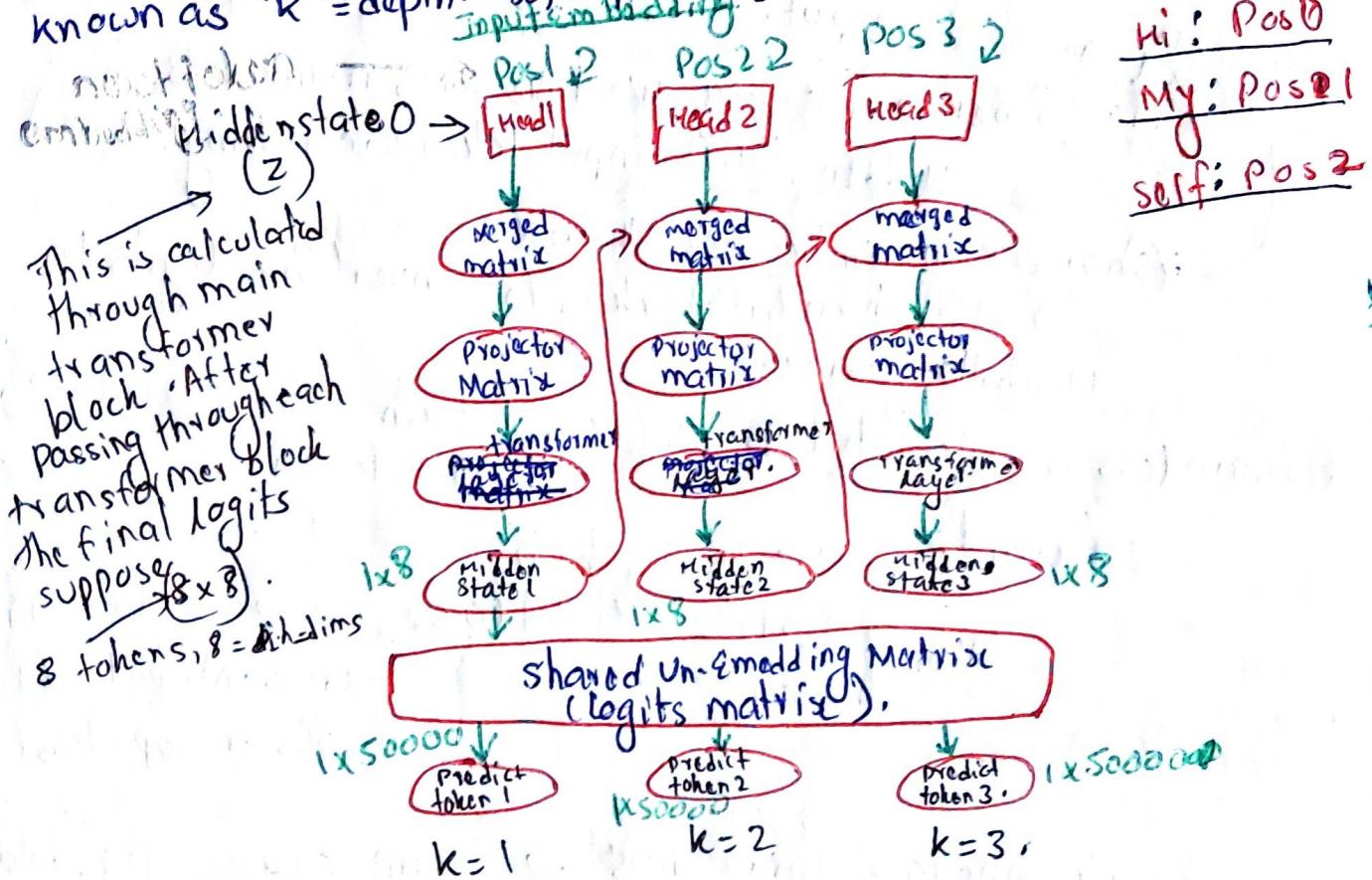
↳ Main LLM

↳ verifies the tokens generated by the draft model & corrects them as needed.

- Deepseek use MTP during pre-training only. During Inference they discarded MTP.

How deep seek implement MTP?

suppose we are predicting for next 3 tokens. in the paper, it is known as ' k ' = depth. So, depth = 3. goes as input



Here, suppose

input sequence length, $T=88$

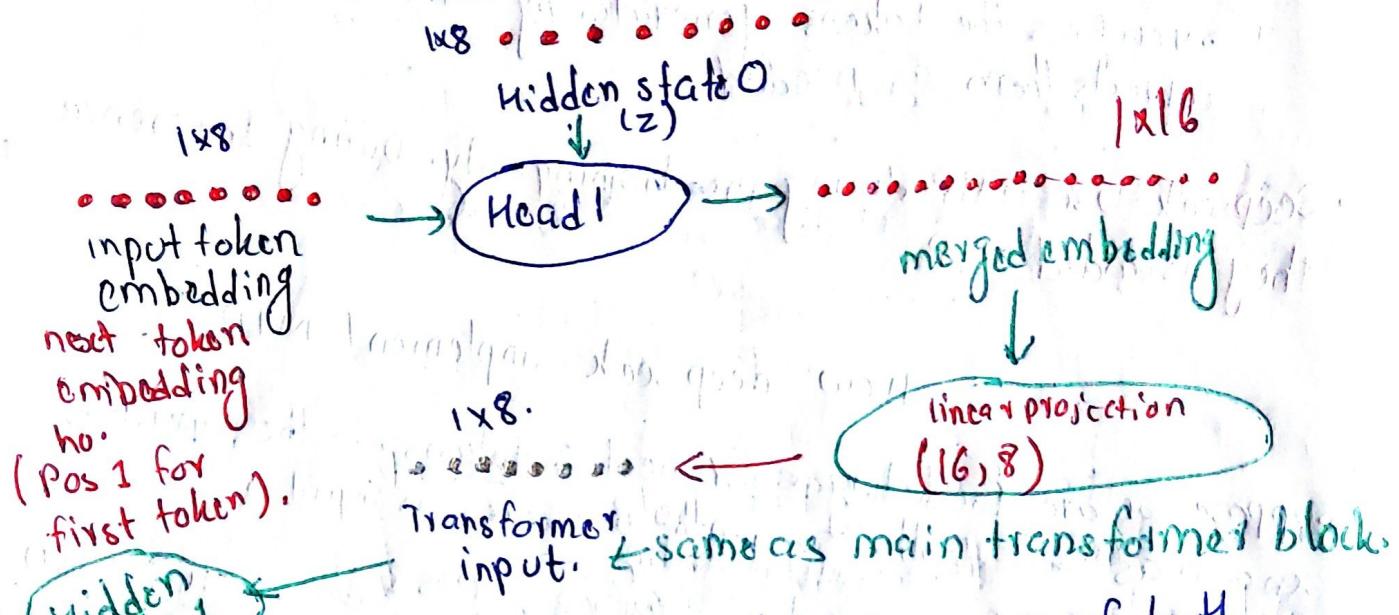
prediction depth $D=3$

This means: we can predict only for $i=4$ as i is from 0-7. because only the last token that most have 3 output token is at index $i=4$.

The first MTP paper predict each next token independently with each other but,

Deepseek introduce sequentially prediction keeping the complete causal chain at each prediction depth.

each head working?

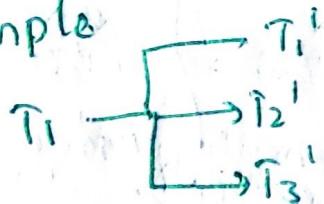


Note: before merging, we take RMS norm of both, hidden state & input token embedding.

↳ shared un-embedding Matrix
↳ it is shared, it work is to convert the logits to vocabulary size.

* Main loss is calculated as \sum (loss of each token predicted from a token)

example

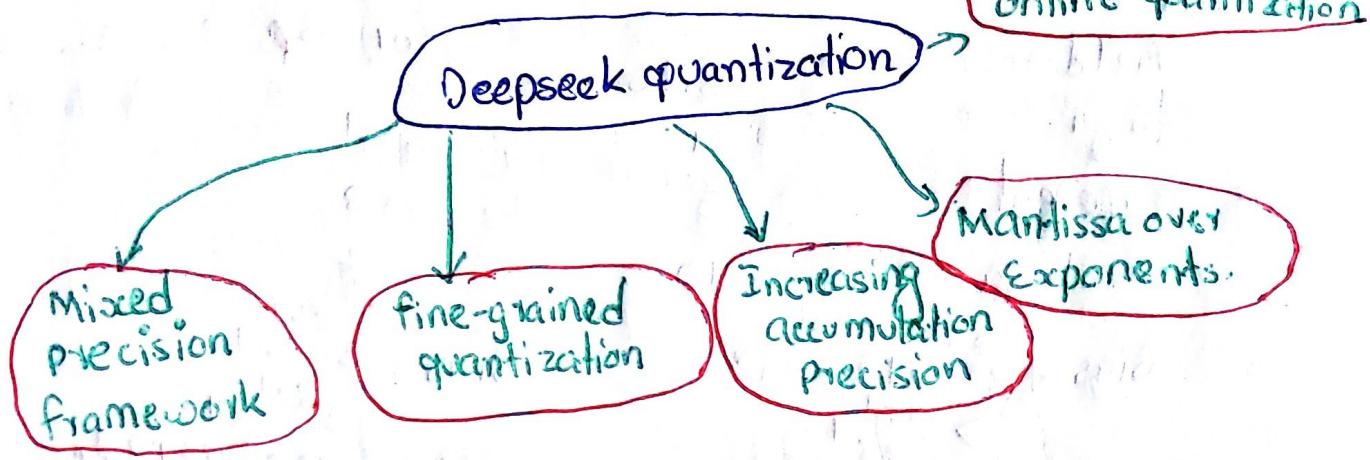


The T_2 ,

total loss = cross entropy loss(T_1, i_1')
+ cross entropy loss(T_2, i_2')
+ cross entropy loss(T_3, i_3')

Deepseek innovate, there must be some causality, hidden state between multi-token prediction.

Quantization



How a number is stored in FP32??

→ According to IEEE-754 standard.

field	bits
Sign	1 bit
Exponent(e)	8 bits
Mantissa	23 bits

for FP32 bias = 127

suppose $x = 3.1415927410125732$.

Binary(approx) = 11.001001000011111

Normalized = $1.100100100011111 \times 2^e$

$$= (e)^s \times (1 + \text{mantissa}) \times 2^{e-127} \quad \text{(formulae)}$$

stored 32bit

0	1	0	0	0	0	0	0
---	---	---	---	---	---	---	---

8bit

1bit

1	0	0	1	0	0	0	1	1	1
---	---	---	---	---	---	---	---	---	---

23 bits

To find the exponent bit we compare the formulae with the Normalized binary of the number.

$$\text{we get, } e-127 = 1$$

$$e = 128$$

we add bias to minimize the negative to occur.

for FP16

field

Sign

Exponent

Mantissa

$$\text{Bias} = 2^{e-1} - 1 \Rightarrow 2^5 - 1 = 15$$

Bits

1

5

10

1

8

7

for Bf16

$$\text{Bias} = 2^{e-1} - 1 = 2^{8-1} - 1 = 128 - 1 = 127$$

Mixed Precision Framework

↪ suppose

$$y = Wx$$

computed initially
in FP32 for numerical
stability then converted
back to BF16.

↪ converted to FP8 from BF16.

↪ maintained in higher precision (BF16 or FP32),
then converted to FP8 on the fly.

(Fprop in paper/figure).

$$y_{\text{FP32}} = W_{\text{FP32}} \times x_{\text{FP8}}$$

↪ stored in BF16 to optimize memory.

In the middle of a neural network,

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} * W^T$$

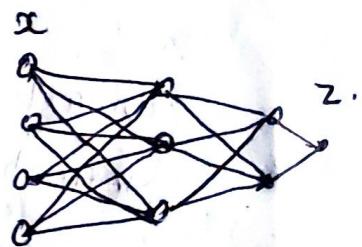
↪ it comes as,
 $z = Wx$

$$\frac{\partial z}{\partial x}$$

computed
initially as FP32, then
converted to BF16.

stored as
BF16 but
converted to

FP8 for
computation.



we know,

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} * \frac{\partial z}{\partial x} = \frac{\partial L}{\partial z}$$

To maintain
dimension in
back propagation

Dgrad
in paper

$$\frac{\partial L}{\partial x_{\text{FP32}}} = \frac{\partial L}{\partial z_{\text{FP8}}} * W^T_{\text{FP8}}$$

↪ stored as BF16.

weights are stored at higher
precision so that we can confirm,
there's update in it. it's learning!!

W grad in paper

$$\frac{dL}{dW} = X^T \frac{\delta L}{\delta Z}$$

↓
stored in
FP8.

computed
+ stored in
FP32 to ensure
numerical
stability for weight
updates.

→ converted from BF16 to FP8

why to do in such approaches?

FP8 Precision

- ↳ significantly reduces computational & memory cost.
- ↳ provides up to 2x speed improvement over BF16 operations.
- ↳ ideal for computationally intensive GEMM operations.

Higher Precision Retained (FP32/BF16):

- 1) Embedding modules, output heads, gating modules, normalization operators, attention operators are sensitive. Hence, these are retained in higher precision.
- 2) Weights & optimizer states also stored in FP32 ensure training stability.

This balance of high-precision storage with low-precision calculation allows the best trade-off between computational efficiency, memory reduction & numerical stability.

b) Fine grained quantization

↳ This method was introduced to protect the precision of weight dimension when they are quantized to lower precision.

* If any outlier occur, it makes the numbers with different values in quantization & dequantization.

$$\text{Suppose: } x = [2, 3, 4, 500] \times \frac{127}{500} = [0.508, 0.762, 1.016, 127.0]$$

Now, quantize to FP8. Assume EAM3 quantization

$$= [0.5, 0.75, 1.0, 127]$$

Dequantize back

$$\text{Recovered Tensor} = [0.5, 0.75, 1.0, 127] \times 500 \approx [1.97, 2.95, 3.94, 500]$$

This occurs due to outlier.

So, what to do??

We have billion parameters in weight, in h-dim also, we have

large no. of values, so,

we will group them in groups & then do quantization,

so, that outlier in one group didn't affect the precision of

another group.

$$x = [1.2, 2.1, 2.001, 3, 76, 0.1, 1.2, 3]$$

group 1 group 2

scaled by factor 1 scaled by factor 2.

$$\Rightarrow = \boxed{x} \times \boxed{W}$$

tensor core. (General matrix mul occurs in FP8 inside NVIDIA Tensor core.)

In CUDA core.

$$\text{FP8} * \text{upscale factor} = \text{FP32}$$



w11 → scaling factor 1

w12 → scaling factor 2

w21 → scaling factor 3

w22 → scaling factor 4.

Increasing accumulation Precision

→ Using higher-precision arithmetic only for the accumulation (summing) step during neural-network computations. Even though the inputs, weights, & outputs are stored in lower precision.

We know, in almost every neural network layer, we have

$$y_i = \sum_{j=1}^d x_j \cdot w_{ji}$$

This has two different operations.

- i) Multiply $\rightarrow x_j \cdot w_{ji}$
- ii) Accumulate (sum) $\rightarrow \sum_j (\cdot)$.

Deepseek-V3 keeps (1) in low precision & (2) in higher precision.

why low-precision accumulation is dangerous??

→ suppose Deepseek uses FP8/BF16 for speed problem:

→ FP8/BF16 has few mantissa bits.

→ when summing thousands of terms, small values get rounded away

- lost entirely.

This causes:

- large numerical error.
- unstable training.
- exploding / vanishing gradients.

WGMM

Low-precision matrix Multiply Accumulate.

→ It's a local-level GPU instruction family introduced with NVIDIA Volta (H100) & it's central to how models like V3 run FP8 efficiently.

• key: multiplication in low precision in tensor core & accumulation is done in FP32 in CUDA core.

Mantissa over Exponents → E4M3 → 3 mantissa bit higher precision
in FP8 there are two format → E5M2.

Deepseek prefer
Always use floating point number
E4M3 format.

but prior works prefers floating point number in percent

* Forward pass: E4M3

* Backward pass: E5M2.

Online quantization

↳ Scaling factor is decided based on the numbers that actually come in iteration (Delayed quantization)

main problem of offline quantization
↳ we run a calibration dataset, and decide the scaling factor at once & after that we do quantization based on that.

* If a rare token appears, overflow or loss.

What deepseek does?
↳ it calculate the scaling factor for the fly

Tensors

* In inference, each weight tensor has one (or a small set of) scale value(s), computed ONCE, & reused for all tokens.

so, in inference, we use offline quantization.