

2025 USA-NA-AIO Round 2, Problem 2, Part 9

 Topics

 My posts

 More

 CATEGORIES

 General

 Site Feedback

 All categories



[Skip to main content](#)

USAAIO 

May 2025

Now, let us study another variant of MHA: **Multi-head Latent Attention (MLA)**. MLA was introduced by **DeepSeek**. It is a core component of DeepSeek's large language model (LLM).

The key intuition of MLA is as follows. In MHA, the key and value projection matrices

$$\mathbf{W}^{\mathbf{K},MHA} \in \mathbb{R}^{H \cdot D_{qk} \times D_2}, \quad \mathbf{W}^{\mathbf{V},MHA} \in \mathbb{R}^{H \cdot D_v \times D_2}$$

may be high dimensional.

For instance, suppose $H \cdot D_{qk} = H \cdot D_v = D_2 = 4096$.

However, it is not necessarily the case that these matrices are with high ranks (such as 4096). Their actual ranks (or top few ranks that make their truncated singular value decomposition (SVD) to be close to the actual matrices) may be much lower than that.

To capture the low-rank feature, MLA proposed the following model:

$$\begin{aligned}\mathbf{W}^{\mathbf{K},MHA} &= \mathbf{W}^{\mathbf{U}\mathbf{K},MLA} \mathbf{W}^{\mathbf{D}\mathbf{K}\mathbf{V},MLA} \\ \mathbf{W}^{\mathbf{V},MHA} &= \mathbf{W}^{\mathbf{U}\mathbf{V},MLA} \mathbf{W}^{\mathbf{D}\mathbf{K}\mathbf{V},MLA}\end{aligned}$$

where

- $\mathbf{W}^{\mathbf{D}\mathbf{K}\mathbf{V},MLA} \in \mathbb{R}^{r \times D_2}$: down-projection matrix for computing keys and values.
- $\mathbf{W}^{\mathbf{U}\mathbf{K},MLA} \in \mathbb{R}^{H \cdot D_{qk} \times r}$: up-projection matrix for computing keys.
- $\mathbf{W}^{\mathbf{U}\mathbf{V},MLA} \in \mathbb{R}^{H \cdot D_v \times r}$: up-projection matrix for computing values.

In practice, rank r is typically much smaller than $\min \{H \cdot D_{qk}, H \cdot D_v, D_2\}$.

In all remaining parts of this problem, to simplify your analysis and highlight the relationships of MHA, GQA and MLA, we make the following assumptions:

- $D_1 = D_2 = D$.
- $D_{qk} = D_v = d$.
- d is a factor of D .

Under these assumptions, the number heads H satisfies

$$H = \frac{D}{d}.$$

Part 9 (10 points, non-coding task)

In this part, you are asked to prove that GQA can be equivalently represented by MLA.

In your solution, it is sufficient for you to prove that for $M \in \{K, V\}$, for matrix

$$\tilde{W}^{M,GQA} = \begin{bmatrix} W^{M,GQA} \\ W^{M,GQA} \\ \vdots \\ W^{M,GQA} \end{bmatrix} \in \mathbb{R}^{D \times D}$$

(defined in Part 6) who is the concatenation of $\frac{H}{G}$ copies of

[Skip to main content](#)

$$\mathbf{W}^{\text{M},GQA} = \begin{bmatrix} \mathbf{W}_0^{\text{M},GQA} \\ \vdots \\ \mathbf{W}_{G-1}^{\text{M},GQA} \end{bmatrix} \in \mathbb{R}^{G \cdot d \times D}$$

matrix $\tilde{\mathbf{W}}^{\text{M},GQA}$ can be decomposed as

$$\tilde{\mathbf{W}}^{\text{M},GQA} = \mathbf{W}^{\text{UM},MLA} \mathbf{W}^{\text{DKV},MLA}$$

where

- $\mathbf{W}^{\text{DKV},MLA} \in \mathbb{R}^{r \times D}$: down-projection matrix for computing keys and values.
- $\mathbf{W}^{\text{UM},MLA} \in \mathbb{R}^{D \times r}$: up-projection matrix for computing \mathbf{M} (keys or values).
- $r = G \cdot d$.



May 2025

Misplaced '#'

We have

$$\begin{aligned} \text{rank}(\tilde{\mathbf{W}}^{\text{M},GQA}) &= \text{rank}(\mathbf{W}^{\text{M},GQA}) \\ &\leq \min\{G \cdot d, D\} \\ &= \min\{r, D\} \\ &= r, \end{aligned}$$

[Skip to main content](#)

where the first equality follows from Part 6.

Therefore, SVD implies

$$\begin{aligned}\tilde{\mathbf{W}}^{\text{M},GQA} &= \sum_{i=0}^{r-1} \sigma_i \mathbf{u}_i \mathbf{v}_i^\top \\ &= \underbrace{[\mathbf{u}_0 \quad \mathbf{u}_1 \quad \cdots \quad \mathbf{u}_{r-1}]}_{\mathbf{W}^{\text{UM},MLA}} \underbrace{\begin{bmatrix} \sigma_0 \mathbf{v}_0^\top \\ \sigma_1 \mathbf{v}_1^\top \\ \vdots \\ \sigma_{r-1} \mathbf{v}_{r-1}^\top \end{bmatrix}}_{\mathbf{WDKV},MLA}.\end{aligned}$$

"" END OF THIS PART """

◆ Related topics

Topic	Replies	Activity
2025 USA-NA-AIO Round 2, Problem 2, Part 12	1	May 2025
2025 USA-NA-AIO Round 2, Problem 2, Part 6	2	May 2025
2025 USA-NA-AIO Round 2, Problem 2, Part 11	1	May 2025
2025 USA-NA-AIO Round 2, Problem 2, Part 1	2	Dec 2025

[Skip to main content](#)

Topic	Replies	Activity
2025 USA-NA-AIO Round 2, Problem 2, Part 10	1	May 2025

 Powered by Discourse