

2025 USA-NA-AIO Round 2, Problem 2, Part 6

USAAIO 

May 2025

Next, let us study a variant of MHA: **Group Query Attention (GQA)**.

Recall that in MHA, the number of heads in queries, keys and values are the same, H . Thus, query $\mathbf{q}_{l_1,h}$ attends to key $\mathbf{k}_{l_2,h}$ with the same head index h .

In GQA, we relax this constraint by allowing keys and values to have G heads ($G \leq H$), where G is factor of H . For instance, if $H = 12$, then $G \in \{1, 2, 3, 4, 6, 12\}$.

In GQA, a query $\mathbf{q}_{l_1,\textcolor{red}{h}}$ with head $\textcolor{red}{h}$ is permitted to attend to a key $\mathbf{k}_{l_2,\textcolor{blue}{g}}$ and use value $\mathbf{v}_{l_2,\textcolor{blue}{g}}$ in computing its output with head $\textcolor{blue}{g}$ if

$$\textcolor{red}{h} \equiv \textcolor{blue}{g} \pmod{G}.$$

Thus, each head in keys and values is mapped to $\frac{H}{G} \geq 1$ heads in queries.

As an example, suppose $H = 12$ and $G = 3$. Then

- Head $\textcolor{blue}{g} = 0$ in keys and values is associated with heads $\textcolor{red}{h} = 0, 3, 6, 9$ in queries.
- Head $\textcolor{blue}{g} = 1$ in keys and values is associated with heads $\textcolor{red}{h} = 1, 4, 7, 10$ in queries.
- Head $\textcolor{blue}{g} = 2$ in keys and values is associated with heads $\textcolor{red}{h} = 2, 5, 8, 11$ in queries.

[Skip to main content](#)
USAAIO 

May 2025

Part 6 (5 points, non-coding task)

For $\mathbf{M} \in \{\mathbf{K}, \mathbf{V}\}$, Denote the \mathbf{M} -projection matrix as

 Topics

 My posts

 More

 CATEGORIES

[General](#)

[Site Feedback](#)

 All categories

$$\mathbf{W}^{\mathbf{M},GQA} = \begin{bmatrix} \mathbf{W}_0^{\mathbf{M},GQA} \\ \vdots \\ \mathbf{W}_{G-1}^{\mathbf{M},GQA} \end{bmatrix}$$

Now, we concatenate $\frac{H}{G}$ copies of the above matrix along axis 0:

$$\tilde{\mathbf{W}}^{\mathbf{M},GQA} = \begin{bmatrix} \mathbf{W}^{\mathbf{M},GQA} \\ \mathbf{W}^{\mathbf{M},GQA} \\ \vdots \\ \mathbf{W}^{\mathbf{M},GQA} \end{bmatrix}$$

What is the relationship between $\text{rank}(\tilde{\mathbf{W}}^{\mathbf{M},GQA})$ and $\text{rank}(\mathbf{W}^{\mathbf{M},GQA})$?

- Reasoning is required.

USAAIO 

May 2025

Misplaced '#'

Let $\{\mathbf{w}_i^* : i \in \{0, 1, \dots, r-1\}\}$ be r linearly independent row vectors that span all row vectors $\mathbf{W}^{\mathbf{M},GQA}$.

Because each row vector in $\mathbf{W}^{\mathbf{M},GQA}$ has $\frac{H}{G}$ copies in $\tilde{\mathbf{W}}^{\mathbf{M},GQA}$, we must have that $\{\mathbf{w}_i^* : i \in \{0, 1, \dots, r-1\}\}$ also spans $\tilde{\mathbf{W}}^{\mathbf{M},GQA}$.



[Skip to main content](#)

Therefore,

$$\text{rank} \left(\tilde{\mathbf{W}}^{\mathbf{M}, GQA} \right) = \text{rank} \left(\mathbf{W}^{\mathbf{M}, GQA} \right).$$

"" END OF THIS PART """

❖ Related topics

Topic	Replies	Activity
2025 USA-NA-AIO Round 2, Problem 2, Part 9	1	May 2025
2025 USA-NA-AIO Round 2, Problem 2, Part 12	1	May 2025
2025 USA-NA-AIO Round 2, Problem 2, Part 1	2	Dec 2025
2025 USA-NA-AIO Round 2, Problem 2, Part 2	2	Dec 2025
2025 USA-NA-AIO Round 2, Problem 2, Part 7	1	May 2025