

- a) Compute the time and space complexities of applying QLoRA and LoRA to LFFN and MHA separately.

LoRA :

LFFN time complexity for LoRA

$$\text{LFFN}(X) = XW + XA(B)^T$$

- XW : $n \times d \times d \times d$ is $O(nd^2)$
- XA : $n \times d \times d \times r$ is $O(n dr)$
- $XA(B)^T$: $n \times r \times r \times d$ is $O(n dr)$
- $XW + XA(B)^T$: $n \times d + n \times d$ is $O(nd)$

Total time complexity is

$$O(nd^2) + 2 \cdot O(n dr) + O(nd), \text{ which is } O(nd^2)$$

LFFN space complexity LoRA

The trainable parameters are

$$A, B \in \mathbb{R}^{d \times r}$$

This is $O(2dr)$, which is $O(d)$

since $d \gg r$

MHA time complexity LoRA

For queries, keys and values we do

$$HW: n \times d \times d \text{ is } O(nd^2)$$

$$HA: n \times d \times d \times r \text{ is } O(ndr)$$

$$HA(B)^T: n \times r \times r \times d \text{ is } O(ndr)$$

$$HW + HA(B)^T: n \times d + n \times d \text{ is } O(nd)$$

This total time complexity is $O(hn^2d)$

because we have h heads and $d \gg r$

Note, if $X \in \mathbb{R}^{a \times b}$ then
 $\text{softmax}(X)$ is $O(a \cdot b)$

time complexity of $\text{head}(H)$:

$$QK^T: n \times d \times d \times n \text{ is } O(n^2 d)$$

$$\frac{QK^T}{\sqrt{d}}: n \times n \text{ is } O(n^2)$$

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right): n \times n \text{ is } O(n^2)$$

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) * V: n \times n \times n \times d \text{ is } O(n^2 d)$$

total time complexity is

$O(h n^2 d)$ for all h heads

Last step is to compute time
complexity of $\text{MHA}(H)$, which
is $O(h n d^2)$

Final time complexity:

$$O(hnd^2 + hn^2d)$$

Space complexity of MHA LoRA

have matrices $A^{i,l_i}, B^{i,l_i} \in \mathbb{R}^{d \times r}$

for each Q, K, V and h heads.

This is: $O(3h(2d \cdot r)) = O(hd)$,

since $d \gg r$.

QLoRA:

Time complexity for LFFN:

Same steps as LoRA but need
to apply a dequantization step
to $W \in \mathbb{R}^{d \times d}$

$$O(nd^2) + O(d^2) = O(nd^2)$$

Space complexity for LFFN:

Same as LoRA: $O(d)$

Time complexity for MHA:

Same steps as LoRA but
with dequantization steps.

$\text{dequant}(\Delta, W)$ is $O(a \cdot b)$,

given $W \in \mathbb{R}^{a \times b}$

dequantization is done once for Q, K, V for each head h . and the final step with $W^{o.e}$.

This is $O(hnd^2)$.

Using time complexity from LoRA, total complexity is

$$O(hnd^2 + hnd^2) + O(hnd^2) = O(hnd^2 + hnd^2)$$

Space complexity MHA QLoRA :

Same as LoRa, $O(hnd)$

b)

	MHA	LFEN
Memory Limit	QLoRA	QLoRA
Faster Inference	LoRA	LoRA
Both	QLoRA	LoRA

With a memory limit, QLoRA is better because dequantizing saves memory.

For faster inference, LoRA is better because dequantizing takes time.

If you want both lower memory and faster inference, it will depend. MHA has many more parameters so maybe in that

Case the inference penalty you pay with QLoRA is worth it.

By similar logic, LoRA might be worth the memory penalty since there are less parameters.