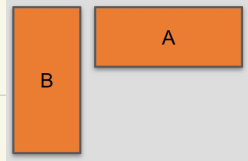


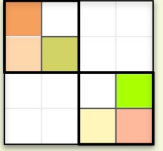
Methods:-

- ① LoRA $\delta W = B \times A$ $B \in \mathbb{R}^{m \times r}$ $A \in \mathbb{R}^{r \times n}$
- ② SVD+LoRA $\delta W = U \Sigma V^T$ $U \in \mathbb{R}^{m \times m}$ $\Sigma \in \mathbb{R}^{r \times r}$ $V \in \mathbb{R}^{n \times n}$



Sparse Fine Tuning:-

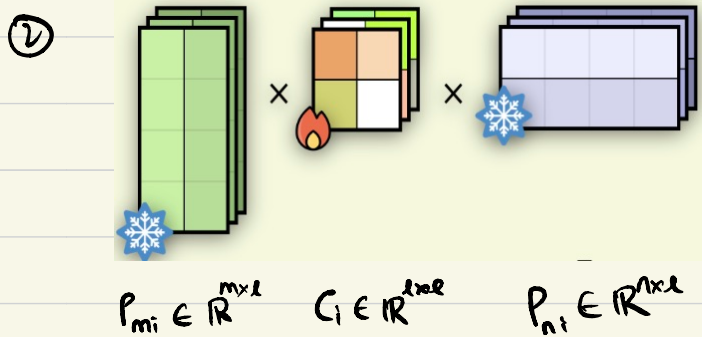
- ① Fourier $\delta W = F C F^T$ $(F) = \text{Inv. Fourier Transform}$
- ② Hadamard $\delta W = H C H^T$ $(H) = \text{Inv. Hadamard Transform}$
- ③ Frame $\delta W = F_r C F_r^T$ $(F_r) = \text{Inv. frame Transform}$
- ④ Random $\delta W = R C R^T$ $(R) = \text{Random values}$



* Computation:-

- ① We will be computing $Y = \delta W X$ $X \in \mathbb{R}^{n \times d}$ (no. of tokens \times hidden dim)

① Direct form : $\delta W X$



$$\sum_{i=1}^l P_{m,i} C_i P_{n,i}^T X$$



$$\sum_i \sum_j P_{m,(i,j)} C_{(i,j)} P_{n,(i,j)}^T X$$

$P_{m,(i,j)} \in \mathbb{R}^{m \times l}$ $C_{(i,j)} \in \mathbb{R}$ $P_{n,(i,j)} \in \mathbb{R}^{l \times d}$

Models:-

- | | | |
|----------------|----------|----------|
| ① Llama-2-7B | $m=2048$ | $n=2048$ |
| ② Gemma-2-2B | $m=2048$ | $n=2048$ |
| | $m=2304$ | $n=1024$ |
| ③ Gemma-2-9B | $m=3584$ | $n=4096$ |
| | $m=3584$ | $n=2048$ |
| ④ Llama-3.1-8B | $m=4096$ | $n=4096$ |
| | $m=4096$ | $n=1024$ |

⑤ Deep Seek !

Number of Tokens:-

→ $n = 1, 256, 1024, 2048, 4096, 8192$

Sparsity blocks:-

- | | |
|-------------------|---------------------------|
| ① Size of C | → $l=2, 4, 16, 64, 256$ |
| ② Sparsity of C | → $75\%, 50\%, 25\%, 0\%$ |

Precision:-

- | |
|--------|
| ① FP32 |
| ② FP16 |
| ③ FP8 |

Threads:-

→ Number of threads - 8, 16, 32, 256, 1024 etc.

