

一、基础方法

1. 直接训练
 - 通过微调直接训练模型
 - 使用预训练模型
2. **SFT**（Supervised Finetune）
 - 通过微调直接训练模型
 - 使用预训练模型
 - 使用LoRA或Adapter进行SFT
3. **RLHF**（Reinforcement Learning with Human Feedback）
 - 通过微调直接训练模型
 - 使用预训练模型
4. 混合方法
 - 混合使用以上方法

SFT（Supervised Finetune）

通过微调直接训练模型

也称为dirty work（脏活）

主要优点：

- **Few-Shot Prompting**：通过提供1-5个示例，让模型学习任务
- **Seed Prompt**：通过设置task_type来指定任务类型
- **参数化**
 - 通过prompt参数
 - task_type：seed prompt和seed参数的pretrain模型
 - seed：prompt
 - 通过answer参数
 - GPT4/Claude3
 - Qwen_72B/deepseek_MoE

主要缺点：



RL（Reinforcement Learning）

通过与环境交互，通过奖励信号进行学习

主要优点：

环境と行動



- **Agent** エージェント
- **action** アクション
- **Environment** エンvironment reward リワード
- **reward** リワード
- **State** ステート

状態と行動のMDP

状態

行動と報酬の組合せ

MDP

1. $\pi: M^2 \rightarrow \{0, 1\}$
2. $R: S \times A \times S \rightarrow \mathbb{R}$
3. $\gamma: M^2 \rightarrow \mathbb{R}$



状態と行動の組合せ

- $s_t, a_t | s_{t+1}, r_{t+1}$

$$\begin{bmatrix} S_1 & \rightarrow & S_1 & \& S_1 & \rightarrow & S_2 & \& S_1 & \rightarrow & S_3 \\ & & & \& & & & \& & & & \end{bmatrix} \begin{bmatrix} S_2 & \rightarrow & S_1 & \& S_2 & \rightarrow & S_2 & \& S_2 & \rightarrow & S_3 \\ & & & \& & & & \& & & & \end{bmatrix} \begin{bmatrix} S_3 & \rightarrow & S_1 & \& S_3 & \rightarrow & S_2 & \& S_3 & \rightarrow & S_3 \end{bmatrix}$$

- $p(s_{t+1}|s_t, a_t) = p(s_{t+1}|s_t)$
- $r_{t+1} = r(s_{t+1})$

$$p(s_{t+1}|s_t) = p(s_{t+1}|s_1, \dots, s_t)$$

- $p(s_{t+1}|s_t) = p(s_{t+1}|s_t)$

- $p(s_{t+1}|s_t) = p(s_{t+1}|s_t)$

$$p(s_{t+1}|s_t) = p(s_{t+1}=s'|s_t=s)$$

価値と行動のMRP

$$\langle S, P, R, \gamma \rangle$$

- $S: \text{状態}$
- $P: \text{遷移確率}$
- $R: \text{報酬} \quad R_S = E[R_{t+1}|S_t=s]$
 - $R_S = \sum_s p(s|s_t) R(s)$
 - $R_S = \sum_s p(s|s_t) R(s)$
 - $R_S = \sum_s p(s|s_t) R(s)$

- γ / $\gamma \in [0,1]$
- Return G_t $t \in \{0, \dots, T\}$
 - $G_0 = R_0$
 - $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$
- $V(s) = E[G_t | S_t = s]$
 - $V(s) = \sum_{t=0}^T \gamma^t P(S_t = s) G_t$
 - $V(s) = \sum_{t=0}^T \gamma^t P(S_t = s) E[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^k R_{t+k+1}]$

$$V(s) = E[G_t | S_t = s]$$

MRPs

$$V(s) = E[G_t | S_t = s] = R_{t+1} + \underbrace{\gamma R_{t+2} + \gamma^2 R_{t+3} + \dots}_{\text{MRP}} = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

MDPs

$$\langle S, A, P, R, \gamma \rangle$$

- S
- A
- $P(s'|s, a) = P(S_{t+1} = s' | S_t = s, A_t = a)$
- $R(s, a) = E[R_{t+1} | S_t = s, A_t = a]$
 - $R(s, a) = \sum_{s'} P(s'|s, a) R(s')$
 - $R(s, a) = \sum_{s'} P(s'|s, a) \sum_{a'} P(a'|a) R(s')$
 - $R(s, a) = \sum_{s'} P(s'|s, a) \sum_{a'} P(a'|a) \sum_{s''} P(s''|s', a') R(s'')$
- $\gamma \in [0, 1]$

□

$$\pi(a|s) = P(A_t = a | S_t = s)$$

- s a
- $\pi(a|s)$
- MDP

$$M = \langle S, A, P, R, \gamma \rangle$$

- M

$$P_{s,s'} = \sum_{a \in A} \pi(a|s) P(s', s | a)$$

- $P(s', s | a) = \sum_{s''} P(s'' | s', a) P(s | s'')$
- $P(s'' | s', a) = \sum_{s''} P(s'' | s', a)$

□

□□□□s□□□□□\$\\pi\\$□□□□

 \$V_{\{\pi\}}(s) = E_{\{\pi\}}[G_t | S_t=s] \quad \text{--- } G_t = \\ \underbrace{R_{\{t+1\}}}_{\{\text{Reward}\}} + \underbrace{\gamma R_{\{t+2\}} + \gamma^2 R_{\{t+3\}} + \dots}_{\{\text{Future Rewards}\}} = \sum_{k=0}^{\infty} \gamma^k R_{\{t+k+1\}} \quad \text{--- } V_{\{\pi\}}(s) = E_{\{\pi\}}[R_{\{t+1\}} + \gamma V_{\{\pi\}}(S_{\{t+1\}}) | S_t=s]

5

□□□□□s□□□□□a□□□□□\$\\pi\\$□□□□□

 \$ q_{\pi}(s,a) = E_{\pi}[G_t | S_t = s, A_t = a] \$
\$ q_{\pi}(s,a) = E_{\pi}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a] \$
\$ V_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s,a) \$

 image-20250817181356989

 image-20250817181417147

10

~~~~~ \$\$ V\_{\{ \}}(s) = \max\_{\{ p | i \}} V\_{\{ \}}(p)(s) ~~~~ \$\$ q\_{\{ \}}(s,a) = \max\_{\{ p | i \}} q\_{\{ \}}(p, a) ~~~~ \$\$

- $\pi_*$  any  $\pi_*$
  - $\pi_*$
  - $\pi_*$
  - $\pi_*$

\$\$ \backslash pi \* \backslash ge any \backslash pi \$\$

$\exists s \forall \pi' (\pi'(s) \geq \pi(s) \wedge \pi' > \pi(s))$

□ □ □ □ □ □ □ □ □ □ □ □

□□□□□□□□□□

```


$$\sum_{s \in S} \sum_{a \in A} p(a|s) q_{\pi}(s, a) = R_s^{\pi} + \gamma \sum_{s' \in S} \sum_{a' \in A} p(s'|s) p(a'|s') q_{\pi}(s', a')$$


```

□□□□□□□□□□□□ \$ V\_{\pi}(s) = \sum\_{a \in A} \pi(a|s) q\_{\pi}(s,a) \$

$\pi(s,a)$

- $\exists \pi_* \forall V_{\{\pi_*\}}(s) \geq V_{\{\pi\}}(s)$
  - $\exists q_{\{\pi_*\}}(s, a) \leq q_{\{\pi\}}(s, a)$
  - $\exists V_{\{\pi\}}(s) \leq q_{\{\pi\}}(s, a)$

```

$$ \pi_{\{ \}}(a | mid s) = \begin{cases} 1 & \text{if } a = \arg \max_{\{a\}} q_{\{ \}}(s, a), \ 0 & \text{otherwise}. \end{cases} \ V_*(s) = \max_a q_{\{ \}}(s, a) \ V_*(s) = R_s^a + \gamma \sum_{\{s'\}} P_{\{s, s'\}}^a V_*(s') \ V_*(s) = \max_a (R_s^a + \gamma \sum_{\{s'\}} P_{\{s, s'\}}^a V_*(s')) \ 

```

□□□□□□□□□□□□□□

```
q_{\{ \}}(s, a) = R_s^a + \gamma \sum_{\{s'\}} P_{\{s, s'\}}^a V_*(s')
```

```
V_{\{ \}}(s) = R_s^a + \gamma \sum_{\{s'\}} P_{\{s, s'\}}^a V_*(s') \ q_{\{ \}}(s, a) = R_s^a + \gamma \sum_{\{s'\}} P_{\{s, s'\}}^a V_*(s') \ q_{\{ \}}(s, a) = R_s^a + \gamma \sum_{\{s'\}} P_{\{s, s'\}}^a \max_{\{a'\}} q_{\{ \}}(s', a')
```

## LLaMA

GPTとLLaMAのTransformerDecoder-only実装

- 
- 
- 

## RMSNorm

□□□□□□□□□□□□□□

```
RMSNorm(x) = \frac{x}{\sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2 + \epsilon}} \cdot \gamma
```

- $x$
- $d$
- $\epsilon$
- $\gamma$

## SwiGLU

□□□□□□□□□□□□□□ Swish □□□□□□□□□□□□

```
Swish(x) = x \cdot \sigma(x)
```

- $\sigma$  Sigmoid

```
GLU(x) = \sigma(W_1x + b_1) \odot (W_2x + b_2)
```

- $\odot$
- $W_1, W_2, b_1, b_2$

```
SwiGLU(x) = Swish(Linear_1(x)) \odot Linear_2(x)
```

□□□

- Swish □□□□□□□□□□□□□□
- GLU □□□□□□□□□□□□□□

## RoPE

□□□□□□□□□□□□□□

## GQA

LLaMA2 3

10

GPT

# Deepseek

MLA

10

1

□□□□□□□□□□

- GPU
  - GPU

□□□□□□□□□□

- DP Data Parallel
  - DDP Distributed Data Parallel
  - FSDP Fully Sharded Data Parallel

## All-Reduce

- $\text{softmax}(\text{GPU})$   $\rightarrow$   $\text{softmax}(\text{CPU})$
  - $\text{softmax}$ 
    - $\text{softmax}(\text{GPU})$
    - $\text{softmax}(\text{CPU})$   $\rightarrow$  **SUM**  $\sqcup$  **MAX**  $\sqcup$
    - $\text{softmax}(\text{CPU})$

### All-Gather

- GPU
  - GPU
  - GPU
  - GPU

1

DP Data Parallel

 Python GIL  CPU



1

1. CPUとGPU間のデータ転送
2. GPU間のデータ転送
3. GPU0による初期化
4. GPUによる実行

問題

問題文:\$\backslash psi\\$問題 \$N\$

- GPU0\$\backslash psi\$(N-1)\$\backslash psi\$
- GPU\$N-1\$:\$\backslash psi\$

問題

- Python GILとCPU
- GPUによる並列化

## DDP(Distributed Data Parallel)

問題

- pytorch rankによる並列化
- MPIによる並列化

問題

- GPU0による初期化
- GPU間のデータ転送
- GPUによる実行
- GPUによるRing-AllReduce
- GPUによる最終結果

問題

問題文:\$\backslash psi\\$問題 \$N\$

問題GPU

- Scatter-Reduce
- All-Gather

問題/問題\$2\backslash psi\\$問題

### Ring-AllReduce

- MPIによる並列化
- GPUによる並列化
- GPUによるGPUによる並列化

問題

### Scatter-Reduce

- nGPUによる並列化
- \$i\$:\$GPU\_j\$ \$(j-i) \% n\$:\$GPU\_{(j+1)}\$:\$GPU\_{(j-1)}\$ \$(j-i)\%n\$

### All-gather

- GPUによる並列化
- \$i\$:\$GPU\_i\$ \$(i+1)\%n\$:\$GPU\_n\$

- $\dots \$i\$ \dots \$GPU\_j\$ \dots \$ (j-i-1) \% n \$ \dots \$ \dots \$ (j-i-2) \% n \$ \dots \$$



## FSDP(Fully Sharded Data Parallel)

မြန်မာ

- မြန်မာရေးရှင်းရွေးကြောင်းရေးရှင်းရွေးGPU
- မြန်မာရေးရှင်းရွေးကြောင်းရေးရှင်းရွေးရှင်းရွေး
- မြန်မာရေးရှင်းရွေးကြောင်းရေးရှင်းရွေးGPUမြန်မာCPU

## DeepSpeed ZeRO-1

မြန်မာ3ဗိုလ်GPUမြန်မာရေးရှင်းရွေးGPUမြန်မာရေးရှင်းရွေး



မြန်မာGPUမြန်မာရေးရှင်းရွေးရှင်းရွေးZeRO-1မြန်မာရေးရှင်းရွေးGPUမြန်မာရေးရှင်းရွေးရှင်းရွေးရှင်းရွေး



မြန်မာ

- မြန်မာ
- မြန်မာ**GPU0****GPU1****GPU2****GPU2**မြန်မာGPU
- မြန်မာ**GPU**မြန်မာ**FP32****FP32****FP16**မြန်မာ
- မြန်**FP16**မြန်မာ**GPU**မြန်မာ

မြန်မာ

မြန်မာ\$\backslash\$psi\$မြန်မာ \$N\$

မြန်မာGPUမြန်မာ

- မြန်မာ\$\psi / \psi \$(N-1)\$\frac{\psi}{\psi} \{N\} \approx \psi\$
- မြန်မာ\$\psi / \psi \$(N-1)\$\frac{\psi}{\psi} \{N\} \approx \psi\$

မြန်မာ/မြန်မာ\$2\psi\$

## DeepSpeed ZeRO-2

ZeRO-2မြန်**FP16**မြန်မာရေးရှင်းရွေးGPUမြန်မာရေးရှင်းရွေး



မြန်မာ

- မြန်မာ
- မြန်မာ**GPU0****GPU1****GPU2****GPU2**မြန်မာGPU
- မြန်မာ**GPU**မြန်မာ**FP32****FP32****FP16**မြန်မာ
- မြန်**FP16**မြန်မာ**GPU**မြန်မာ

မြန်မာ

\$\$\psi \approx \frac{1}{N} \sum\_{n=1}^{N-1} \psi\_n

GPUによる実装

- $\frac{\partial \psi}{\partial w_i} \approx \frac{1}{N-1} \left( \frac{\psi_1 - \psi_N}{w_1 - w_N} \right)$
- $\frac{\partial \psi}{\partial w_i} \approx \frac{1}{N-1} \left( \frac{\psi_1 - \psi_N}{w_1 - w_N} \right) \approx \frac{2\psi}{N}$

$\frac{\partial \psi}{\partial w_i} \approx \frac{2\psi}{N}$

### DeepSpeed ZeRO-3

ZeRO-3のFP16実装



GPUによる実装

- GPUによる実装
- GPUによる実装
- GPUによる実装
- GPUによる実装

GPUによる実装

\$\$\psi \approx \frac{1}{N} \sum\_{n=1}^{N-1} \psi\_n

GPUによる実装

- $\frac{\partial \psi}{\partial w_i} \approx \frac{1}{N-1} \left( \frac{\psi_1 - \psi_N}{w_1 - w_N} \right) \approx \frac{2\psi}{N}$
- $\frac{\partial \psi}{\partial w_i} \approx \frac{1}{N-1} \left( \frac{\psi_1 - \psi_N}{w_1 - w_N} \right) \approx \frac{2\psi}{N}$

$\frac{\partial \psi}{\partial w_i} \approx \frac{2\psi}{N}$

GPUによる実装

TPUによる実装

GPUによる実装

TPUによる実装

TPUによる実装

$Y = XW$

- $X \in \mathbb{R}^{2 \times 2}, W \in \mathbb{R}^{2 \times 2}$

$$\begin{aligned} Y &= \begin{pmatrix} y_1 & y_2 \\ y_3 & y_4 \end{pmatrix} = \begin{pmatrix} x_1 & x_2 \\ x_3 & x_4 \end{pmatrix} \begin{pmatrix} w_1 & w_2 \\ w_3 & w_4 \end{pmatrix} \\ W &= \begin{pmatrix} w_0 & w_1 \\ w_2 & w_3 \end{pmatrix} \\ Y_0 &= XW_0 = \begin{pmatrix} x_1 & x_2 \\ x_3 & x_4 \end{pmatrix} \begin{pmatrix} w_0 & w_1 \\ w_2 & w_3 \end{pmatrix} \\ Y_1 &= XW_1 = \begin{pmatrix} x_1 & x_2 \\ x_3 & x_4 \end{pmatrix} \begin{pmatrix} w_1 & w_2 \\ w_3 & w_4 \end{pmatrix} \\ Y_2 &= \begin{pmatrix} y_1 & y_2 \\ y_3 & y_4 \end{pmatrix} \begin{pmatrix} w_2 & w_4 \\ w_4 & w_1 \end{pmatrix} \\ Y_3 &= \begin{pmatrix} y_1 & y_2 \\ y_3 & y_4 \end{pmatrix} \begin{pmatrix} w_3 & w_1 \\ w_1 & w_2 \end{pmatrix} \end{aligned}$$

$\psi \approx \frac{1}{N} \sum_{n=1}^{N-1} \psi_n$



□□□  
□□□□□□□□□□

□□□ \$\\$ Y = XW \\$\\$ □□□

- $X \in R^{2 \times 2} \cap W \in R^{2 \times 2}$

```

\begin{bmatrix} y_1 & y_2 \\ y_3 & y_4 \end{bmatrix} = \begin{bmatrix} x_1 & x_2 \\ x_3 & x_4 \end{bmatrix}
\begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix} $ GPU $ 
W = \begin{bmatrix} W_0 & W_1 \end{bmatrix}
$ X = \begin{bmatrix} X_1 & X_2 \end{bmatrix} $ Y_1 = 
\begin{bmatrix} x_1 & x_3 \end{bmatrix} \begin{bmatrix} w_1 & w_2 \end{bmatrix} = 
\begin{bmatrix} y_{11} & y_{12} \\ y_{31} & y_{32} \end{bmatrix} Y_2 = \begin{bmatrix} x_2 & x_4 \end{bmatrix}
\begin{bmatrix} w_3 & w_4 \end{bmatrix} = \begin{bmatrix} y_{23} & y_{24} \\ y_{43} & y_{44} \end{bmatrix} $ Y = Y_1 + Y_2 $

```



PP

GPU

10

- Stage Stage GPU
  - Stage Stage



1

- GPU2
  - GPU2
  - GPU1

5

- GPU
  -

□ □ □ □ □ □

## F-then-B

- 亂序输入
  - 乱序输入mini-batch乱序输出mini-batch

1F1B

- `stage4`
  - `stage4`  $\sqcap$  `F42`  $\sqcap$  `stage4`  $\sqcap$  2  $\sqcap$  `micro-batch`  $\sqcap$  `stage4`  $\sqcap$  `F42`  $\sqcap$  `stage4`  $\sqcap$  `F41`  $\sqcap$  `B41`  $\sqcap$  `stage4`  $\sqcap$  1  $\sqcap$  `micro-batch`  $\sqcap$  `stage4`  $\sqcap$  `F41`  $\sqcap$  `stage4`  $\sqcap$  **`F42`**  $\sqcap$  `stage4`  $\sqcap$  **`F41`**  $\sqcap$  `stage4`



- 

## Gpipe

**F-then-B**



Mini-batch → Micro-batches → Mini-batch → 4 Micro-batches

Micro-batch1 → GPU0 → GPU1 → GPU0 → Micro-batch2

Batch Normalization → Gpipe → micro-batch → mini-batch

### PipeDream --- DeepSpeed

1F1B

Gpipe

- mini-batch → micro-batch → GPU
  - Pipeline Flush
    - 数据流
    - 内存/缓存
    - 线程同步
- mini-batch → micro-batch → memory

PipeDream

- micro-batch → micro-batch → micro-batch
- machine1 → micro-batch → forward → machine2
- machine4 → machine1 → 1F1B
- Bubble → GPU



FW

- Machine1 FW5 → Micro-batch 1 → Micro-batch 2 → Micro-batch 3 → Micro-batch 4 → Micro-batch 5
  - Micro-batch 5 → Micro-batch 1 → Micro-batch 5
  - Micro-batch 4 → Micro-batch 4
- Machine2 FW5 → batch1-2 → Machine1 → batch1 → Machine1 → Machine2

Machine

1F1B → PipeDream → Weight stashing → Vertical Sync

- Weight stashing
  - Machine 1
    - Machine 1 → Micro-batch 5 → **Micro-batch 5** → **W<sub>5</sub>**
    - Machine 1 → **Micro-batch 5** → **W<sub>5</sub>**
    - W<sub>5</sub>** → **W<sub>5</sub>**
- Vertical Sync
  - Machine 1 → Micro-batch 1 → **Machine 1** → **W<sub>1</sub>** → Stage

2. **Stage** Micro-batch **W<sub>1</sub>**
  3. **W<sub>1</sub>**

1

1 device stage

1

- Device 1 123456789012345678901234567890
  - Device 1 123456789012345678901234567890



## Weight Stashing

device minibatch

1

- Device 1 □ 00005000 0000000000100000000100000000
  - Device 1 □ 00005000 0000000000100000000100000000

2 device stage

1



Vertical Sync

pipeline device device

1

- Device 1 5 1 1
  - Device 2 5 1 1

10 of 10

float16 float32 Master-Weight

1

1. **Master-Weight** **fp16**
  2. **inputs** **FP16**
  3. **outputs** **fp32**



## Loss Scaling

- fp16とfp32
- fp32とfp32
- fp16とfp32とfp32
- fp32**とfp32

FP16とFP16



### RNNとfp16とfp32

fp32

fp32とfp32



fp32



fp32

- fp32
- fp32とfp32



fp32

### MHAとMQAとGQA



## MHA

fp32

### MQAとMulti-Query Attention

fp32

fp32

- TransformerとkeyとvalueとKとVとQ
- KとV
- KV

fp32

- 
  - 

# GQA Grouped-Query Attention

7 of 7

三

- **Query** **Key**/**Value**
  - **N**=**1** **MQA** **N** **Query** **MHA**

三

-  **MHA** 

## MLA Multi-Head Latent Attention

Deepseek-V3

1

- KV Cache
  - MOA / GOA / KV Cache

1

- KV Cache
  - MOANGOA

1

- **Key****Value**KVKVKVKVKVKVKVKVKVKVKVKV
  - KVKVMLAOKVKVKV



10

Query Key \$ \begin{bmatrix} q\_t^R & q\_t^C \end{bmatrix} \$ \$ \begin{bmatrix} k\_t^R & k\_t^C \end{bmatrix} \$ \$ \begin{bmatrix} q\_t^C & k\_t^C \end{bmatrix} \$ RoPE \$ (q\_t^R, k\_t^R) \$

1. **Key-Value**\$h\_t\$
  2. **Key-Value**\$c\_t^{\{KV\}}\$c\_t^Q\$
    - \$c\_t^{\{KV\}}\$KV
    - \$c\_t^Q\$Q
  3. **KeyValue**
    - **Key** \$ Value \$k\_{\{t,i\}}^C, v\_{\{t,i\}}^C\$
    - **Value**\$k\_{\{t,i\}}^R\$R
  4. **Query**\$q\_{\{t,i\}}^R\$**Query**\$q\_{\{t,i\}}^C\$
  5. **QKV**



KV



$\text{W}^{\text{KV}} \text{c\_t}^{\text{KV}} \text{k\_t}^{\text{C}} \text{v\_t}^{\text{C}} = \text{W}^{\text{UK}} \text{c\_t}^{\text{KV}} \text{v\_t}^{\text{C}}$

███████████████████████████████████\$c\_t^{\{KV\}}\$█████████████████████**K|V**████

Q



██████ \$ c\_t^Q = W^{\{DQ\}} h\_t \backslash q\_t^C = W^{\{UQ\}} c\_t^Q \$ KV Cache ██████████

MLA|Q|K|RoPE||



$\$ \$ K_{rot} = R_n(W_{QK} \cdot c_t^{KV}) \$ \$$

-  \$W\_{\{QK\}}\$ 

\$\$ S = (W^{\{UQ\}})^{\wedge}T(c\\_t^{\wedge}Q)^{\wedge}TR\\_m^{\wedge}TR\\_nc\\_t^{\wedge}\{KV\}W^{\wedge}\{UK\} \$\$  
\$= (c\\_t^{\wedge}Q)^{\wedge}TR\\_m^{\wedge}T(W^{\{UQ\}})^{\wedge}TW^{\wedge}\{UK\}R\\_nc\\_t^{\wedge}\{KV\}) = (c\\_t^{\wedge}Q)^{\wedge}TR\\_m^{\wedge}TW\\_{merge}R\\_nc\\_t^{\wedge}\{KV\} \$\$\$\$ W\\_{\{QK\}}W

# MLA RoPE \$c\_t^{\{KV\}}\$



MFA | Multi-matrix Factorization Attention |

8 / 10

A horizontal row of fifteen empty square boxes, intended for children to write their names in, likely as part of a classroom activity.

ANSWER

## Paged Attention

1

A horizontal row of 24 small, identical rectangular boxes arranged in a single row.

GPU

- GPU/CPU
- 亂序/順序GPU

缓存行对齐

- 缓存行对齐  
◦ 乱序缓存行的大小为1024，顺序缓存行的大小为1024 token，4 token的大小为1024
- token对齐  
◦ 乱序token的大小为**1000 token**，5 token的大小为**10 token**，995 token的大小为**1000 token**
- 缓存行对齐  
◦ 乱序缓存行的大小为[500MB] [1GB] [300MB] [2GB]，顺序缓存行的大小为500MB，600MB，300MB
- 乱序缓存行的大小为500MB，600MB，300MB



LLM

- 预训练LLM



- 乱序token对齐token对齐
- 乱序token对齐Paged Attention对齐KV Cache



- 乱序缓存行对齐CPU
- KV Cache对齐Prompt对齐Paged Attention对齐KV Cache  
◦ 乱序缓存行对齐KV对齐
- 乱序KV Cache对齐Copy-on-Write对齐，3对齐5对齐

## Flash Attention

IO SRAM IO

GPU SRAM HBM Flash attention Attention HBM



闪存

- 乱序缓存行对齐HBM对齐
- 乱序缓存行对齐GPU对齐
- 乱序缓存行对齐SRAM对齐Pytorch对齐

## pytorch attention

- HBM Q K SRAM
- $S = QK^T$
- $S$  HBM
- HBM S HRAM
- $P = \text{softmax}(S)$
- $P$  HBM
- HBM P SRAM
- $O = PV$
- $O$  HBM
- $O$

## Compute Bound

- **Compute-Bound**
  - 矩阵乘法
  - GPU
- **Memory-Bound**
  - ReLU Softmax Sum Dropout
  - GPU

## Safe-Softmax

softmax \$\$ softmax(x\_i) = \frac{e^{x\_i}}{\sum\_{j=1}^n e^{x\_j}} \$\$ ~~max(x) \$x\_i\$~~  
Safe Softmax

Safe Softmax \$\$ softmax(x\_i) = \frac{e^{x\_i - \max(x)}}{\sum\_{j=1}^n e^{x\_j - \max(x)}} \$\$

- $\max(x) $x_i$$
- $e^{\max(x)} $e^{x_i - \max(x)}$

## Tiling Algorithm

- Tiling Algorithm
  - SRAM SRAM
  - SRAM SRAM Flash Attention SRAM
- Recomputation
  - GPU HBM
  - $S = QK^T$  Flash Attention HBM SRAM
  - $P = \text{softmax}(S)$  HBM
- Kernal Fusion
  - Q K V HBM
  - Kernal

## SRAM

- Q K V SRAM
- $S = QK^T$
- $P = \text{softmax}(S)$ 
  - softmax softmax

- Safe-Softmax
  - Safe-Softmax $\rightarrow$ Safe-Softmax $\rightarrow$ Safe-Softmax $\rightarrow$ Safe-Softmax
  - $x_1 = [1,2]$  $x_2 = [3,4]$
  - $m(x_1) = 2 = m(x)$
  - $f(x_1) = [e^{1-m(x_1)}, e^{2-m(x_1)}] = [e^{-1}, e^0]$
  - $l(x_1) = e^{-1} + e^0$
  - $m(x) = \max(m(x_1), m(x_2)) = 4$
  - $f(x_2) = [e^{3-m(x_2)}, e^{4-m(x_2)}] = [e^{-1}, e^0]$
  - $l(x_2) = e^{-1} + e^0$
  - $f(x) = l(x)$
- $$f(x) = [e^{m(x_1)-m(x)}f(x_1), e^{m(x_2)-m(x)}f(x_2)] \quad f(x) = [e^{-2}, (e^{-1}, e^0), e^0(e^{-1}, e^0)] \quad l(x) = e^{m(x_1)-m(x)} * l(x_1) + e^{m(x_2)-m(x)} * l(x_2) \quad l(x) = e^{-3} + e^{-2} + e^{-1} + e^0$$

- \$O=PV\$
- SRAM $\rightarrow$ \$O\rightarrow\$HBM

## GPU

◦ SIMD 亂序执行并行处理

GPU $\rightarrow$ Warp $\rightarrow$ 32 $\rightarrow$ Warp

SM $\rightarrow$ GPU

- CUDA Cores $\rightarrow$ CUDA 亂序执行并行处理
- Tensor Cores $\rightarrow$ Tensor 亂序执行并行处理
- Warp Scheduler $\rightarrow$ Warp 亂序执行并行处理
- Shared Memory 亂序执行并行处理
- L1 Cache 亂序执行并行处理
- Load/Store Units $\rightarrow$ HBM $\rightarrow$ Shared Memory 亂序执行并行处理
- Register File 亂序执行并行处理
- SFU 亂序执行并行处理



GPU $\rightarrow$ L1 Cache $\rightarrow$ SM $\rightarrow$ L2 Cache $\rightarrow$ SM $\rightarrow$ Shared Memory

- SRAM $\rightarrow$ L1 Cache $\rightarrow$ SM $\rightarrow$ L2 Cache $\rightarrow$ SM $\rightarrow$ Shared Memory
- HBM $\rightarrow$ SRAM $\rightarrow$ L1 Cache $\rightarrow$ SM $\rightarrow$ L2 Cache $\rightarrow$ SM $\rightarrow$ Shared Memory

◦ NVLink



- NVLink $\rightarrow$ GPU
- PCIe $\rightarrow$ CPU-GPU

GPU HBM

GPU HBM

1. GPU HBM
2. GPU HBM
3. GPU HBM
4. GPU HBM