

📌 重要概念

1. 模型
 - 模型训练与推理
 - 模型部署
2. **SFT** (Supervised Finetune)
 - 模型训练与推理
 - 模型部署
 - 模型训练与推理 LoRA Adapter 模型训练与推理 SFT 模型
3. **RLHF** (Reinforcement Learning)
 - 模型训练与推理
 - 模型部署
4. 模型
 - 模型训练与推理

SFT (Supervised Finetune)


模型训练与推理

模型训练与推理 dirty work 模型训练与推理

模型训练与推理

- **Few-Shot Prompting** 模型训练与推理 1-5 模型训练与推理
- **Seed Prompt** 模型训练与推理 task_type 模型训练与推理
- 模型
 - 模型训练与推理
 - answer 模型训练与推理
 - task_type 模型训练与推理 sft 模型训练与推理
- 模型
 - 模型 prompt
 - 模型 task_type 模型训练与推理 seed prompt 模型训练与推理 seed 模型训练与推理 pretrain 模型训练与推理
 - 模型 seed 模型训练与推理 prompt
 - 模型 answer
 - 模型 GPT4/Claude3
 - 模型 Qwen_72B/deepseek_MoE

模型训练与推理

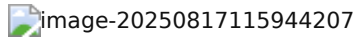
 image-20250817002125325

RL (Reinforcement Learning)

模型训练与推理

模型训练与推理

강화학습



- **Agent** 에이전트
- **action** 행동
- **Environment** 환경, reward
- **reward** 보상
- **State** 상태

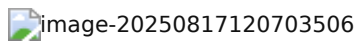
강화학습 MDP

정의

강화학습은 에이전트가 환경과 상호작용하여 보상을 최대화하는 것을 목표로 한다.

강화학습 MDP

1. M^2 : 상태 공간, 행동 공간, 전이 확률, 보상 함수
2. 에이전트
3. M^2 : 상태 공간, 행동 공간



강화학습 MDP의 수학적 표현

- 상태 공간 S , 행동 공간 A , 전이 확률 P , 보상 함수 R

$$\begin{matrix} S_1 & \text{to} & S_1 & \& S_1 & \text{to} & S_2 & \& S_1 & \text{to} & S_3 \\ S_2 & \text{to} & S_1 & \& S_2 & \text{to} & S_2 & \& S_2 & \text{to} & S_3 \\ S_3 & \text{to} & S_1 & \& S_3 & \text{to} & S_2 & \& S_3 & \text{to} & S_3 \end{matrix}$$

- 상태 s_t 에서 행동 a_t 을 취하면 다음 상태 s_{t+1} 와 보상 r_{t+1} 을 얻는다.
 $p(s_{t+1}|s_t, a_t)$
- 보상 함수 $R(s, a, s')$

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

- 가치 함수 $V(s)$
 - $V(s) = E[R_t | s_t = s]$

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

강화학습 MRP

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

- S : 상태 공간
- P : 전이 확률
- R : 보상 함수 $R_S = E[R_{t+1} | S_t = s]$
 - R_S : 상태 s 에 대한 기대 보상
 - $R_{ss'} = E[R_{t+1} | s_t = s, s_{t+1} = s']$
 - $R_{ss'} = E[R_{t+1} | s_t = s, a_t = a, s_{t+1} = s']$

- $\gamma \in [0,1]$

- Return G_t to the environment

- $G_t = R_{t+1} + \gamma G_t$
- $G_t = R_{t+1} + \gamma G_t$

$$G_t = R_{t+1} + \gamma G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

- $V(s)$ is the value of state s

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

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

MRPs

$$V(s) = E[G_t | S_t = s] = E[R_{t+1} + \underbrace{\gamma V(S_{t+1})}_{\text{value of next state}} | S_t = s] = E[R_{t+1} + \gamma V(S_{t+1}) | S_t = s]$$

MDP

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

- S is the set of states
- A is the set of actions
- P is the transition probability function $P_{ss'}^a = p[S_{t+1}=s' | S_t=s, A_t=a]$
- R is the reward function $R_s^a = E[R_{t+1} | S_t=s, A_t=a]$
 - R_s^a is the expected reward
 - $R_s^a = E[R_{t+1} | S_t=s, A_t=a]$
 - $R_s^a = E[R_{t+1} | S_t=s, A_t=a]$
- $\gamma \in [0,1]$

π

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

- $\pi(a|s)$ is the probability of taking action a in state s
- $\pi(a|s)$ is the probability of taking action a in state s
- MDP is a tuple $\langle \mathcal{M}, \pi \rangle$

$$\mathcal{M} = \langle S, A, P, R, \gamma \rangle, \pi$$


- $\pi(a|s)$ is the probability of taking action a in state s

$$P_{ss'}^a = \sum_{a \in A} \pi(a|s) P_{ss'}^a \quad R_s^a = \sum_{a \in A} \pi(a|s) R_s^a$$

- $\pi(a|s)$ is the probability of taking action a in state s
- $\pi(a|s)$ is the probability of taking action a in state s

π


□□□□□s□□□□□□ π □□□□□

 $V_{\pi}(s) = E_{\pi}[G_t|S_t=s]$ □□□□□ $G_t =$

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

□□□□□


□□□□□s□□□□□a□□□□□□ π □□□□□

 $q_{\pi}(s,a) = E_{\pi}[G_t|S_t = s, A_t = a]$ □□□□□□□□□□□□□□□□

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

□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□ $V_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s,a)$



□□□□□□□□□□□□□□□□□□□□□□ $q_{\pi}(s,a) = R_s + \gamma \sum_{s' \in S} P_{ss'} a V_{\pi}(s')$



□□□□□□□□

□□□□□□□ $V_{\pi}(s) = \max_{\pi} V_{\pi}(s)$ □□□□□□□ $q_{\pi}(s,a) = \max_{\pi} q_{\pi}(s,a)$ □□□□□

- $\pi^* \geq$ any π
-
-
-

$\pi^* \geq$ any π

□□□ $V_{\pi'}(s) \geq V_{\pi}(s)$ □□□ $\pi' > \pi$

□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□

□□□□ $V_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s,a)$ □□□□ $q_{\pi}(s,a) = R_s + \gamma \sum_{s' \in S} P_{ss'} a V_{\pi}(s')$ □□□□ $V_{\pi}(s) = \sum_{a \in A} \pi(a|s) (R_s + \gamma \sum_{s' \in S} P_{ss'} a V_{\pi}(s'))$ □□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□ $V_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s,a)$

□□□□ π^* □□□□□□□□□□□□□□□□ $q_{\pi^*}(s,a)$

- $\pi^* \geq V_{\pi^*}(s) \geq V_{\pi}(s)$
- $q_{\pi^*}(s,a) \geq q_{\pi}(s,a)$
- $V_{\pi}(s) \geq q_{\pi}(s,a)$

$$q_{\pi}(a \mid s) = \begin{cases} 1 & \text{if } a = \arg \max_a q_{\pi}(s, a), \\ 0 & \text{otherwise.} \end{cases}$$

$$V_*(s) = \max_a q_*(s, a)$$

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

$$V_*(s) = \max_a (R_s + \gamma \sum_{s'} P_{s, s'} a V_*(s'))$$
$$q_{\pi}(s, a) = R_s + \gamma \sum_{s' \in S} P_{ss'} a V_{\pi}(s')$$

$$V_{\pi}(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} V_{\pi}(s')$$

$$q_{\pi}(s, a) = R_s + \gamma \sum_{s' \in S} P_{ss'} a \max_{a'} q_{\pi}(s', a')$$

LLaMA

GPT LLaMA Transformer Decoder-only

-
-
-

RMSNorm

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

- x
- d
- ϵ
- γ

SwiGLU

Swish

$$\text{Swish}(x) = x \cdot \sigma(x)$$

- σ Sigmoid

$$\text{GLU}(x) = \sigma(W_1 x + b_1) \odot (W_2 x + b_2)$$

- \odot
- W_1, W_2, b_1, b_2

$$\text{SwiGLU}(x) = \text{Swish}(\text{Linear}_1(x)) \odot \text{Linear}_2(x)$$

- Swish
- GLU

RoPE

GQA

LLaMA2 3

GPT

Deepseek

MLA

- GPU
- GPU

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

All-Reduce

- GPU
- GPU
 - GPU SUM MAX
 - GPU

All-Gather

- GPU GPU
- GPU
 - GPU
 - GPU

DP Data Parallel

Python GIL CPU

image-20250814165126760

1. CPU端を主としてGPU端
2. GPU端を主としてGPU0端
3. GPU0端を主としてGPU端
4. GPU端を主としてGPU端

GPU

GPU端を主としてGPU端

- GPU0端を主としてGPU端
- GPU端を主としてGPU端

GPU

- GPU端を主としてGPU端
- GPU0端を主としてGPU端

DDP Distributed Data Parallel

GPU

- GPU端を主としてGPU端
- GPU0端を主としてGPU端

GPU

- GPU0端を主としてGPU端
- GPU端を主としてGPU端
- GPU端を主としてGPU端
- GPU端を主としてGPU端
- GPU端を主としてGPU端
- GPU端を主としてGPU端

GPU

GPU端を主としてGPU端

GPU端を主としてGPU端

- Scatter-Reduce端を主としてGPU端
- All-Gather端を主としてGPU端

GPU端を主としてGPU端

Ring-AllReduce

- GPU端を主としてGPU端
- GPU端を主としてGPU端
- GPU端を主としてGPU端

GPU

Scatter-Reduce

- GPU端を主としてGPU端
- GPU端を主としてGPU端

All-gather

- GPU端を主としてGPU端
- GPU端を主としてGPU端

- ψ_i GPU j $(j-i-1) \bmod n$ $(j-i-2) \bmod n$



FSDP(Fully Sharded Data Parallel)

- GPU
-
- GPU CPU

DeepSpeed ZeRO-1

3 GPU GPU



GPU ZeRO-1 GPU



-
- GPU0 GPU1 GPU2 GPU2 GPU
- GPU FP32 FP16
- FP16 GPU

ψ N

GPU

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

2ψ

DeepSpeed ZeRO-2

ZeRO-2 FP16 GPU



-
- GPU0 GPU1 GPU2 GPU
- GPU FP32 FP16
- FP16 GPU

$$\frac{\psi(N)}{N} \approx \psi(N)$$

GPU

- $$\frac{\psi(N)}{N} \approx \psi(N)$$
- $$\frac{\psi(N)}{N} \approx \psi(N)$$

$$\frac{\psi(N)}{N} \approx \psi(N)$$

DeepSpeed ZeRO-3

ZeRO-3 FP16



- GPU
- GPU
- GPU
- GPU

$$\frac{\psi(N)}{N} \approx \psi(N)$$

GPU

- $$\frac{\psi(N)}{N} \approx \psi(N)$$
- $$\frac{\psi(N)}{N} \approx \psi(N)$$

$$\frac{\psi(N)}{N} \approx \psi(N)$$

TP

GPU

$$Y = XW$$

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

$$\begin{bmatrix} y_1 & y_2 & y_3 & y_4 \end{bmatrix} = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 \\ w_1 & w_2 & w_3 & w_4 \end{bmatrix}$$

$$W = \begin{bmatrix} W_0 & W_1 \\ W_2 & W_3 \end{bmatrix}$$

$$Y_0 = XW_0 = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 \\ w_1 & w_2 & w_3 & w_4 \end{bmatrix}$$

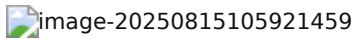
$$Y_1 = XW_1 = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 \\ w_1 & w_2 & w_3 & w_4 \end{bmatrix}$$

$$Y = \begin{bmatrix} Y_0 & Y_1 \end{bmatrix}$$

□□□□□□□□

$$Y = XW$$

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

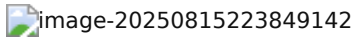
$$\begin{aligned} \begin{pmatrix} y_1 & y_2 & y_3 & y_4 \end{pmatrix} &= \begin{pmatrix} x_1 & x_2 & x_3 & x_4 \end{pmatrix} \begin{pmatrix} w_1 & w_2 & w_3 & w_4 \end{pmatrix} \\ W &= W_0 W_1 \\ X &= \begin{pmatrix} X_1 & X_2 \end{pmatrix} \quad Y_1 = \begin{pmatrix} x_1 & x_3 \end{pmatrix} \begin{pmatrix}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


□□□□□**PP**□

□□□□□□□□□□□□□□□□□□□□**GPU**□

□□□□□□

- Stage Stage GPU
-



1111

- GPU2 is not available
- GPU2 is not available
- GPU1 is not available

5/5

- [illegible]

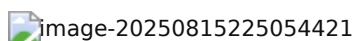
□□□□□□

F-then-B

- `tf.nn.conv2d`
- `tf.nn.conv2d` mini-batch `tf.nn.conv2d` mini-batch

1F1B

- 00000000000000000000000000000000
- 0000000000 stage4 0 F42 stage4 00 2 0 micro-batch 0000000000 F42 000000 F41 0000 B41 stage4 00 1 0 micro-batch 00000000000000000000000000000000 F41 00000000 **F42** 0000 **F41** 00000000



- | | | | | | | | |
|--|--|--|--|--|--|--|--|
| | | | | | | | |
|--|--|--|--|--|--|--|--|

Gpipe

□□□□□□□□□□□□□□□□□□□□**F-then-B**□□



□□□□□□Mini-bacth□□□□□□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 vs micro-batch
 - Pipeline Flush
 -
 -
 -
- mini-batch vs micro-batch

PipeDream□□□□□□

- micro-batch mirco-batch
- machine1 micro-batch forward machine2
- machine4 machine1 **1F1B**
- Bubble GPU



 image-20250815235148856

103

- 5Micro-batchMachine1Micro-batch 1FW2-4
 - **Micro-batch 5** **Micro-batch 1** **Micro-batch 5** **Micro-batch 4**
- Machine2FW5batch1-2Machine1batch1Machine1Machine2

5/5

[illegible]

- **Weight sharding**
 - 1. **Machine 1** 用 **Micro-batch 5** 个数据 计算 W_5
 2. **Machine 1** 用 **Micro-batch 5** 个数据 计算 W_5
 3. 所有 **Machine** 计算完 W_5 后 再求和
- **Vertical Sync**
 - 1. **Micro-batch** 用 **Machine 1** 个数据 计算 W_1 后 再求和

1. fp16
2. fp32
3. fp16fp32
4. fp32

fp16FP16

image-20250813121114207

RNNfp16fp32

image-20250813115950009

image-20250813115500851

- 1.
- 2.

image-20250813115927153

MHAMQAGQA

image-20250816120438735

MHA

MQAMulti-Query Attention

- TransformerkeyvalueKVQ
- KV
- KV

- □□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□

GQA Grouped-Query Attention

□□□□□□

111

- **Query** **Key** **Value**
- $N \times N = 1$ **MQA** N **Query** **MHA**

111

- **MHA**

MLA Multi-Head Latent Attention

Deepseek-V3

555

- KV Cache
- MQA GQA KV Cache

□□□□□□

- 降低KV Cache的内存占用
- 支持MQA和GQA的并行化

□□□□□

- **Key Value** $\text{KV} \rightarrow \text{KV}$
- $\text{KV} \rightarrow \text{MLA} \rightarrow \text{KV}$



 image-20250816120905940

□□□□□□□□

$$\text{QueryKey} \begin{bmatrix} q_t^R & q_t^C \end{bmatrix} \begin{bmatrix} k_t^R & k_t^C \end{bmatrix} (q_t^C, k_t^C) \text{RoPE} (q_t^R, k_t^R)$$

1. hash_t
2. c_t^{KV}
 - c_t^{KV}
 - c_t^{Q}
3. $\text{Key} \oplus \text{Value}$
 - $\text{Key} \oplus \text{Value} \oplus k_{t,i}^{\text{C}}, v_{t,i}^{\text{C}}$
 - $\text{Key} \oplus \text{Value} \oplus k_{t,i}^{\text{R}}$
4. $\text{Query} \oplus q_{t,i}^{\text{R}}$
5. $\text{Q} \oplus \text{V}$



 image-20250816131927886

KV



计算 $c_t^{KV} k_t^C v_t^C$ 的公式为 $c_t^{KV} = W^{DKV} h_t \cdot k_t^C = W^{UK} c_t^{KV} \cdot v_t^C = W^{UV} c_t^{KV}$ 的公式为 $c_t^{KV} = W^{UV} c_t^{KV}$

计算 c_t^{KV} 的公式为 $c_t^{KV} = W^{UV} c_t^{KV}$

计算 Q 的公式为 $Q = W^{DQ} h_t$



计算 $c_t^Q = W^{DQ} h_t$ 的公式为 $c_t^Q = W^{UQ} c_t^Q$ 的公式为 $c_t^Q = W^{UQ} c_t^Q$

计算 Q 的公式为 $Q = W^{DQ} h_t$

计算 k_t^R 的公式为 $k_t^R = W^{DR} h_t$

- 计算 Q 的公式为 $Q = W^{DQ} h_t$
- 计算 c_t^{KV} 的公式为 $c_t^{KV} = W^{UV} c_t^{KV}$

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

- 计算 W_{QK} 的公式为 $W_{QK} = W_{QK}$

计算 $S = (W^{UQ} c_t^Q)^T (W^{UR} c_t^R)$ 的公式为 $S = (W^{UQ} c_t^Q)^T (W^{UR} c_t^R)$

计算 Q 的公式为 $Q = W^{DQ} h_t$

- 计算 k_t^C 的公式为 $k_t^C = W^{DK} h_t$
- 计算 Q 的公式为 $Q = W^{DQ} h_t$
- 计算 K 的公式为 $K = W^{DK} h_t$

计算 Q 的公式为 $Q = W^{DQ} h_t$

计算 Q 的公式为 $Q = W^{DQ} h_t$

计算 Q 的公式为 $Q = W^{DQ} h_t$

计算 Q 的公式为 $Q = W^{DQ} h_t$

计算 Q 的公式为 $Q = W^{DQ} h_t$

计算 Q 的公式为 $Q = W^{DQ} h_t$


计算 Q 的公式为 $Q = W^{DQ} h_t$

📄

- GPU vs CPU 推理性能对比
- 推理/训练 GPU 显存需求对比


📄

- 推理性能对比
 - 推理性能对比 推理 吞吐量 1024 token 推理 1024 token 推理 4 token 推理
- 推理性能对比 token 推理性能对比
 - 推理性能对比 1000 token 推理 5 token 推理 10 token 推理 995 token 推理
- 推理性能对比
 - 推理性能对比 推理性能对比 [推理 500MB] [推理 1GB] [推理 300MB] [推理 2GB] 推理性能对比 500MB 推理性能对比 600MB 推理性能对比
 - 推理性能对比 推理性能对比 推理性能对比 推理性能对比


 image-20250819121745375

📄

- 推理性能对比 vLLM 推理性能对比

 image-20250819121817645

- 推理性能对比 token 推理性能对比 token 推理性能对比
- 推理性能对比 推理性能对比 Paged Attention 推理性能对比 KV Cache


 image-20250819121841203

- 推理性能对比 CPU 推理性能对比
- KV Cache 推理性能对比 Prompt 推理性能对比 Paged Attention 推理性能对比 KV Cache 推理性能对比
 - 推理性能对比 推理性能对比 KV 推理性能对比
 - 推理性能对比 Copy-on-Write 推理性能对比 推理性能对比 3 推理性能对比 5 推理性能对比

Flash Attention

推理性能对比 IO 推理性能对比 SRAM 推理性能对比 IO

推理性能对比 GPU 推理性能对比 SRAM 推理性能对比 HBM 推理性能对比 Flash attention 推理性能对比 Attention 推理性能对比 HBM

 image-20250820134630626

📄

- 推理性能对比 HBM 推理性能对比
- 推理性能对比 GPU 推理性能对比
- 推理性能对比 推理性能对比 SRAM 推理性能对比 Pytorch 推理性能对比

pytorch attention

- HBM QK 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**
 -
 - GPU
- **Memory-Bound**
 - ReLU, Softmax, Sum, Dropout
 - GPU

Safe-Softmax

softmax $\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$ x_i

Safe Softmax $\text{softmax}(x_i) = \frac{e^{x_i - \max(x)}}{\sum_{j=1}^n e^{x_j - \max(x)}}$

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

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

SRAM

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

- Safe-Softmax
- Safe-Softmax Safe-Softmax

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

$$f(x) = [e^{\{m(x_1)-m(x)\}}f(x_1), e^{\{m(x_2)-m(x)\}}f(x_2)] \setminus f(x) = [e^{\{-2\}}(e^{\{-1\}}, e^{\{0\}}), e^{\{0\}}(e^{\{-1\}}, e^{\{0\}})] \setminus l(x) = e^{\{m(x_1)-m(x)\}} \cdot l(x_1) + e^{\{m(x_2)-m(x)\}} \cdot l(x_2) \setminus l(x) = e^{\{-3\}} + e^{\{-2\}} + e^{\{-1\}} + e^{\{0\}}$$

- $O = PV$
- $\text{SRAM} \cdot O \cdot \text{HBM}$

GPU

SIMT SIMT

GPU Warp 32

SM GPU


- **CUDA Cores** CUDA Cores
- **Tensor Cores** Tensor Cores
- **Warp Scheduler** Warp Warp
- **Shared Memory** Shared Memory
- **L1 Cache** L1 Cache
- **Load/Store Units** HBM Shared Memory
- **Register File** Register File
- **SFU** SFU

 image-20250819141437101

GPU

- **SRAM** L1 Cache SM L2 Cache SM Shared Memory
- **HBM** HBM

GPU

 image-20250819132629103

- **NVLink** GPU
- **PCIe** CPU-GPU

GPU 的 HBM 是通過 PCIe 與 CPU 共享的 GPU HBM

GPU 的 HBM 是通過 PCIe 與 CPU 共享的 GPU HBM

1. GPU 的 HBM 是通過 PCIe 與 CPU 共享的 GPU HBM
2. GPU 的 HBM 是通過 PCIe 與 CPU 共享的 GPU HBM
 - GPU 的 HBM 是通過 PCIe 與 CPU 共享的 GPU HBM
3. GPU 的 HBM 是通過 PCIe 與 CPU 共享的 GPU HBM
4. GPU 的 HBM 是通過 PCIe 與 CPU 共享的 GPU HBM