

Optimizing Efficiency in Large-Scale Transformer Models

Haitian Jiang 2025.02.26

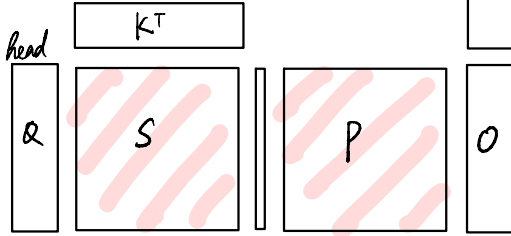
Outline :

- Memory and Computation Challenges in Transformers (background)
 - Memory challenge:
 - Training:
 - Long context makes attention hard and slow (for FlashAttn)
 - Parameter size (for Mixed Precision)
 - Inference: KV cache and its memory burden (for MLA)
 - Computation challenge:
 - Heavy MLP and the MoE basics (for DeepSeek MoE)
- FlashAttention for Long Context Training
 - GPU Architecture Basics
 - FlashAttention 2 Algorithm
- DeepSeek V3 as a Bag of Optimizations
 - Mixed Precision Training
 - Precision and range in FP, different FP representations
 - Why it's fast and (memory) efficient
 - FP16 and BF16 mixed precision training
 - DeepSeek FP8 quantized mixed precision training
 - Multi-Latent Attention
 - MHA, GQA, MQA, MLA
 - DeepSeek MoE

DO NOT DISTRIBUTE

Motivation

Attention: $\text{softmax}(QK^T)V$



Model Name	n_{params}	n_{layers}	d_{model}	n_{heads}	d_{head}
GPT-3 Small	125M	12	768	12	64
GPT-3 Medium	350M	24	1024	16	64
GPT-3 Large	760M	24	1536	16	96
GPT-3 XL	1.3B	24	2048	24	128
GPT-3 2.7B	2.7B	32	2560	32	80
GPT-3 6.7B	6.7B	32	4096	32	128
GPT-3 13B	13.0B	40	5140	40	128
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128

• Size of attn score :

$$4K \times 4K \times 96 \times 96 \times 4B = 576 \text{ GiB}$$

seq seq n_{head} n_{layer} $fp32$

How to store ? So much to r/w ?

• Parameter size :

$$175 \times 10^9 \times 4B \approx 700 \text{ GiB}$$

n_{params} $fp32$

• KV cache size: (inference)

$$2 \times 1 \times 4K \times 96 \times 12K \times 4B = 36 \text{ GiB}$$

KV $batch$ seq n_{layer} d_{model} $fp32$

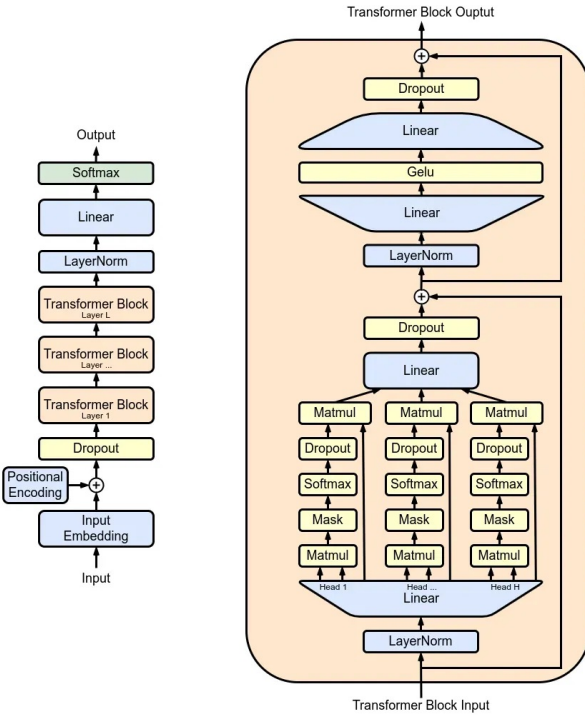
For 1 query !

• Computation

$$\text{Attn} : 4b \cdot seq \cdot d^2 + 2b \cdot seq^2 \cdot d$$

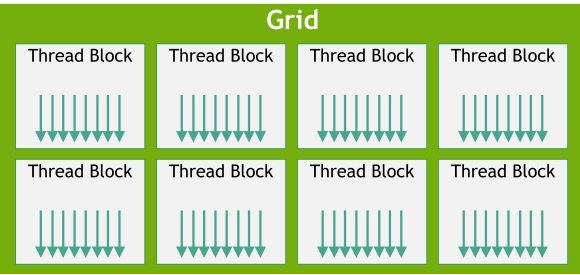
$$\text{MLP} : 8b \cdot seq \cdot d^2$$

$$seq < d, \text{ flops (MLP)} > \text{ flops (Attn)}$$

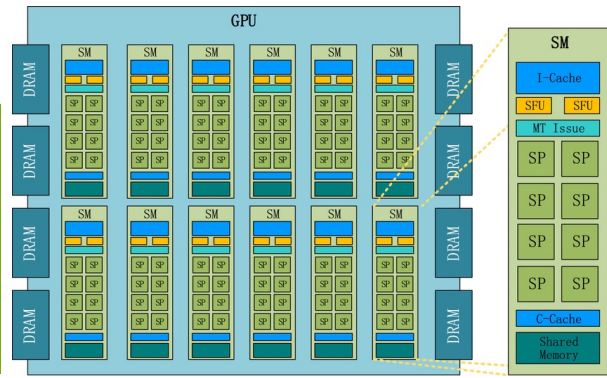


A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
FFN expand ratio	Seq len	Params w/o embed	N layers	Embed dim	Params w embed	Matmul %	Attn %	Matmul total	Attn total	FFN flops	Attn matmul flops	Attn O ²	Attn A ^V	Attn softmax
4	1024	84934656	12	768	123.6M	81.81%	18.19%	2087354105856	464158457856	115964116992	57982058496	19327352832	19327352832	25165824
4	1024	301989888	24	1024	353.6M	85.71%	14.29%	14843406974972475109122048	412316860416	206158430208	51539607552	51539607552	51539607552	50331648
4	1024	707788000	36	1280	772.1M	88.23%	11.77%	52183852646406960564928512	966367641600	483183820800	96636764160	96636764160	96636764160	75497472
4	1024	1474560000	48	1600	1.6B	90.36%	9.84%	1449551462400	15466714103802013269520000	1006632960000	161061273600	161061273600	161061273600	100663296
4	1024	1207959552	24	2048	1.3B	94.12%	5.88%	59373627899903712059703296	1649267441664	824633720832	51539607552	103079215104	50331648	
4	1024	2516582400	32	2560	2.6B	95.24%	4.76%	16492674416648248484691968	3435973836800	1717986918400	85899345920	171798691840	67108864	
4	1024	6442450944	32	4096	6.6B	96.97%	3.03%	4222124650689	13196237016968796093022208	439804651104	137438953472	274877906944	67108864	
4	1024	12681408000	40	5140	12.9B	97.57%	2.43%	1.03886E+15	2587382251520	17313430950000	8657174528000	215587225600	431174451200	83886080
4	1024	173946175488	96	12288	174.6B	98.97%	1.03%	3.41992E+16	3562610947522	2374945115996	11874725799808	12369950581248	2473901162496	201326592

GPU Architecture

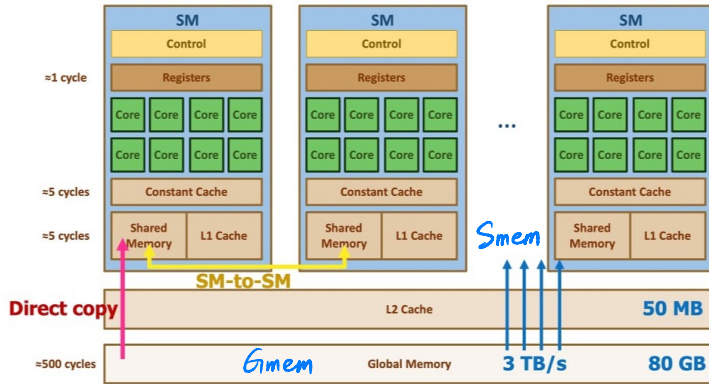


software view

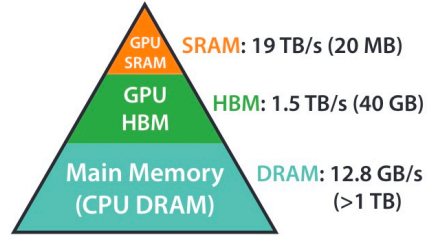


hardware view

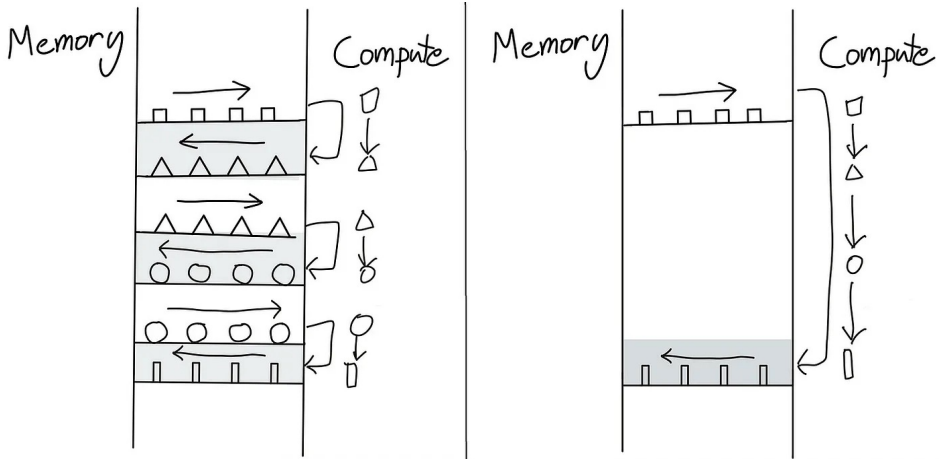
Memory in the H100 GPU Architecture



memory hierarchy



Memory Hierarchy with Bandwidth & Memory Size



Operator Fusion Simplified

FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness

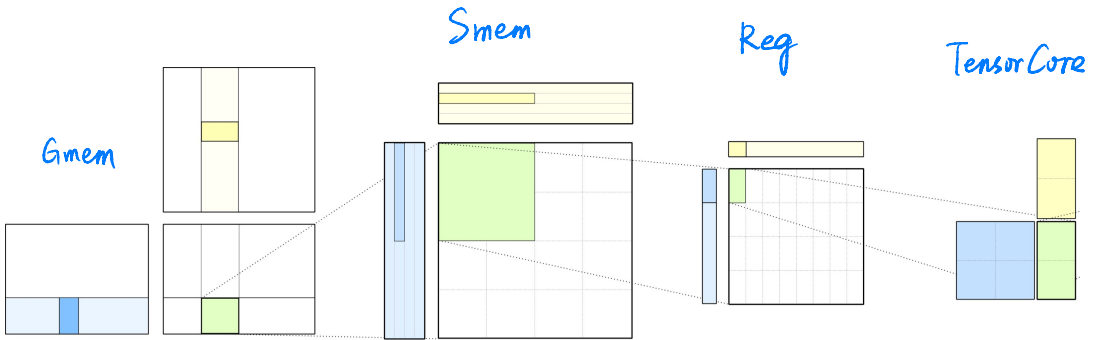
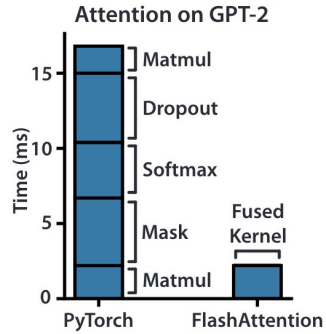
Which takes longer time? Computation? Communication?
(Memory loading)

Algorithm 0 Standard Attention Implementation

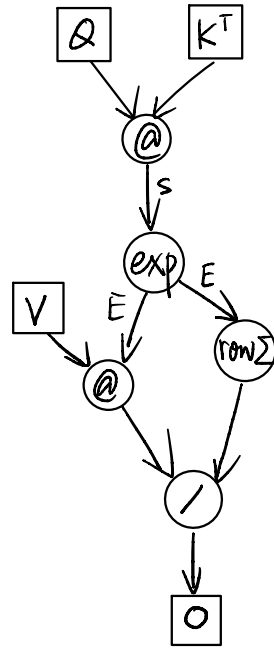
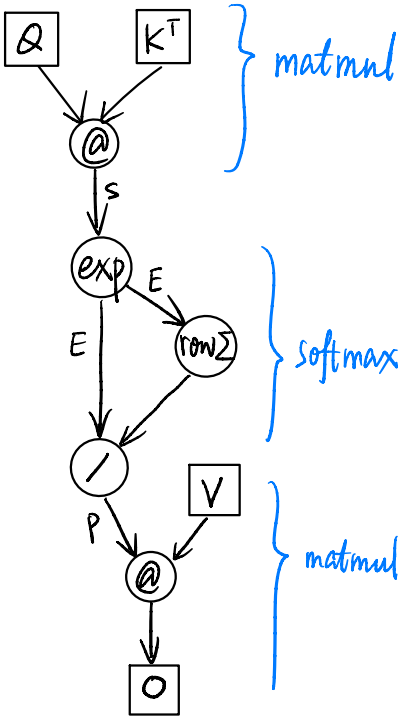
Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM.

- 1: Load \mathbf{Q}, \mathbf{K} by blocks from HBM, compute $\mathbf{S} = \mathbf{Q}\mathbf{K}^T$, write \mathbf{S} to HBM.
- 2: Read \mathbf{S} from HBM, compute $\mathbf{P} = \text{softmax}(\mathbf{S})$, write \mathbf{P} to HBM.
- 3: Load \mathbf{P} and \mathbf{V} by blocks from HBM, compute $\mathbf{O} = \mathbf{P}\mathbf{V}$, write \mathbf{O} to HBM.
- 4: Return \mathbf{O} .

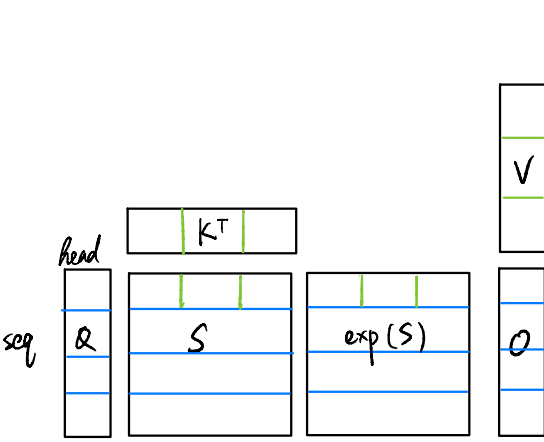
$S = Q @ K^T$: load Q, K from Gmem, write S to Gmem
 $E = \exp(S)$: load S from Gmem, write E to Gmem
 $\gamma = \text{rowsum}(E)$: load E from Gmem, writing γ negligible
 $P = E / \gamma$: load E from Gmem, write P to Gmem
 $O = P @ V$: load P, V from Gmem, write O to Gmem



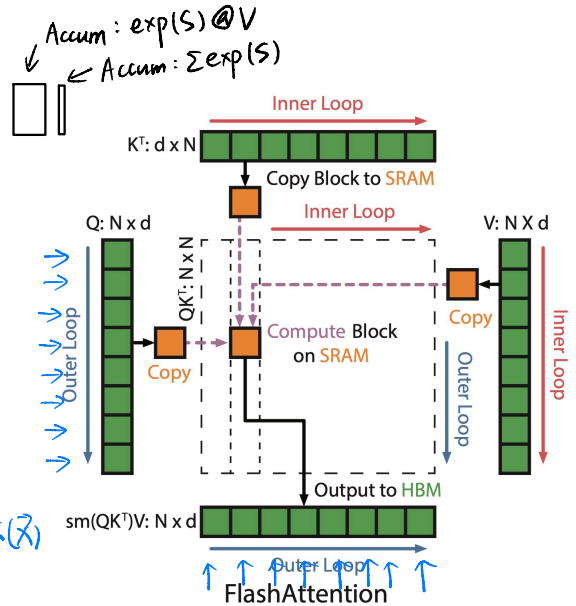
$$(v_1, v_2, \dots, v_n) \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} = \sum_i v_i \cdot u_i$$



$$\left[\exp(QK^T) / \sum \exp(QK^T) \right] @ V \rightarrow \left[\exp(QK^T) @ V \right] / \sum \exp(QK^T)$$



$$\text{softmax}(\vec{x}) = \frac{e^{x-m}}{\sum e^{x-m}} \quad m = \max(\vec{x})$$



Algorithm 1 FLASHATTENTION-2 forward pass

Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM, block sizes B_c, B_r .

- 1: Divide \mathbf{Q} into $T_r = \left\lceil \frac{N}{B_r} \right\rceil$ blocks $\mathbf{Q}_1, \dots, \mathbf{Q}_{T_r}$ of size $B_r \times d$ each, and divide \mathbf{K}, \mathbf{V} into $T_c = \left\lceil \frac{N}{B_c} \right\rceil$ blocks $\mathbf{K}_1, \dots, \mathbf{K}_{T_c}$ and $\mathbf{V}_1, \dots, \mathbf{V}_{T_c}$, of size $B_c \times d$ each.
- 2: Divide the output $\mathbf{O} \in \mathbb{R}^{N \times d}$ into T_r blocks $\mathbf{O}_1, \dots, \mathbf{O}_{T_r}$ of size $B_r \times d$ each, and divide the logsumexp L into T_r blocks L_1, \dots, L_{T_r} of size B_r each.
- 3: **for** $1 \leq i \leq T_r$ **do** *outer loop: parallel*
- 4: Load \mathbf{Q}_i from HBM to on-chip SRAM.
- 5: On chip, initialize $\mathbf{O}_i^{(0)} = (\mathbf{0})_{B_r \times d} \in \mathbb{R}^{B_r \times d}$, $\ell_i^{(0)} = (\mathbf{0})_{B_r} \in \mathbb{R}^{B_r}$, $m_i^{(0)} = (-\infty)_{B_r} \in \mathbb{R}^{B_r}$.
- 6: **for** $1 \leq j \leq T_c$ **do** *inner loop: sequential*
- 7: Load $\mathbf{K}_j, \mathbf{V}_j$ from HBM to on-chip SRAM.
- 8: On chip, compute $\mathbf{S}_i^{(j)} = \mathbf{Q}_i \mathbf{K}_j^T \in \mathbb{R}^{B_r \times B_c}$.
- 9: On chip, compute $m_i^{(j)} = \max(m_i^{(j-1)}, \text{rowmax}(\mathbf{S}_i^{(j)})) \in \mathbb{R}^{B_r}$, $\tilde{\mathbf{P}}_i^{(j)} = \exp(\mathbf{S}_i^{(j)} - m_i^{(j)}) \in \mathbb{R}^{B_r \times B_c}$ (pointwise), $\ell_i^{(j)} = e^{m_i^{(j-1)} - m_i^{(j)}} \ell_i^{(j-1)} + \text{rowsum}(\tilde{\mathbf{P}}_i^{(j)}) \in \mathbb{R}^{B_r}$.
- 10: On chip, compute $\mathbf{O}_i^{(j)} = \text{diag}(e^{m_i^{(j-1)} - m_i^{(j)}})^{-1} \mathbf{O}_i^{(j-1)} + \tilde{\mathbf{P}}_i^{(j)} \mathbf{V}_j$.
- 11: **end for**
- 12: On chip, compute $\mathbf{O}_i = \text{diag}(\ell_i^{(T_c)})^{-1} \mathbf{O}_i^{(T_c)}$.
- 13: On chip, compute $L_i = m_i^{(T_c)} + \log(\ell_i^{(T_c)})$.
- 14: Write \mathbf{O}_i to HBM as the i -th block of \mathbf{O} .
- 15: Write L_i to HBM as the i -th block of L .
- 16: **end for**
- 17: Return the output \mathbf{O} and the logsumexp L .

Complexity	Normal Attn	Flash Attn	omitting batch size,
Computation	$O(N^2 d)$	$O(N^2 d)$	$N = \text{sequence length}$
GPU memory	$O(N^2 d)$	$O(Nd)$	$d = \text{head dim}$
I/O (access to Gmem)	$O(Nd + N^2)$	$O\left(\frac{N^2 d^2}{M}\right)$	$M = \text{Smem size}$
			$N \gg d, M \gg d^2$

Skipped: safe softmax (numerical stability)

online softmax (FlashAttn 1, FA2 reorder computation)

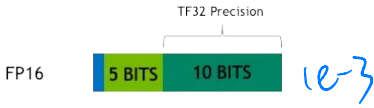
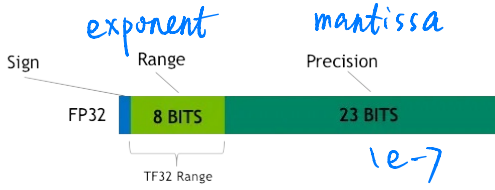
backward computation

mask, dropout

hardware utilization techniques

Mixed-Precision Training

$$f = (-1)^s \times 1.m \times 2^{e - \text{bias}}$$



FP8 e4m3 better in precision

FP8 e5m2 larger range

16 bit: w g

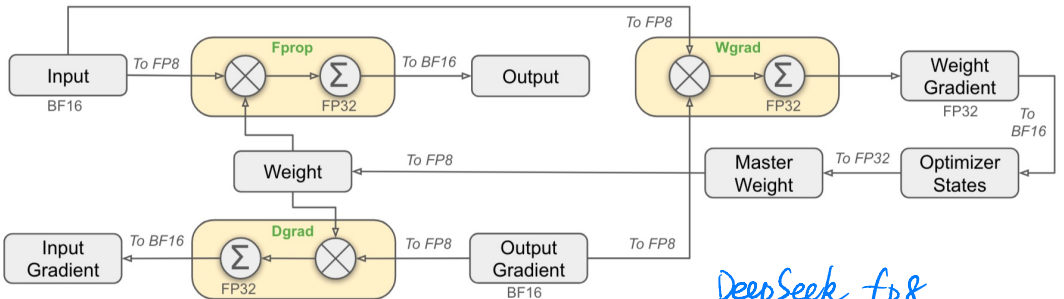
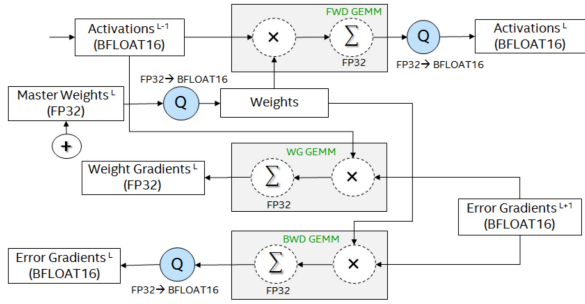
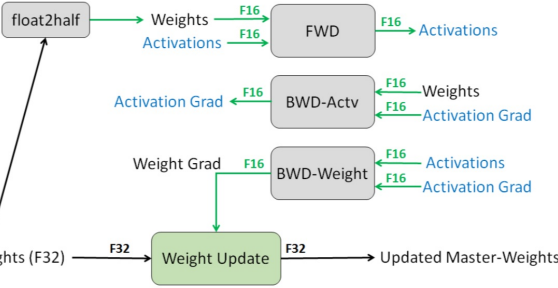
32 bit: w g \bar{g} \bar{g}^2

Why?

	Hopper	Ampere	Turing	Volta
Supported Tensor Core precisions	FP64, TF32, bfloat16, FP16, FP8, INT8	FP64, TF32, bfloat16, FP16, INT8, INT4, INT1	FP16, INT8, INT4, INT1	FP16
Supported CUDA Core precisions	FP64, FP32, FP16, bfloat16, INT8	FP64, FP32, FP16, bfloat16, INT8	FP64, FP32, FP16, INT8	FP64, FP32, FP16, INT8

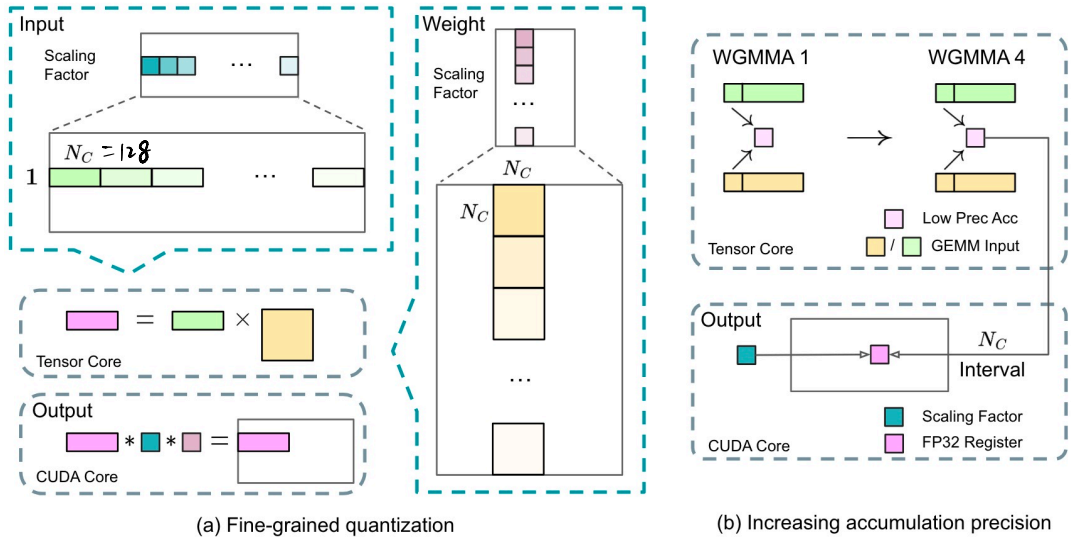
f_{p16} : grad scaler (prevent underflow)

b_{f16}



DeepSeek fp8

DeepSeek FP8 Quantization



Transformer Engine FP8 (from NVIDIA)

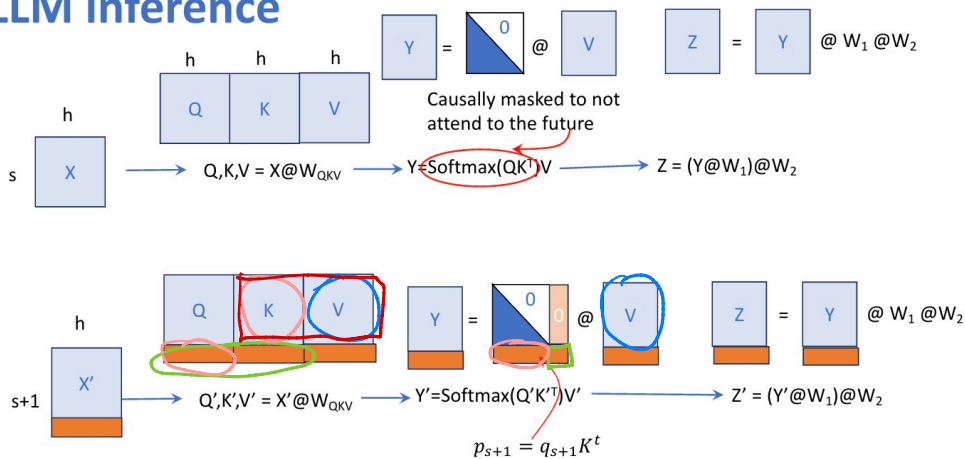
hybrid, no promotion

original precision: (bf(b/fp32))

- embedding & lm head
- MoE gating module
- normalizations
- attention

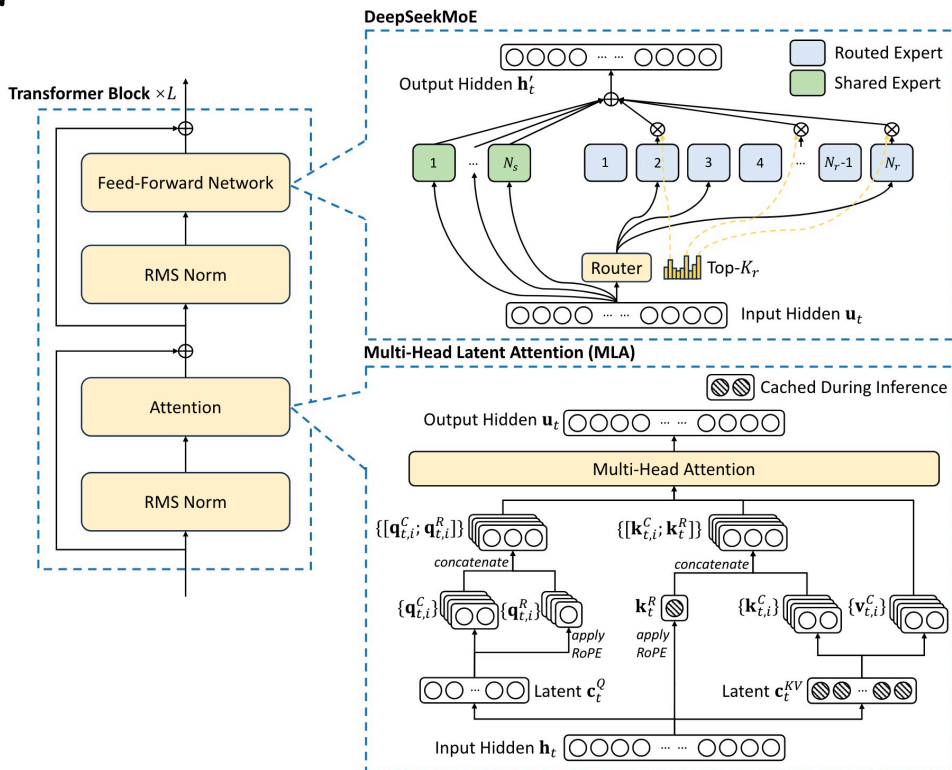
KV
Cache

LLM inference




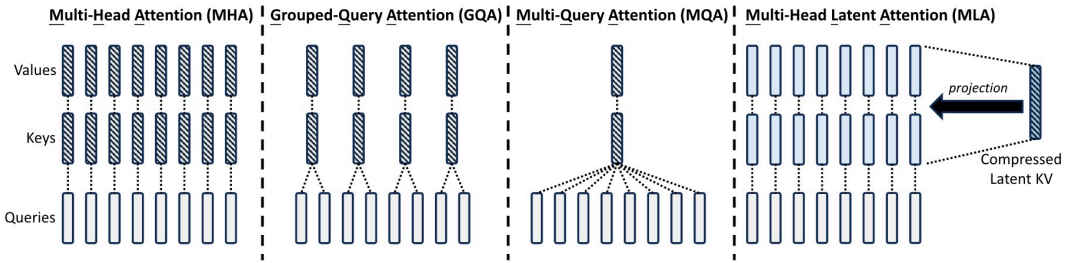
What to save across successive generation iterations for the same request?

DeepSeek Architecture



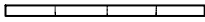
Attention: Multi-Latent Attention

 Cached During Inference



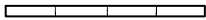
MHA

h_t



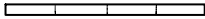
GQA

h_t



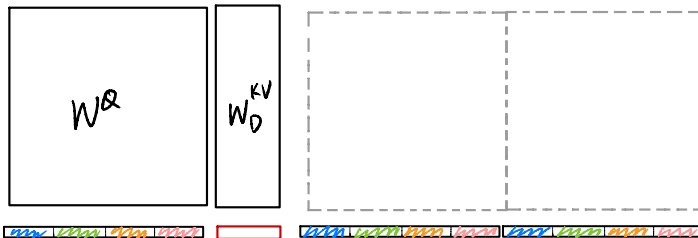
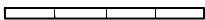
MQA

h_t

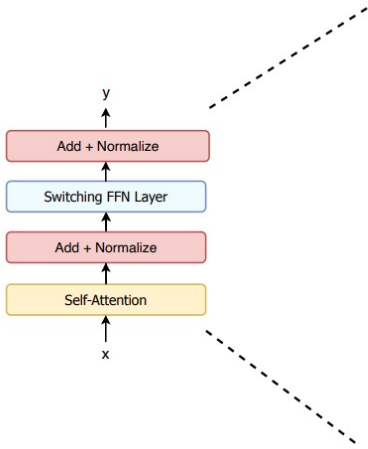


MLA

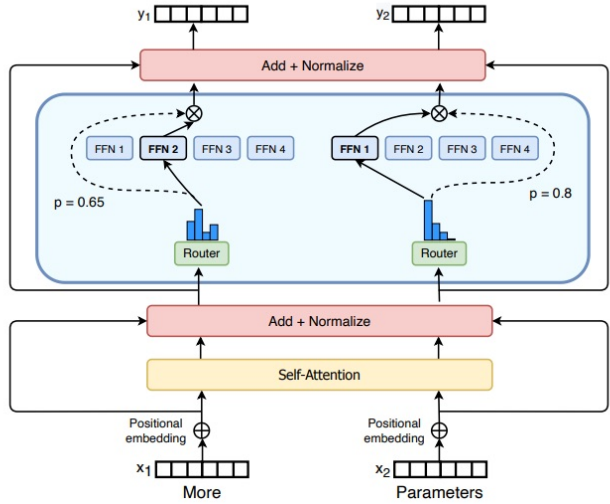
h_t



MLP: Mixture of Expert

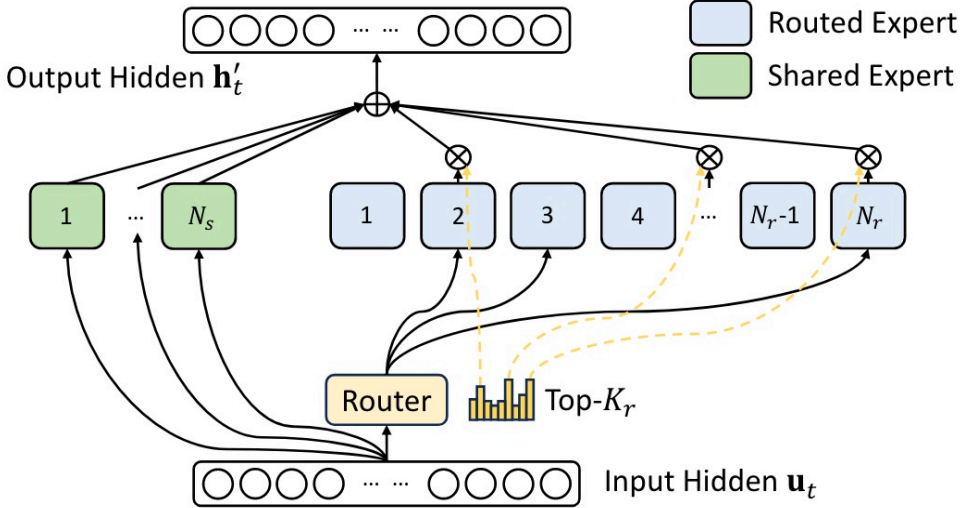


MoE from switch transformer



DeepSeek MoE

671 B v.s. 37B



Differences

- More experts
 - Shared experts
- } \Rightarrow outperform GShard by lot

routing: reducing cross-device communication

naive: top-K experts among all experts

deepseek: each token: select $M \geq 3$ devices (EP),
top-K experts among these M devices

Load balancing: routing collapse & efficiency

V2: expert, device, communication loss; token drop

V3: adding bias to affinity score to control; no token drop

training paradigm: MTP (V3)

