

DeepSeek Architecture

Changshi Li

Wuhan University

2025.3.6



Contents

1 Architecture

2 Loss Functions

3 Training Framework

4 FP8 Training



Math format

Let $X \in \{0, 1\}^{T \times d}$ for input, where d is the dimension of tokens and T is the length of inputs. Then the **embedding layer** is defined as:

$$\text{Enc}(X) = XE, \quad (1)$$

where $E \in \mathbb{R}^{d \times D}$.

Remark: $\text{Enc}(X) : \{0, 1\}^{T \times d} \rightarrow \mathbb{R}^{T \times D}$.

Let $X \in \mathbb{R}^{T \times D}$ be the hidden input and $X_{i \cdot}^T \in \mathbb{R}^D$, where D is the dimension of middle layer. Then **RMSNorm** layer is defined as:

$$\text{RMSNorm}(X_{i \cdot}) := \frac{X_{i \cdot}}{\sqrt{\frac{1}{D} \|X_{i \cdot}\|_2^2}}, \forall i \in \{0, 1, \dots, T\} \quad (2)$$

Remark: $\text{RMSNorm}(X) : \mathbb{R}^{T \times D} \rightarrow \mathbb{R}^{T \times D}$.



Math format

Let $X \in \mathbb{R}^{T \times D}$ be the hidden input, then the Rotary Positional Embedding (**RoPE**) is defined as:

$$\text{RoPE}(X; W_p, W_R) = XW_p W_R^T, \quad (3)$$

where

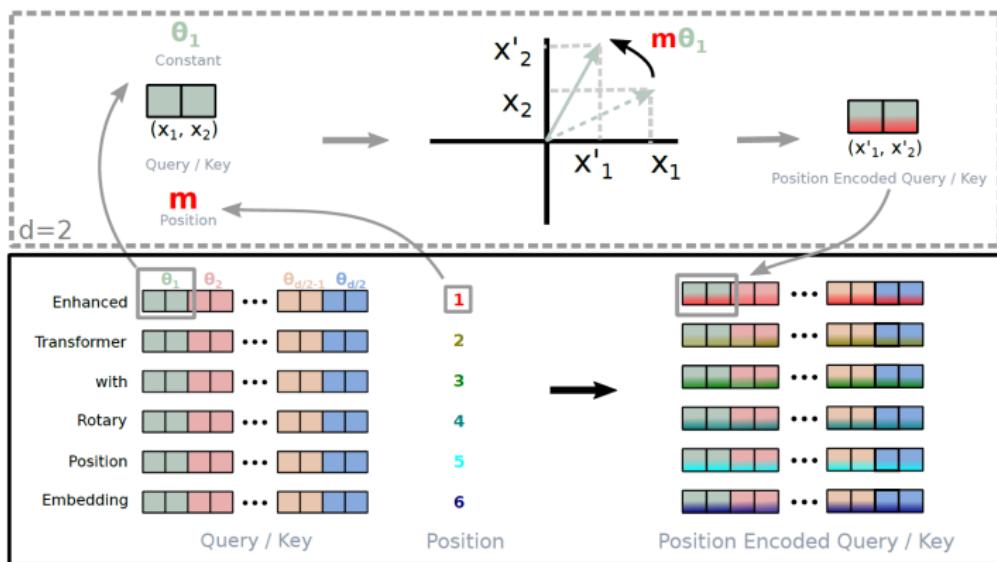
$$W_R = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{D^R/2} & -\sin m\theta_{D^R/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{D^R/2} & \cos m\theta_{D^R/2} \end{pmatrix}$$

, $\theta_i = 10000^{-2(i-1)/d}$, $i \in [1, 2, \dots, D^R/2]$ and $W_p \in \mathbb{R}^{D \times D^R}$.

Remark: $\text{RoPE}(X) : T \times D \rightarrow T \times D^R$.



RoPE



Math format

Let $X \in \mathbb{R}^{T \times D}$ be the hidden input, then the **LoRA** layer is defined as:

$$\text{LoRA}(X; W_{down}, W_{up}) := XW_{down}W_{up}, \quad (4)$$

where $W_{down} \in \mathbb{R}^{D \times d_{down}}$, $W_{up} \in \mathbb{R}^{d_{down} \times d_h h_n}$, $d_{down} = d_c$ for keys and values, $d_{down} = d'_c$ for queries, d_h is the hidden dimension before attention layer and h_n is the number of heads in multi-head attention layer.

Remark: $\text{LoRA}(X; W_{down}, W_{up}) : T \times D \rightarrow T \times d_h h_n$.

Let $L \in T \times d_h h_n$, Define an extract operator as follows:

$$L_{..h} = L_{i,j \times h_n + h}, \forall i \in \{1, \dots, T\}, j \in \{1, \dots, d_h\} \quad (5)$$

where $h \in \{1, \dots, h_n\}$.

Remark: $L_{..h} : T \times d_h h_n \rightarrow T \times d_h$



Math format

Let $X \in \mathbb{R}^{T \times D}$ be the hidden input, $W_O \in \mathbb{R}^{d_h n_h \times D}$, $\text{Softmax}(X)$ is a row-wise operator, then the multi-head attention (**MLA**) is defined as follows:

- ▶ LoRA layer:

$$\begin{cases} L_Q = \text{LoRA}(X; W^{DQ}, W^{UQ}), \\ L_K = \text{LoRA}(X; W^{DKV}, W^{UK}), \\ V = \text{LoRA}(X; W^{DKV}, W^{UV}). \end{cases} \quad (6)$$

- ▶ RoPE layer:

$$R_Q = \text{RoPE}(X; W_{DQ}, W_{QR}), R_K = \text{RoPE}(X; W_{DKV}, W_{KR}) \quad (7)$$

- ▶ Cat the output of LoRA and ROPE for one head:

$$Q_h = [L_{Q..h}, R_{Q..h}], K_h = [L_{k..h}, R_{K..h}], \quad (8)$$

where $Q_h, R_K \in \mathbb{R}^{T \times (d_h + D^R)}$.



Math format

- ▶ Single head attention:

$$S_h = \text{Softmax}\left(\frac{Q_h K_h^T + M}{\sqrt{D^R + d_h}}\right) V_{..h}, \quad (9)$$

where $M \in \{-\infty, 0\}^{T \times T}$, $M_{ij} = -\infty$ if $i < j$ else $M_{ij} = 0$.

Remark: $S_h \in \mathbb{R}^{T \times T}$.

- ▶ Concatenate & Residual

$$\text{MLA}(X) := [S_1, \dots, S_h, \dots, S_{h_n}] W_O, \quad (10)$$

Remark: $\text{MLA} : \mathbb{R}^{T \times D} \rightarrow \mathbb{R}^{T \times D}$.



Math format

The **MLP** layer is defined as:

$$\text{MLP}(X) := f_3 \circ (\sigma(f_2(X)) \odot f_1(X)), \quad (11)$$

where $f_i(X) := XW_i + b_i$, $b_i \in \mathbb{R}^{D_i}$, $W_1 \in \mathbb{R}^{D \times D_1}$, $W_i \in \mathbb{R}^{D_{i-1} \times D_i}$ and $W_L \in \mathbb{R}^{D_{L-1} \times D}$.

Remark: $\text{MLP} : \mathbb{R}^{T \times D} \rightarrow \mathbb{R}^{T \times D}$.

For each element x , the **active function** σ is defined as:

$$\text{SiLU}(x) := \frac{x}{1 + e^{-x}}. \quad (12)$$



Math format

Let $X \in \mathbb{R}^{T \times D}$ be the hidden input, the **Gate** is defined as follows:

- ▶ Obtain the weights of N_r experts for every tokens:

$$G = \text{Sigmoid}(XW_{experts}), \quad (13)$$

where $W_{experts} \in \mathbb{R}^{D \times N_r}$ and Sigmoid is a row-wise operator.

- ▶ Select top K_r experts for every tokens:

$$g'_{i,t} = \begin{cases} G_{i,t}, & G_{i,t} \in \text{Topk}(\{G_{i,k} \mid 1 \leq k \leq N_r\}, K_r) \\ 0, & \text{otherwise,} \end{cases} \quad (14)$$

- ▶ Normalization:

$$\text{Gate}(X)_{it} := \frac{g'_{i,t}}{\sum_{t=1}^{N_r} g'_{j,t}} \quad (15)$$

Remark: $\text{Gate}(X) : \mathbb{R}^{T \times D} \rightarrow \mathbb{R}^{T \times N_r}$. In the code only need $T \times K_r$ because the remains are zero.



Math format

Let $X \in \mathbb{R}^{T \times D}$ be the hidden input. There are N_r routed experts (\mathbf{MLP}^r) and N_s shared experts (\mathbf{MLP}^s), then the **MOE** layer is defined as:

$$\text{MoE}(X) = \sum_{i=1}^{N_s} \text{MLP}_i^s(X) + \sum_{j=1}^{N_r} \text{Gate}(X)_{\cdot j} \text{MLP}_j^r(X), \quad (16)$$

Remark: $\text{MoE}(X) : \mathbb{R}^{T \times D} \rightarrow \mathbb{R}^{T \times D}$.



Math format

The **Block** in DeepSeek is defined as follows:

- ▶ The attention layer

$$\text{ResMLA}(X) := X + \text{MLA} \circ \text{RMSNorm}(X) \quad (17)$$

- ▶ If the block number less than the number of dense layers (`n_dense_layers`).

$$\text{ResMLP}(X) := X + \text{MLP} \circ \text{RMSNorm}(X), \quad (18)$$

- ▶ For the other blocks

$$\text{ResMoE}(X) := X + \text{MoE} \circ \text{RMSNorm}(X), \quad (19)$$

- ▶ The i -th block is defined as:

$$\text{Block}_i(X) = \begin{cases} \text{ResMLP} \circ \text{ResMLA}(X), & i < \text{n_dense_layers}, \\ \text{ResMoE} \circ \text{ResMLA}(X), & \text{others} \end{cases} \quad (20)$$

The **embedding** is defined as:

$$X^e := \text{RMSNorm} \circ \text{Block}_M \cdots \text{Block}_1 \circ \text{Enc}(X), \quad (21)$$

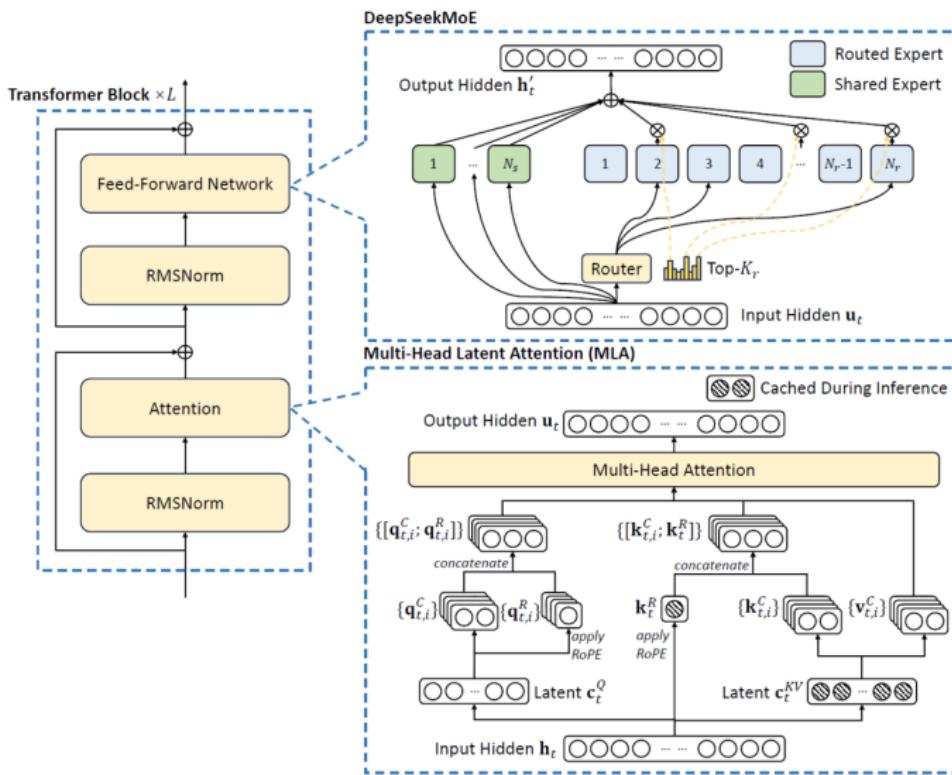
where $X^e \in \mathbb{R}^{T \times D}$. The output of **DeepSeek** is defined as:

$$\text{DeepSeek} := \arg \max_{\text{index}} X^i = X^e W_{\text{head}}, \quad (22)$$

where $W_{\text{head}} \in \mathbb{R}^{D \times d}$.



MLA and MOE



DeepSeek parameters

Parameter	168B	236B	671B
vocab_size (d)	102400	102400	129280
dim (D)	2048	5120	7168
inter_dim (MLP)	10944	12288	18432
moe_inter_dim (MoE)	1408	1536	2048
n_layers (M)	27	60	61
n_dense_layers	1	1	3
n_heads (n_h)	16	128	128
n_routed_experts (N_r)	64	160	256
n_shared_experts (N_s)	2	2	1
n_activated_experts (K_r)	6	6	8
q_lora_rank (LoRA)	0	1536	1536
kv_lora_rank (LoRA)	512	512	512
v_head_dim (h_n)	128	128	128



Contents

1 Architecture

2 Loss Functions

3 Training Framework

4 FP8 Training



Main Loss Function

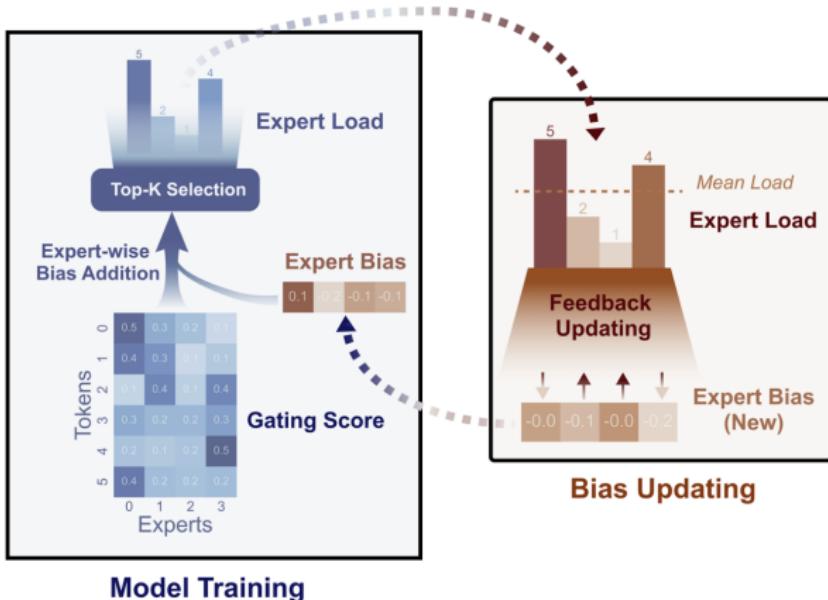
The main loss function is:

$$L_{main} = - \sum_{i=1}^T \log \text{Softmax}(X')_{ij_{\text{true}}} \quad (23)$$

, where j_{true} is the next token in the directory.



Auxiliary-Loss-Free Load Balancing



An unbalanced expert load will lead to routing collapse and diminish computational efficiency in scenarios with expert parallelism.

Auxiliary-Loss-Free Load Balancing

Algorithm 1: Adjusting the per-expert bias b_i during training

Input: MoE model θ , training batch iterator B , bias update rate u .

1. Initialize $b_i = 0$ for each expert;

for a batch $\{(\mathbf{x}_k, \mathbf{y}_k)\}_k$ in B **do**

2. Train MoE model θ on the batch data $\{(\mathbf{x}_k, \mathbf{y}_k)\}_k$, with gating scores calculated according to Eq. (3);

3. Count the number of assigned tokens c_i for each expert, and the average number \bar{c}_i ;

4. Calculate the load violation error $e_i = \bar{c}_i - c_i$;

4. Update b_i by $b_i = b_i + u * \text{sign}(e_i)$;

end

Output: trained model θ , corresponding bias b_i

$$g'_{i,t} = \begin{cases} s_{i,t}, & s_{i,t} + b_i \in \text{Topk}(\{s_{j,t} + b_j | 1 \leq j \leq N_r\}, K_r), \\ 0, & \text{otherwise.} \end{cases}$$



Complementary Sequence-Wise Auxiliary Loss

$$\mathcal{L}_{\text{Bal}} = \alpha \sum_{i=1}^{N_r} f_i P_i,$$

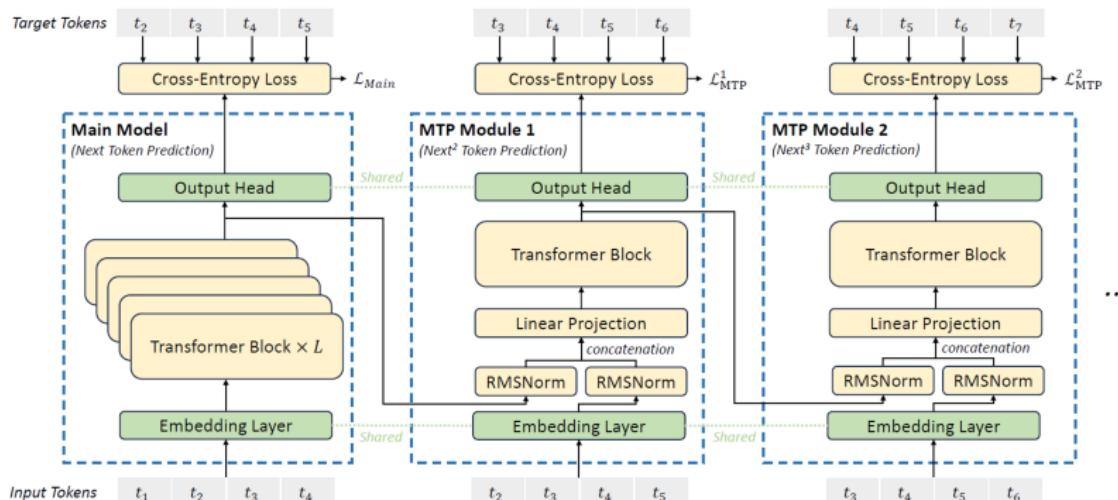
$$f_i = \frac{N_r}{K_r T} \sum_{t=1}^T \mathbb{1} (s_{i,t} \in \text{Topk}(\{s_{j,t} | 1 \leq j \leq N_r\}, K_r)),$$

$$s'_{i,t} = \frac{s_{i,t}}{\sum_{j=1}^{N_r} s_{j,t}},$$

$$P_i = \frac{1}{T} \sum_{t=1}^T s'_{i,t},$$



Multi-Token Prediction



$$\mathcal{L}_{\text{MTP}}^k = \text{CrossEntropy}(P_{2+k:T+1}^k, t_{2+k:T+1}) = -\frac{1}{T} \sum_{i=2+k}^{T+1} \log P_i^k[t_i],$$



Contents

1 Architecture

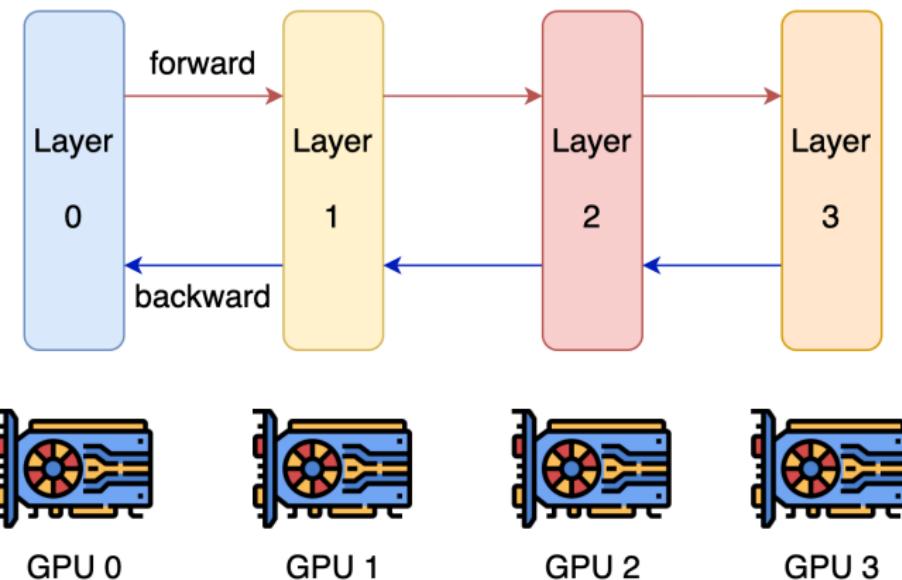
2 Loss Functions

3 Training Framework

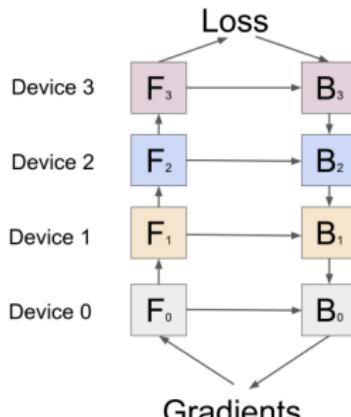
4 FP8 Training



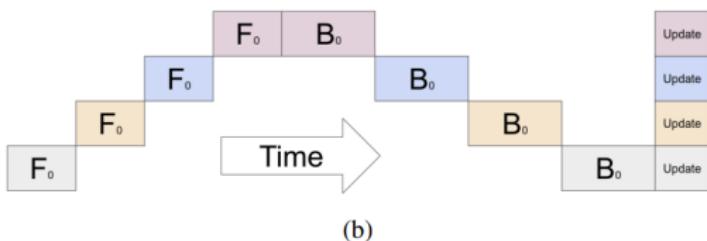
Background – Model Parallel



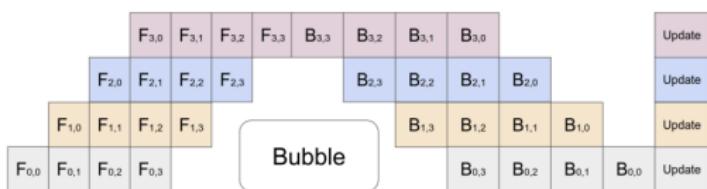
Background – Pipeline Parallel



(a)



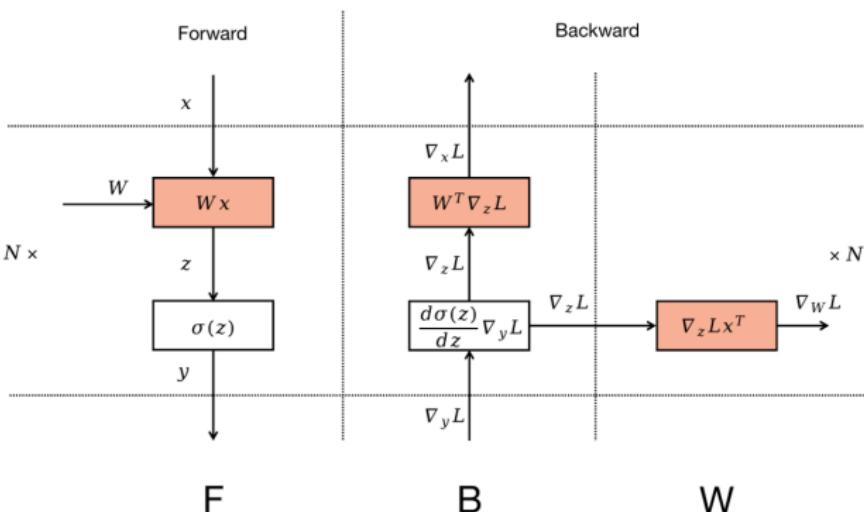
(b)



(c)

DeepSeek divides each chunk into four components: attention, all-to-all dispatch, MLP, and all-to-all combine

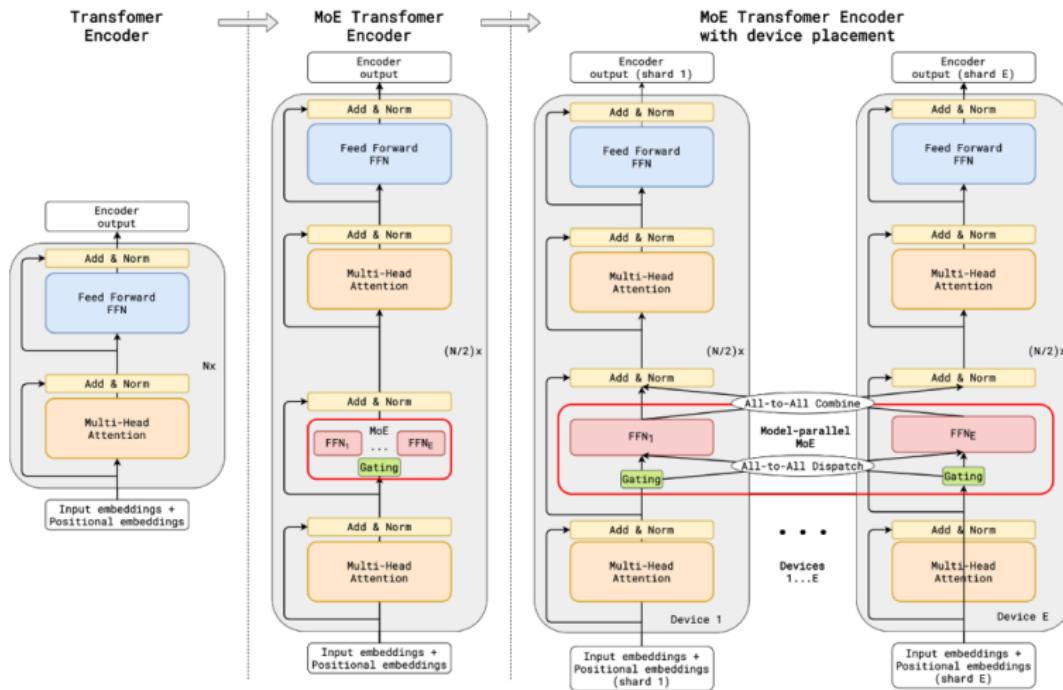
Attention and MLP



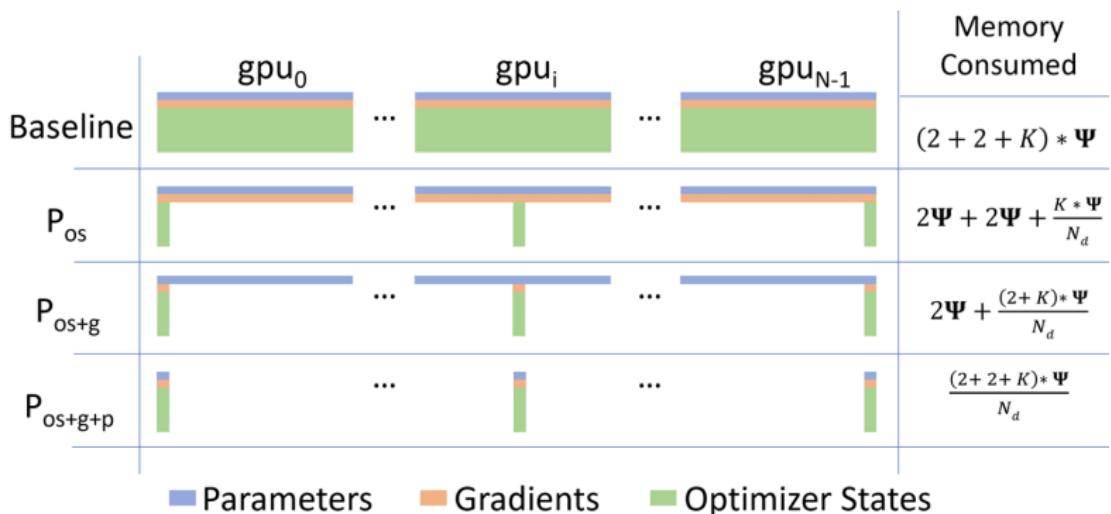
ZERO BUBBLE PIPELINE PARALLELISM



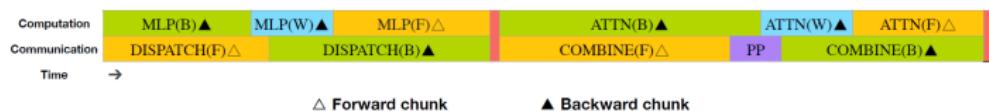
All-to-All Combine/Dispatch



ZeRO-123



Overlapping strategy



DualPipe



Contents

1 Architecture

2 Loss Functions

3 Training Framework

4 FP8 Training



FP8 Format

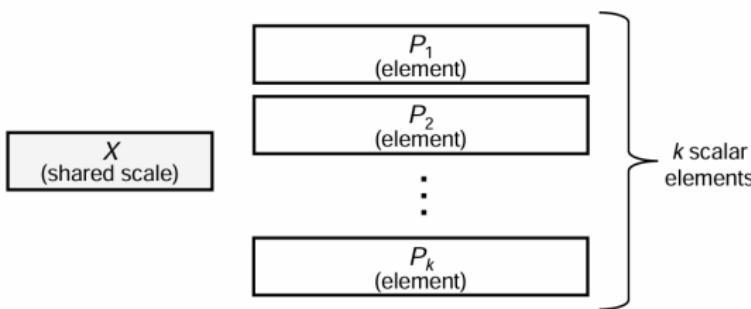


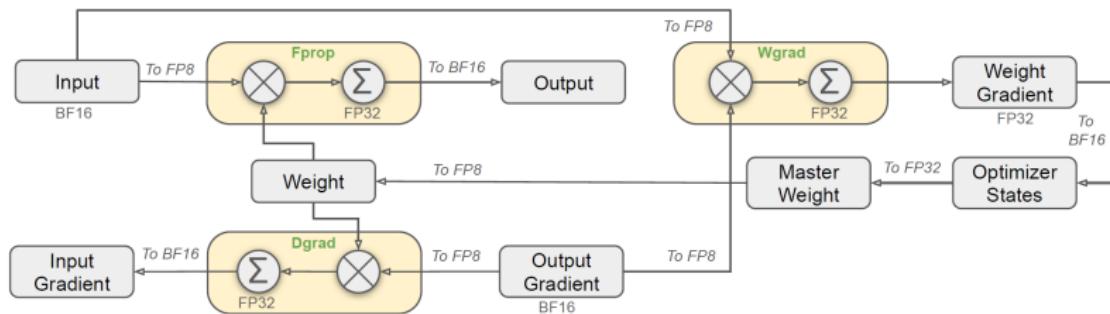
Figure 1: A single block in a Microscaling data format. The block encodes a vector of k numbers, each with value XP_i .

Let $\{v_i\}_{i=1}^k$ be the k real numbers represented in an MX block. The value of each number can be inferred as follows:

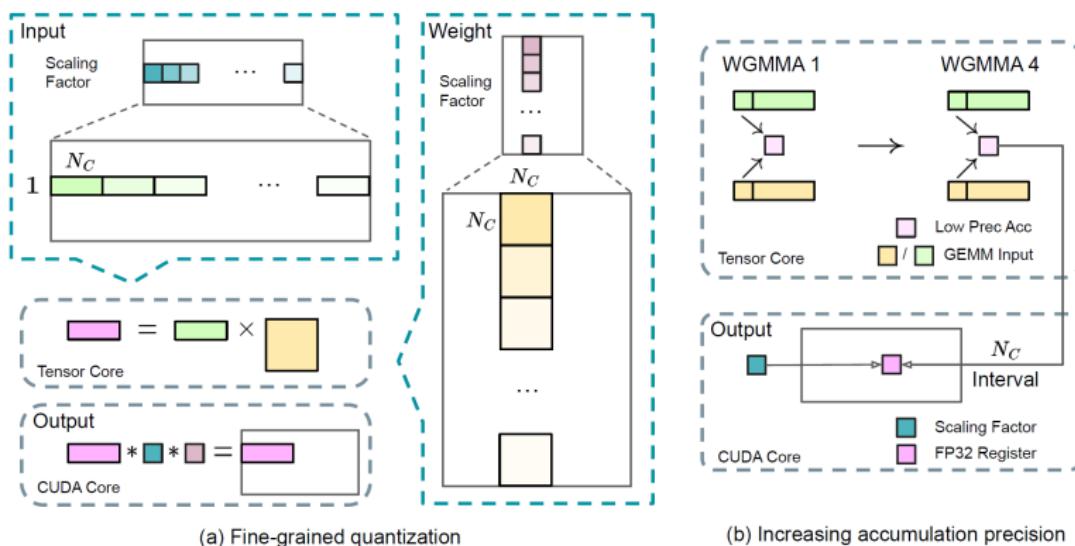
- If $X = \text{NaN}$, then $v_i = \text{NaN}$ for all i
- If $|XP_i| > V_{maxFloat32}$ then v_i is implementation-defined
- Otherwise, $v_i = XP_i$



FP8 Format



FP8 Format



THANKS!

