



Block-wise Training in LLMs

Kun Yuan (袁 坤)

Center for Machine Learning Research @ Peking University

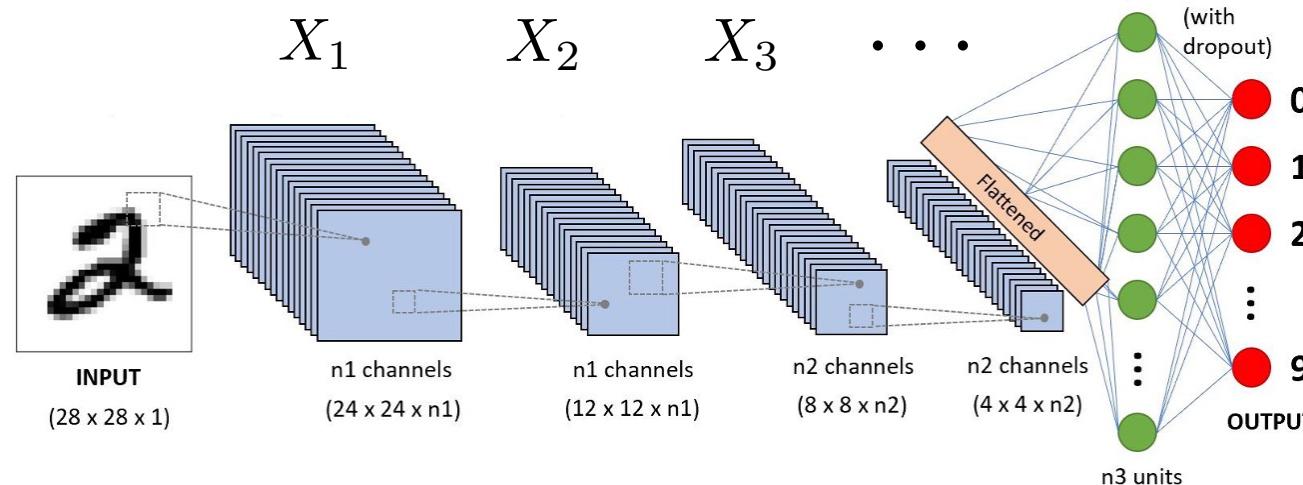


PART 01

Basics and Motivation

LLM pretraining is essentially solving stochastic optimization

- The model weights in neural networks are a set of matrices $\mathbf{X} = \{\mathbf{X}_\ell\}_{\ell=1}^L$



- Let $h(\mathbf{X}; \xi)$ be the language model; $\hat{y} = h(\mathbf{X}; \xi)$ is the predicted token

cross entropy

LLM cost function:

$$\mathbf{X}^* = \arg \min_{\mathbf{X}} \left\{ \mathbb{E}_{\xi \sim \mathcal{D}} \left[L(h(\mathbf{X}; \xi), y) \right] \right\}$$

↑
↓ ↓ ↓
data distribution **pred. token** **real token**

LLM pretraining is essentially solving stochastic optimization

- If we define $\xi = (\xi, y)$ and $F(\mathbf{X}; \xi) = L(h(\mathbf{X}; \xi), y)$, the LLM problem becomes

Stochastic optimization: $\mathbf{X}^* = \arg \min_{\mathbf{X}} \left\{ \mathbb{E}_{\xi \sim \mathcal{D}} [F(\mathbf{X}; \xi)] \right\}$

- In other words, LLM pretraining is essentially solving a stochastic optimization problem
- Adam is the standard approach in LLM pretraining

Optimizer states

$$\begin{aligned}
 \mathbf{G}_t &= \nabla F(\mathbf{X}_t; \xi_t) && \text{(stochastic gradient)} \\
 \boxed{\mathbf{M}_t} &= (1 - \beta_1) \mathbf{M}_{t-1} + \beta_1 \mathbf{G}_t && \text{(first-order momentum)} \\
 \boxed{\mathbf{V}_t} &= (1 - \beta_2) \mathbf{V}_{t-1} + \beta_2 \mathbf{G}_t \odot \mathbf{G}_t && \text{(second-order momentum)} \\
 \mathbf{X}_{t+1} &= \mathbf{X}_t - \frac{\gamma}{\sqrt{\mathbf{V}_t} + \epsilon} \odot \mathbf{M}_t && \text{(adaptive SGD)}
 \end{aligned}$$

Memory = **Model** + **Gradient** + **Optimizer states** + **Activations**

- Given a **model** with **P** parameters, **gradient** will consume **P** parameters, and **optimizer states** will consume **2P** parameters; **4P parameters in total**.

$$\mathbf{P} \quad \mathbf{G}_t = \nabla F(\mathbf{X}_t; \boldsymbol{\xi}_t)$$

$$2\mathbf{P} \quad \left\{ \begin{array}{l} \mathbf{M}_t = (1 - \beta_1)\mathbf{M}_{t-1} + \beta_1 \mathbf{G}_t \\ \mathbf{V}_t = (1 - \beta_2)\mathbf{V}_{t-1} + \beta_2 \mathbf{G}_t \odot \mathbf{G}_t \end{array} \right.$$

$$\mathbf{P} \quad \mathbf{X}_{t+1} = \mathbf{X}_t - \frac{\gamma}{\sqrt{\mathbf{V}_t} + \epsilon} \odot \mathbf{M}_t$$

Optimizer states introduces significant memory cost

Memory = **Model** + **Gradient** + **Optimizer states** + **Activations**

- Activations are auxiliary variables to facilitate the gradient calculations

Consider a linear neural network

$$z_i = X_i z_{i-1}, \forall i = 1, \dots, L$$

$$f = \mathcal{L}(z_L; y)$$

The gradient is derived as follows

$$\frac{\partial f}{\partial X_i} = \frac{\partial f}{\partial z_i} z_{i-1}^\top$$

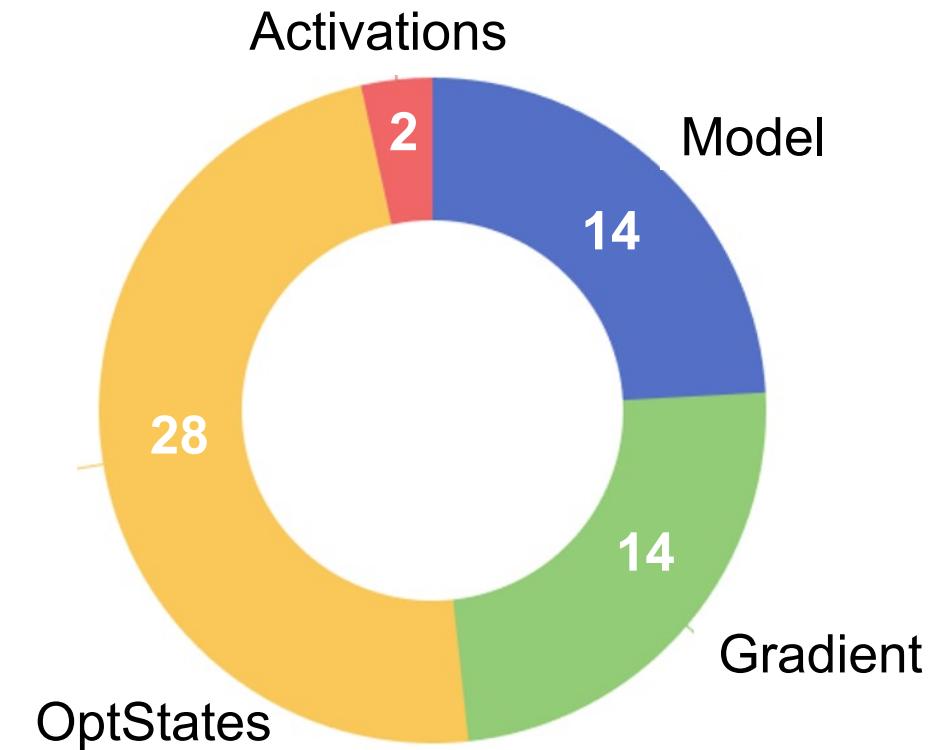
Need to store activations z_1, z_2, \dots, z_L

- The size of activations depends on sequence length and batch size

Minimum memory requirement

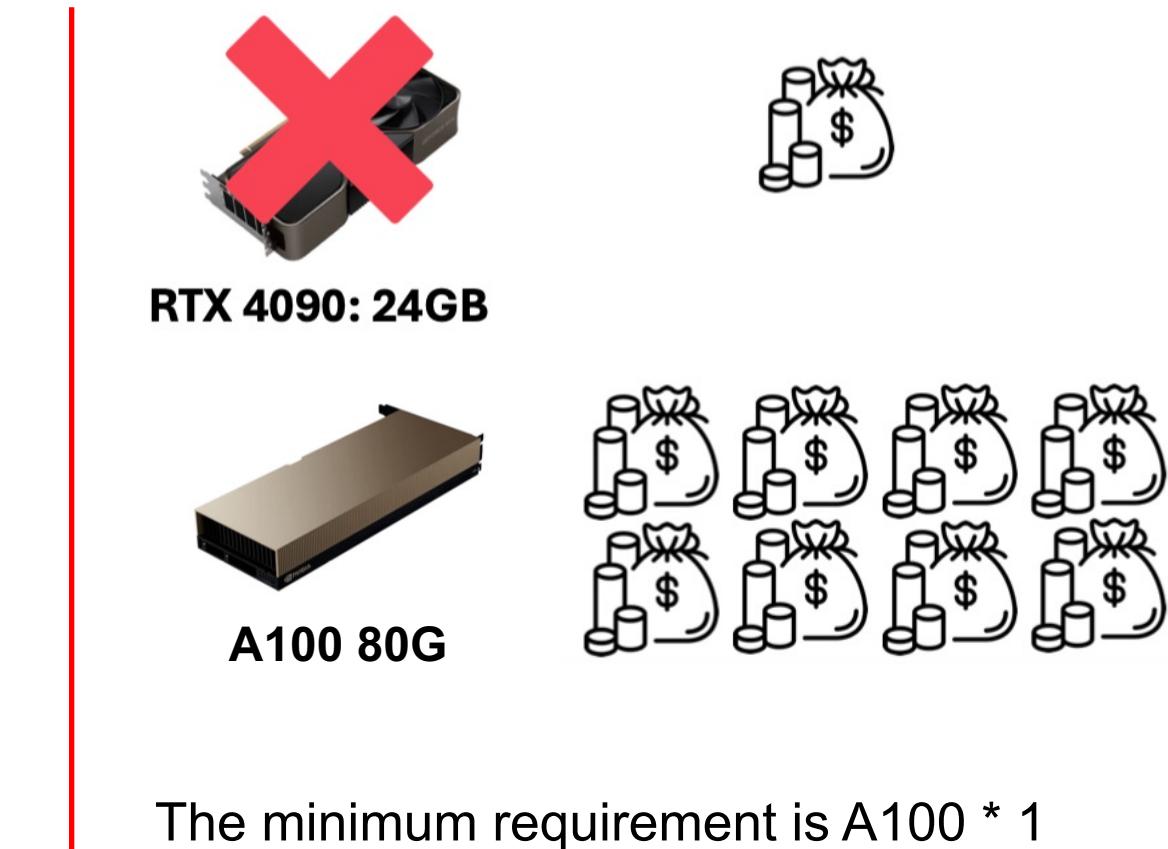
- Pretrain LLaMA 7B model (BF16) from scratch with **a single batch size** requires

- Parameters: 7B
- Model storage: $7B * 2 \text{ Bytes} = 14 \text{ GB}$
- Gradient storage: 14 GB
- Optimizer states: 28 GB (using Adam)
- Activation storage: 2 GB [Zhao et. al., 2024]
- In total: **58 GB**



Minimum memory requirement: LLaMA 7B

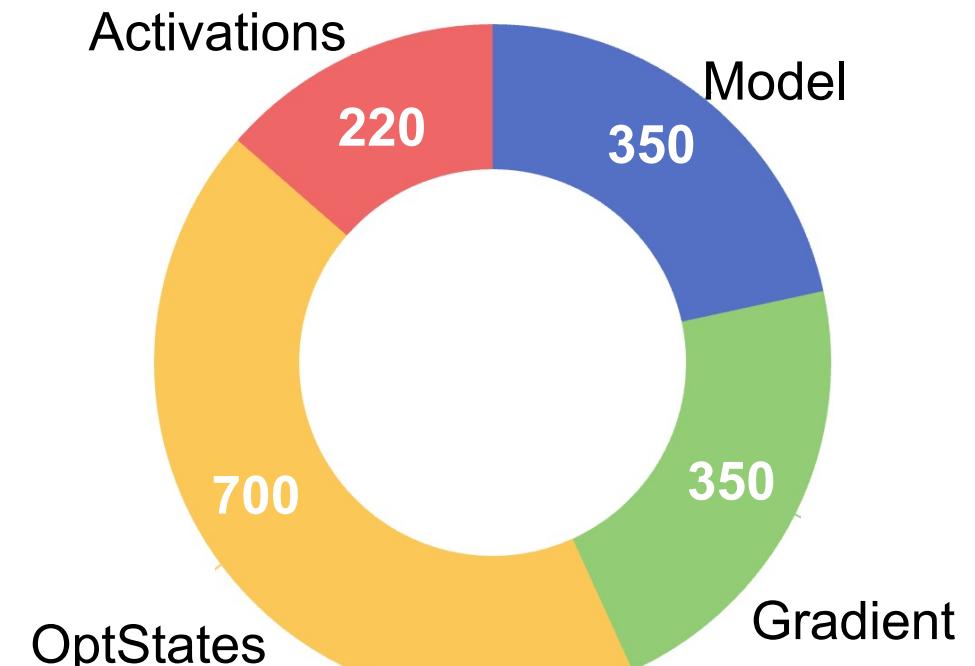
- Pretrain LLaMA 7B model (BF16) from scratch with **a single batch size** requires
 - Parameters: 7B
 - Model storage: $7B * 2 \text{ Bytes} = 14 \text{ GB}$
 - Gradient storage: 14 GB
 - Optimizer states: 28 GB (using Adam)
 - Activation storage: 2 GB [Zhao et. al., 2024]
 - In total: **58 GB**



Minimum memory requirement: GPT-3

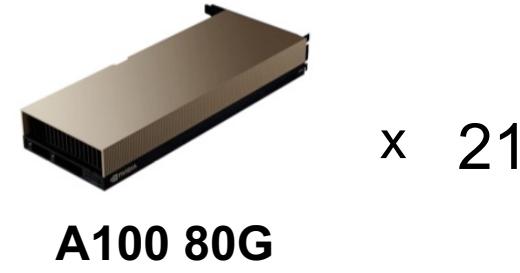
- Pretrain GPT-3 model (BF16) from scratch with **a single batch size** requires

- Parameters: 175B
- Model storage: $175\text{B} * 2 \text{ Bytes} = 350 \text{ GB}$
- Gradient storage: 350 GB
- Optimizer states: 700 GB (using Adam)
- Activation storage: ~220 GB
- In total: **1620 GB**



Minimum memory requirement: GPT-3

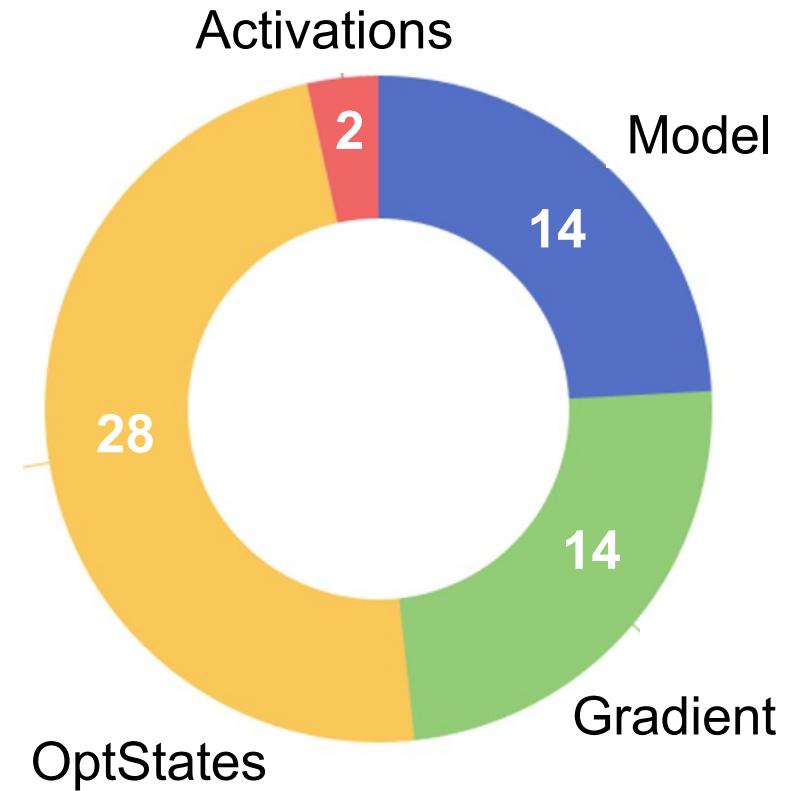
- Pretrain GPT-3 model (BF16) from scratch with **a single batch size** requires
 - Parameters: 175B
 - Model storage: $175\text{B} * 2 \text{ Bytes} = 350 \text{ GB}$
 - Gradient storage: 350 GB
 - Optimizer states: 700 GB (using Adam)
 - Activation storage: ~220 GB
 - In total: **1620 GB**



The minimum requirement is A100 * 21
Very expensive!

Memory-efficient algorithm is in urgent need

- With memory-efficient algorithms, we can
 - Train larger models** on limited computing resources
 - Use a larger training batch size to **improve throughput**
- Activation-incurred memory is **relatively minor** when using a single batch size
- Gradient-incurred memory can be **removed** by layer-wise calculation and dropping
- How to save memory caused by optimizer states?





PART 02

Block-wise training

Block-wise training

- Main idea: Layer-wise finetuning

$$\min_W \mathbb{E}_{\xi \sim \mathcal{D}}[F(W; \xi)] = \mathbb{E}_{\xi \sim \mathcal{D}}[F(W_1, W_2, \dots, W_L; \xi)]$$

where $W := \{W_1, W_2, \dots, W_L\}$. When update W_ℓ , the other layers remain freeze

Algorithm 1 Layerwise Importance Sampling AdamW (LISA)

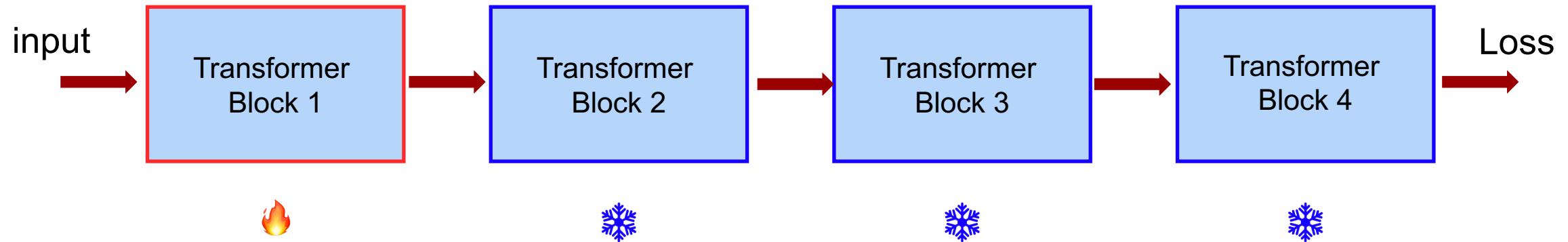
Require: number of layers N_L , number of iterations T , sampling period K , number of sampled layers γ , initial learning rate η_0

- 1: **for** $i \leftarrow 0$ to $T/K - 1$ **do**
 - 2: Freeze all layers except the embedding and language modeling head layer
 - 3: Randomly sample γ intermediate layers to unfreeze
 - 4: Run AdamW for K iterations with $\{\eta_t\}_{t=iK}^{ik+k-1}$
 - 5: **end for**
-

[LISA: Layerwise Importance Sampling for Memory-Efficient Large Language Model Fine-Tuning]

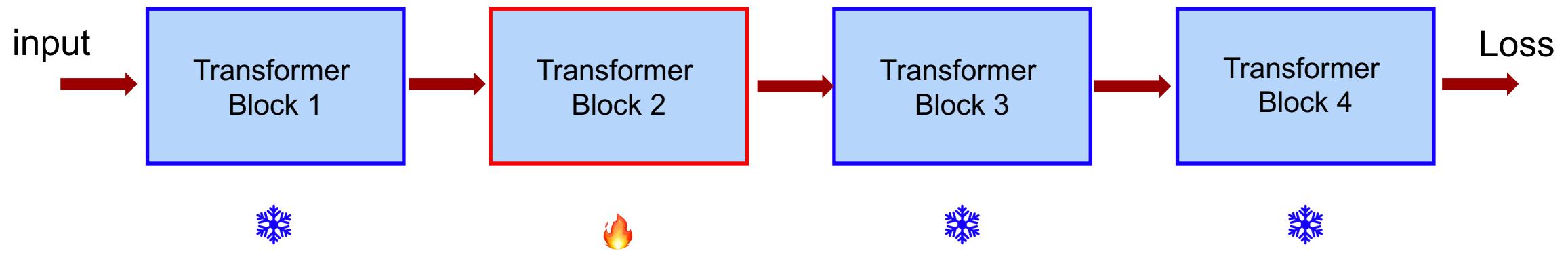
Block-wise training

Update the first block



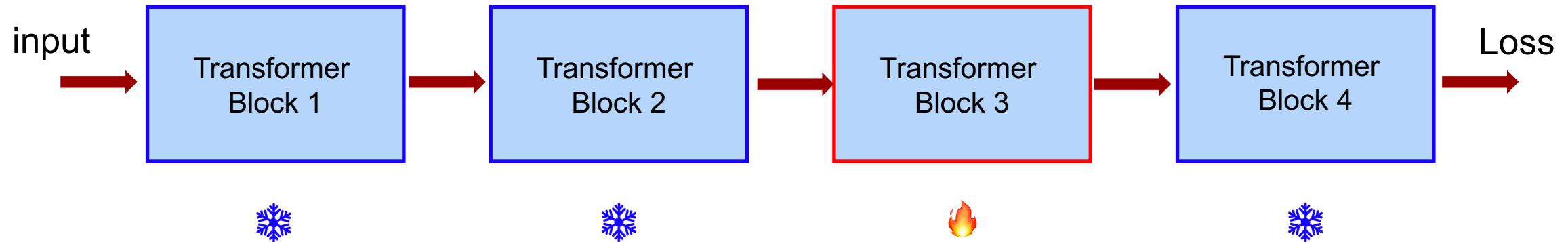
Block-wise training

Update the second block



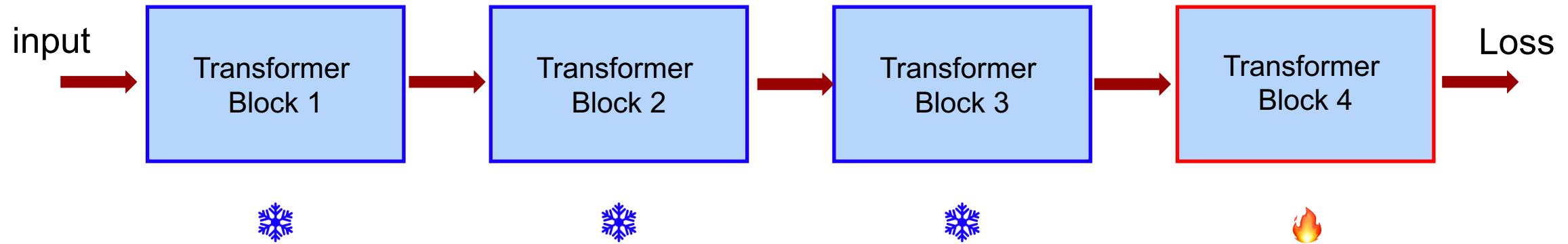
Block-wise training

Update the third block



Block-wise training

Update the forth block



Memory = **Model** + **Gradient** + **Optimizer states** + **Activations**

Original

$$\mathbf{P} \quad \mathbf{G}_t = \nabla F(\mathbf{X}_t; \xi_t)$$

$$2\mathbf{P} \quad \begin{cases} \mathbf{M}_t = (1 - \beta_1)\mathbf{M}_{t-1} + \beta_1 \mathbf{G}_t \\ \mathbf{V}_t = (1 - \beta_2)\mathbf{V}_{t-1} + \beta_2 \mathbf{G}_t \odot \mathbf{G}_t \end{cases}$$

$$\mathbf{P} \quad \mathbf{X}_{t+1} = \mathbf{X}_t - \frac{\gamma}{\sqrt{\mathbf{V}_t} + \epsilon} \odot \mathbf{M}_t$$

Block-wise

$$G_{t,i} = \nabla_i F(\mathbf{X}_t; \xi_t) \quad \mathbf{P/L}$$

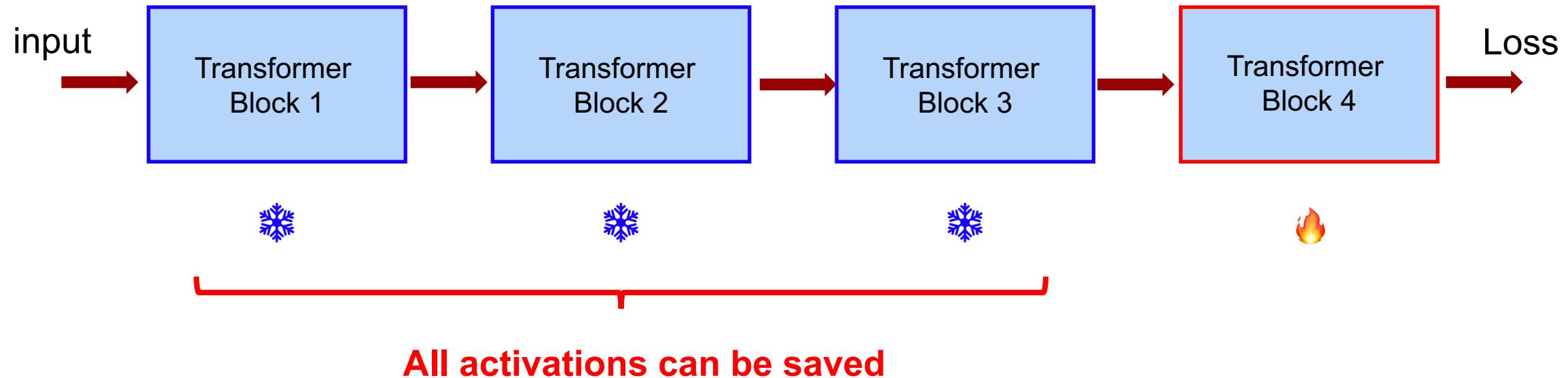
$$\boxed{\begin{aligned} M_{t,i} &= (1 - \beta_1)M_{t-1,i} + \beta_1 G_{t,i} & \mathbf{2P/L} \\ V_{t,i} &= (1 - \beta_2)V_{t-1,i} + \beta_2 G_{t,i} \odot G_{t,i} \end{aligned}}$$

$$X_{t+1,i} = X_{t,i} - \frac{\gamma}{\sqrt{V_{t,i}} + \epsilon} \odot M_{t,i} \quad \mathbf{P}$$

Block-wise training can save gradients and optimization states significantly

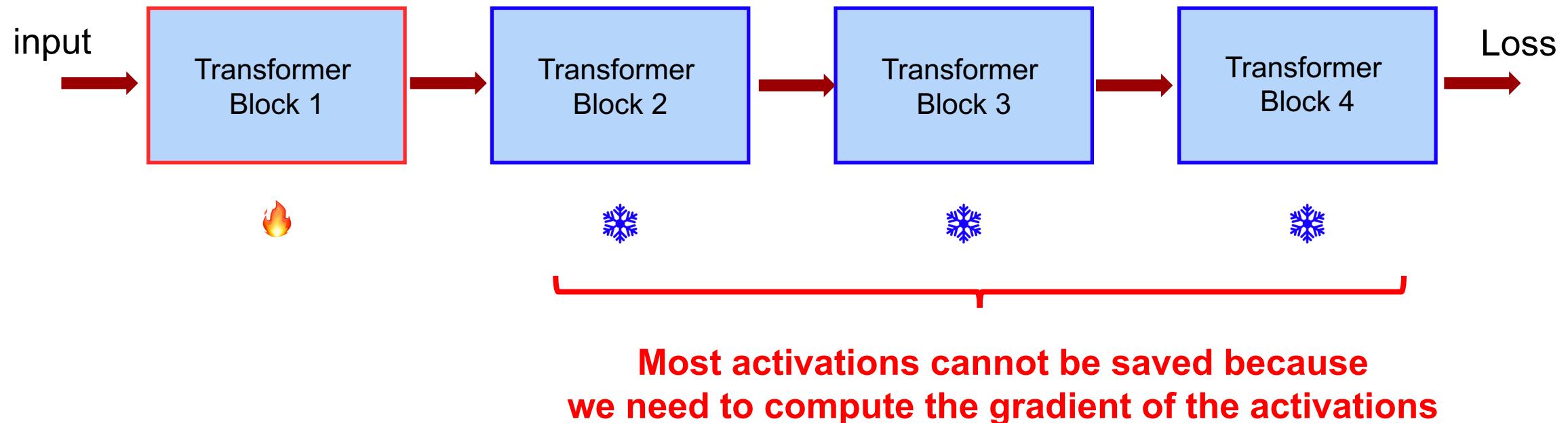
Memory cost

- Block-wise training can also save activations, but depends on which block has been activated

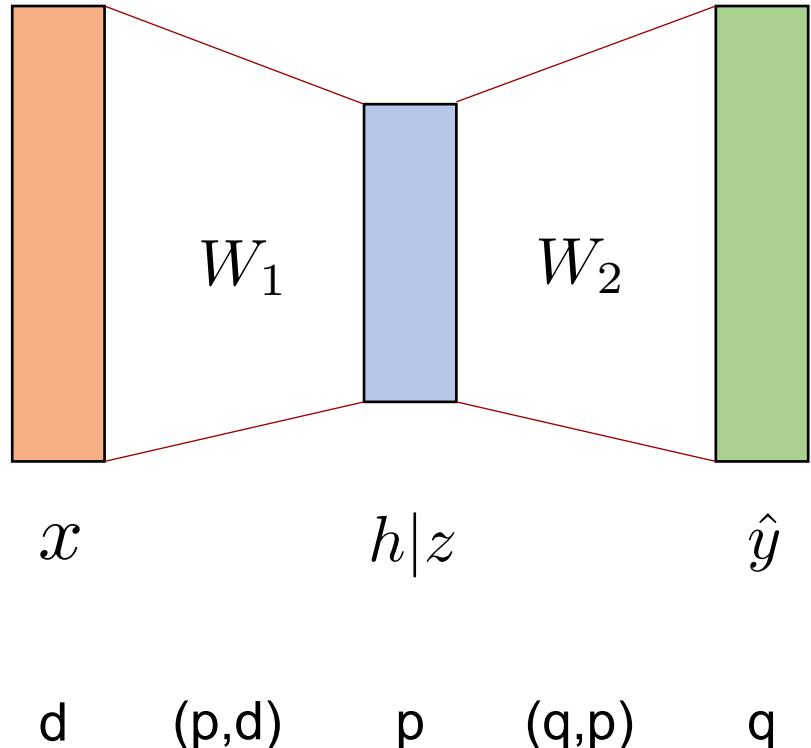


Memory cost

- Block-wise training can also save activations, but depends on which block has been activated



Example: 2-layer MLP



$$h = W_1 x$$

$$z = \sigma(h)$$

$$\hat{y} = W_2 z$$

$$f = L(\hat{y})$$

Forward

$$\frac{\partial f}{\partial W_1} = \frac{\partial f}{\partial h} x^T$$

$$\frac{\partial f}{\partial h} = \frac{\partial f}{\partial z} \odot \nabla \sigma(h)$$

$$\frac{\partial f}{\partial W_2} = \frac{\partial L}{\partial \hat{y}} z^T, \quad \frac{\partial f}{\partial z} = W_2^\top \frac{\partial L}{\partial \hat{y}}$$

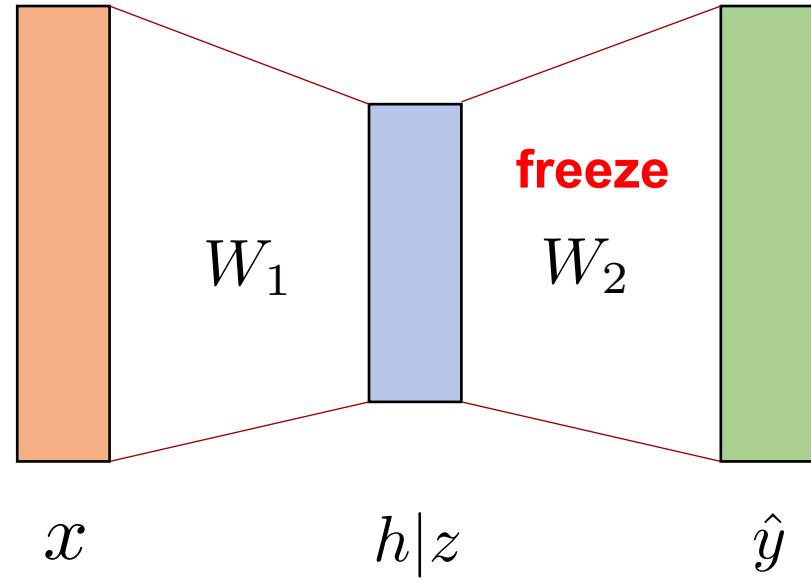
$$\frac{\partial f}{\partial \hat{y}} = \nabla L(\hat{y})$$

Backward

Store h , z and \hat{y}

Store $\nabla_{W_1} f(W_1)$ and $\nabla_{W_2} f(W_2)$

Example: 2-layer MLP



Save activations

$$h = W_1 x$$

$$z = \sigma(h)$$

$$\hat{y} = W_2 z$$

$$f = L(\hat{y})$$

Forward

$$\frac{\partial f}{\partial W_1} = \frac{\partial f}{\partial h} x^T$$

$$\frac{\partial f}{\partial h} = \frac{\partial f}{\partial z} \odot \nabla \sigma(h)$$

$$\frac{\partial f}{\partial z} = W_2^T \frac{\partial L}{\partial \hat{y}}$$

$$\frac{\partial f}{\partial \hat{y}} = \nabla L(\hat{y})$$

Backward

Activation h is still needed

Example: Attention

self attention

$$(i) \quad Q = X \underbrace{WQ}_{(N, d_{\text{model}})}$$

$$K = X \underbrace{W_K}_{(d_{\text{model}}, d_k)} \quad (d_{\text{model}}, d_k)$$

$$V = X \underbrace{W_V}_{(N, d_v)}$$

$$A = \frac{QK^T}{\sqrt{d_k}}$$

Save Q, K, V

Do not save X

$$(ii) \quad S_{i,j} = \frac{e^{A_{i,j}}}{\sum_{k=1}^N e^{A_{ik}}} \quad (\text{get } S)$$

$$O = SV \quad (N, d_k)$$

$$\text{Attn} = O \underbrace{W_O}_{(N, d_{\text{model}})}$$

Save S, V

Do not save A , O

Example: Attention

$$B^{(i)}_{ij} \frac{\partial s_{ij}}{\partial A_{mn}} = \begin{cases} 0, & \text{if } m \neq i \\ -\frac{e^{A_{im}} e^{A_{ij}}}{(\sum_{k=1}^N e^{A_{ik}})^2} = -S_{in} S_{ij}, & \text{if } m=i \text{ 且 } n \neq j \\ \frac{e^{A_{ij}} (\sum_{k \neq j} e^{A_{ik}})}{(\sum_{k=1}^N e^{A_{ik}})^2} = S_{ij}(1-S_{ij}), & \text{if } m=i \text{ 且 } n=j \end{cases}$$

故 $\frac{\partial L}{\partial A_{ij}} = \sum_{k \neq j} \frac{\partial L}{\partial s_{ik}} (-S_{ik} S_{ij}) + \frac{\partial L}{\partial s_{ij}} S_{ij}(1-S_{ij})$
 $= S_{ij} \left(\frac{\partial L}{\partial s_{ij}} - \sum_{k=1}^N \frac{\partial L}{\partial s_{ik}} S_{ik} \right)$ 要同时用到 $\frac{\partial L}{\partial s}$ 与 S 在，大小 $N \times N$ (B1)

$(N, d_k) \quad (N, N) \quad (N, d_k)$
 $\frac{\partial L}{\partial V} = S^T \frac{\partial L}{\partial O} \quad (B2)$

且 S 的稀疏化会
影响 Back Prop

Experiments

While Block-wise training is memory-efficient, it converges slower than standard training.

Existing literature uses it to fine-tune, not to pretrain.

Fine-tune on a single RTX 3090

Method	Parameter	Gradient	Optimizer states	Memory consumption
Adam	16.1GB	32.1GB	96.6GB	144.8GB+
LOMO	16.1GB	0.5GB	—	21.5GB
LoRA-rank100	16.7GB	1.0GB	3.1GB	OOM
LoRA-rank8	16.2GB	0.1GB	0.3GB	22.3GB
Block-wise	BAdam	0.9GB	2.6GB	23.5GB

[BAdam: A Memory Efficient Full Parameter Optimization Method for Large Language Models]

Experiments

Model: Llama 3-8B							
Method	GSM8K	Aqua	MMLU-Math	SAT-Math	MATH	NumGLUE	Average
Base model	25.9	22.8	33.7	39.5	12.8	34.5	28.2
Adam	54.5	40.5	44.3	51.4	18.4	55.4	44.1
LOMO	32.1	28.0	40.0	39.5	13.1	37.1	31.6
LoRA	47.5	44.9	45.3	50.9	14.5	56.9	43.3
Galore	33.1	37.4	41.2	42.7	15.0	36.9	34.4
BAdam	48.1	42.5	50.5	56.8	15.7	53.0	44.4

Model: Llama 3-70B							
Method	GSM8K	Aqua	MMLU-Math	SAT-Math	MATH	NumGLUE	Average
Base model	52.4	46.5	52.2	58.2	21.2	37.9	44.7
LoRA	73.3	59.5	58.3	64.1	34.2	64.8	59.0
BAdam	78.8	63.4	64.2	76.4	26.2	67.3	62.7

Table 6: Zero-shot math benchmark scores of the finetuned Llama 3-8B and Llama 3-70B on MathInstruct dataset by different methods.