

Attentions and How to Run Them Fast

By Kaizhao

Agenda

- Basic math stuffs
- Exact Attentions
 - Flash Attention 1-3
 - Ring Attention
 - Tree Attention
- Attention as RNN
- Approximate versions
 - StreamingLLM
 - Linear Attention (fast weight perspective)
- Speculative Decoding
- Trend and future

Basic Math stuffs

$$\vec{S} = [s_0, s_1, s_2, \dots, s_T]$$

$$\text{Softmax}(\vec{S})_t = \frac{e^{\vec{S}_t}}{\sum_t^T e^{\vec{S}_t}}$$

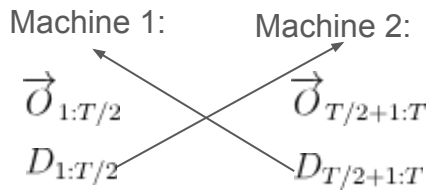
Stare at it for a few seconds, tell me what you see?

always has been

wait, it's just RNN with growing memory?

Basic MLE interview question 2:

Now imagine you have two machines, and a very long \vec{S} how do you parallelize your softmax?



Corrected version:

$$\vec{O}'_{1:T/2} = \vec{O}_{1:T/2} \cdot \frac{D_{1:T/2}}{D_{1:T/2} + D_{T/2+1:T}}$$
$$\vec{O}'_{T/2+1:T} = \vec{O}_{T/2+1:T} \cdot \frac{D_{T/2+1:T}}{D_{1:T/2} + D_{T/2+1:T}}$$

Basic MLE interview question 2.5:

What if you have N machines, how do you build an attention mechanism that can run at $O(T/N \log(T/N))$

What is the communication cost?

What is the memory complexity?

Exercise left for audience offline

Basic MLE interview question 3:

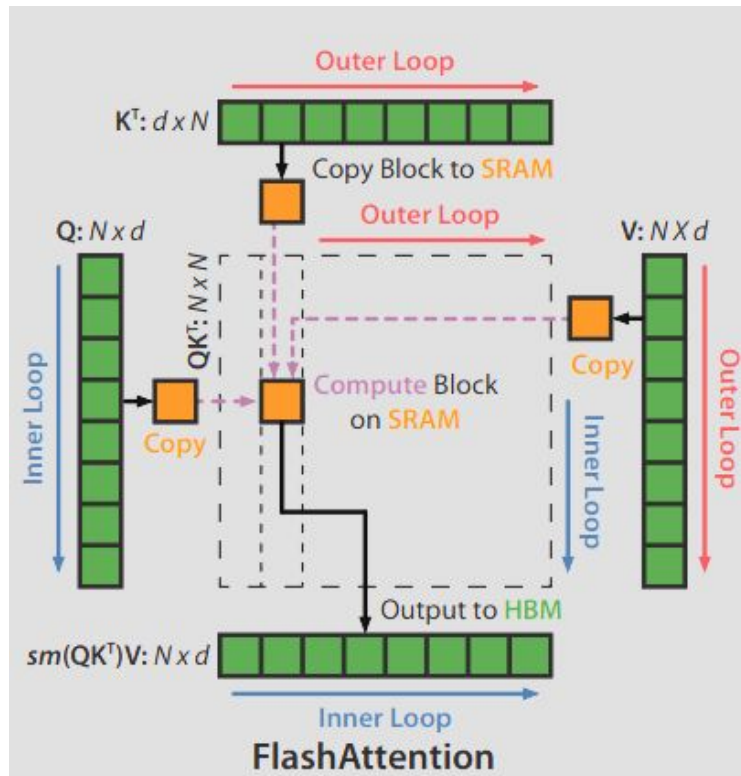
What function's derivative is Softmax?

$$\frac{\partial \log \text{sumexp}(S + x)}{\partial x} \Big|_{x=0} = \text{Softmax}(S)$$

Remember what we've learnt so far and let's get started



Flash Attention

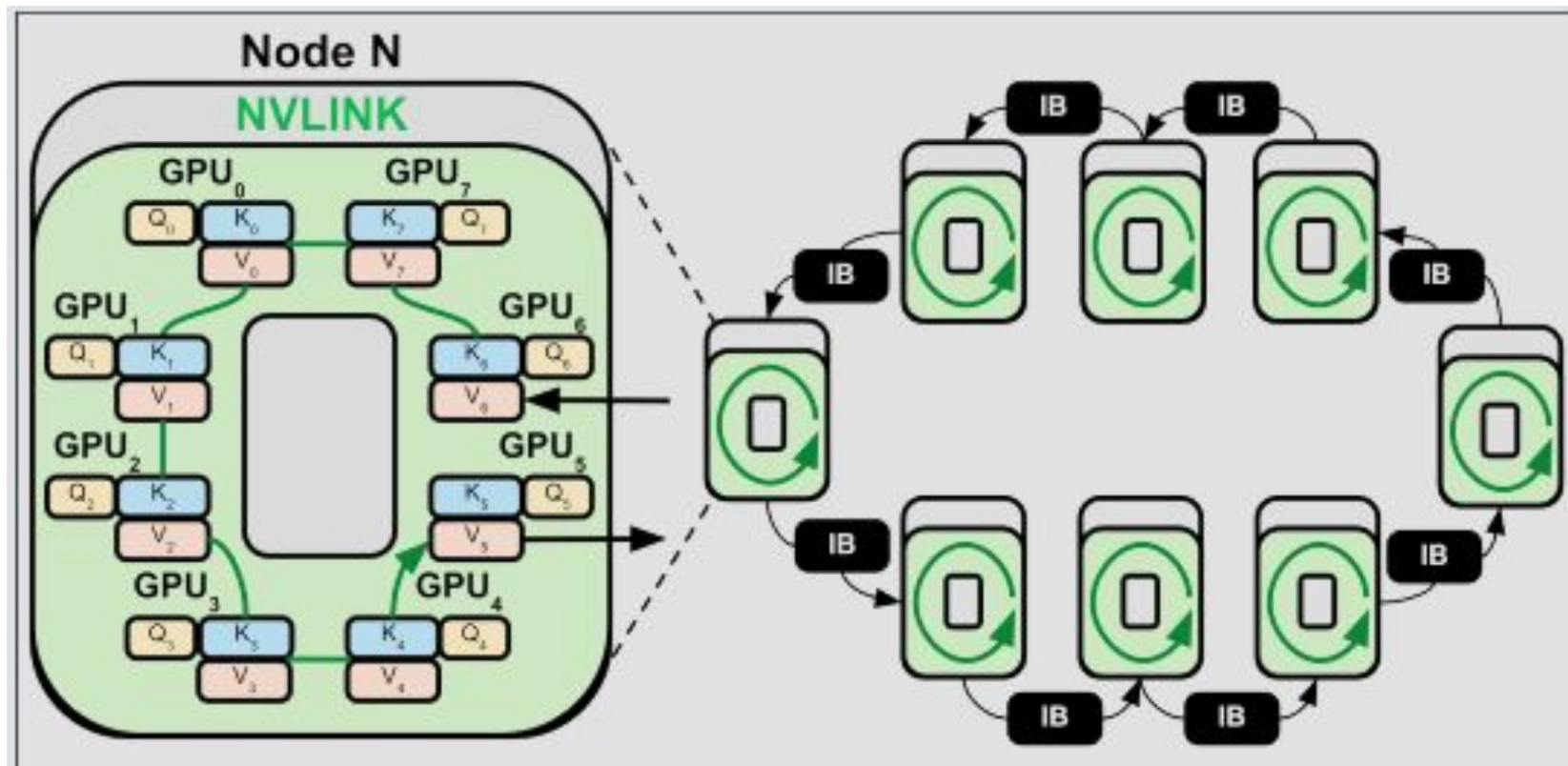


Algorithm 1 FLASHATTENTION

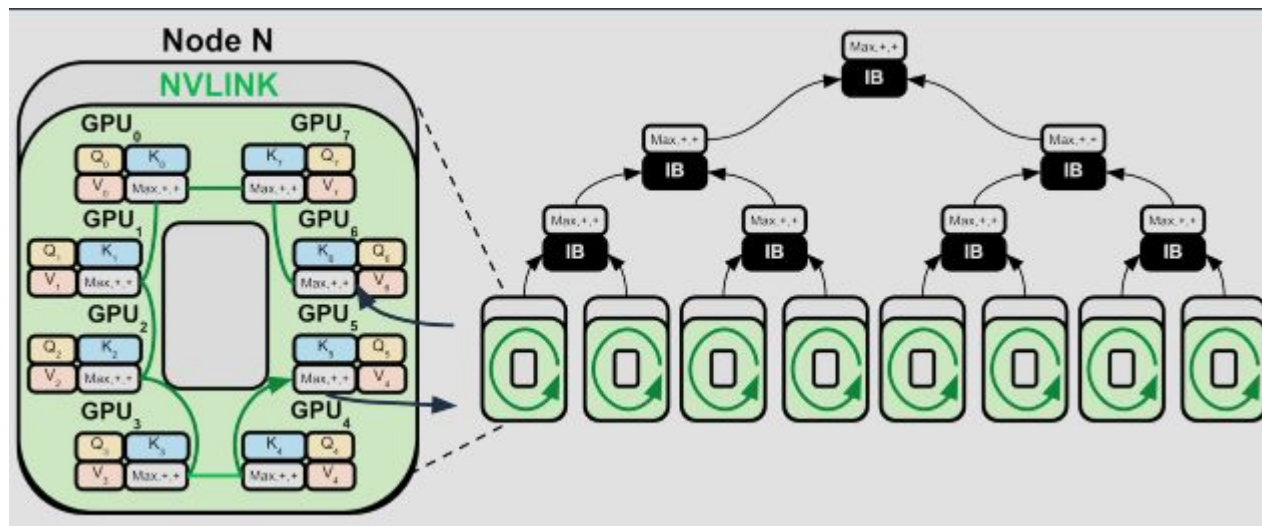
Require: Matrices $Q, K, V \in \mathbb{R}^{N \times d}$ in HBM, on-chip SRAM of size M .

- 1: Set block sizes $B_c = \lceil \frac{M}{4d} \rceil$, $B_r = \min(\lceil \frac{M}{4d} \rceil, d)$.
- 2: Initialize $O = (0)_{N \times d} \in \mathbb{R}^{N \times d}$, $\ell = (0)_N \in \mathbb{R}^N$, $m = (-\infty)_N \in \mathbb{R}^N$ in HBM.
- 3: Divide Q into $T_r = \lceil \frac{N}{B_r} \rceil$ blocks Q_1, \dots, Q_{T_r} of size $B_r \times d$ each, and divide K, V into $T_c = \lceil \frac{N}{B_c} \rceil$ blocks K_1, \dots, K_{T_c} and V_1, \dots, V_{T_c} , of size $B_c \times d$ each.
- 4: Divide O into T_r blocks O_1, \dots, O_{T_r} of size $B_r \times d$ each, divide ℓ into T_r blocks $\ell_1, \dots, \ell_{T_r}$ of size B_r each, divide m into T_r blocks m_1, \dots, m_{T_r} of size B_r each.
- 5: **for** $1 \leq j \leq T_c$ **do**
- 6: Load K_j, V_j from HBM to on-chip SRAM.
- 7: **for** $1 \leq i \leq T_r$ **do**
- 8: Load Q_i, O_i, ℓ_i, m_i from HBM to on-chip SRAM.
- 9: On chip, compute $S_{ij} = Q_i K_j^T \in \mathbb{R}^{B_r \times B_c}$.
- 10: On chip, compute $\tilde{m}_{ij} = \text{rowmax}(S_{ij}) \in \mathbb{R}^{B_r}$, $\tilde{P}_{ij} = \exp(S_{ij} - \tilde{m}_{ij}) \in \mathbb{R}^{B_r \times B_c}$ (pointwise), $\tilde{\ell}_{ij} = \text{rowsum}(\tilde{P}_{ij}) \in \mathbb{R}^{B_r}$.
- 11: On chip, compute $m_i^{\text{new}} = \max(m_i, \tilde{m}_{ij}) \in \mathbb{R}^{B_r}$, $\ell_i^{\text{new}} = e^{m_i - m_i^{\text{new}}} \ell_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\ell}_{ij} \in \mathbb{R}^{B_r}$.
- 12: Write $O_i \leftarrow \text{diag}(\ell_i^{\text{new}})^{-1} (\text{diag}(\ell_i) e^{m_i - m_i^{\text{new}}} O_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{P}_{ij} V_j)$ to HBM.
- 13: Write $\ell_i \leftarrow \ell_i^{\text{new}}$, $m_i \leftarrow m_i^{\text{new}}$ to HBM.
- 14: **end for**
- 15: **end for**
- 16: **Return** O .

Ring Attention



Tree Attention



Summary

All the modern memory optimization/parallelization leverage the mathematical property of softmax, which you learnt today, with some twists to the hardware/software settings.

I also encourage you to work out the numerically stable version of it yourself. You just need keep one more variable the **max** and correct for it.

In case we have more time: speculative decoding

Sample from a smaller model $p(x)$ to guess what a larger target model $q(x)$ wants to say

- If $p(x') > q(x')$, accept the token
- else reject it, and sample again from $p - q$ renormalized

Proof of correctness

We will now show that for any distributions $p(x)$ and $q(x)$, the tokens sampled via *speculative sampling* from $p(x)$ and $q(x)$ are distributed identically to those sampled from $p(x)$ alone. Let β be the acceptance probability (Definition 3.1).

Note that as $p'(x) = \text{norm}(\max(0, p(x) - q(x))) = \frac{p(x) - \min(q(x), p(x))}{\sum_{x'} (p(x') - \min(q(x'), p(x')))} = \frac{p(x) - \min(q(x), p(x))}{1 - \beta}$, the normalizing constant for the adjusted distribution $p'(x)$ is $1 - \beta$, where the last equation follows immediately from Lemma 3.3 and Theorem 3.5.

Now:

$$P(x = x') = P(\text{guess accepted}, x = x') + P(\text{guess rejected}, x = x')$$

Where:

$$P(\text{guess accepted}, x = x') = q(x') \min(1, \frac{p(x')}{q(x')}) = \min(q(x'), p(x'))$$

And:

$$P(\text{guess rejected}, x = x') = (1 - \beta)p'(x') = p(x') - \min(q(x'), p(x'))$$

Overall:

$$P(x = x') = \min(p(x'), q(x')) + p(x') - \min(p(x'), q(x')) = p(x').$$

As desired. \square

Next episode preview? If you still let me present

[AngelBeats!]
次回予告
Next Episode Preview