

Paper

<https://arxiv.org/abs/2205.14135>

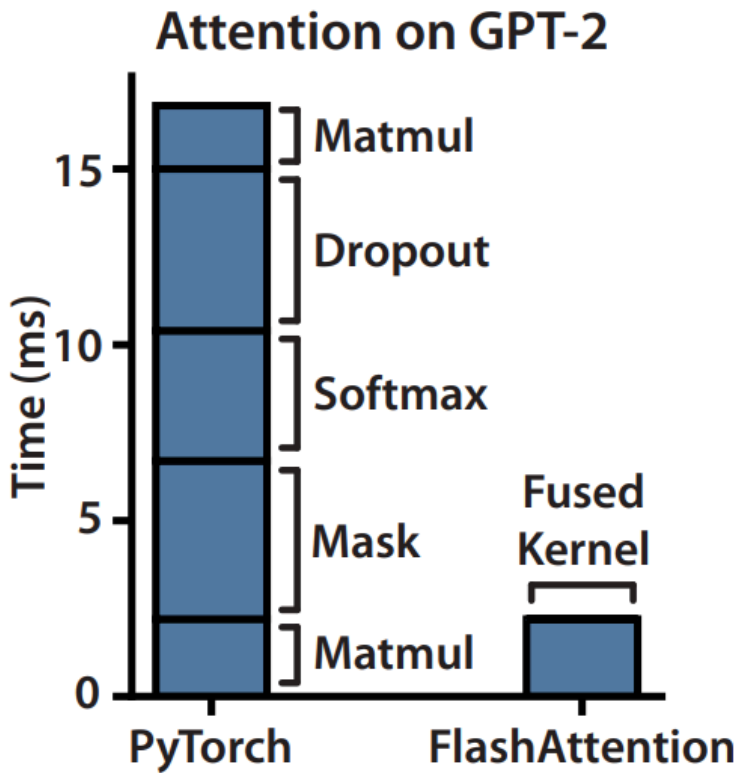
Goal

- Attention is slow because of reads and writes
- Use tiling and recomputation to reduce IO (inputs/outputs)
- Application: transformers that can be faster and have longer context

Background papers

- Approximate attention
 - Sparse
 - Low rank
 - Trade quality for speed
 - Not widely adapted

Tldr



- Attention takes way too much time

Original Attention Algorithm

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 $\mathbf{S} = \mathbf{Q}\mathbf{K}^\top$, 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} .
-

- Given $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ (since attention deals with Query, key value vectors)
 - Want to get an \mathbf{O}
- $\mathbf{S} = \mathbf{Q}\mathbf{K}^\top$
- $\mathbf{P} = \text{Softmax}(\mathbf{S})$
- $\mathbf{O} = \mathbf{P}\mathbf{V}$
- HBM = high bandwidth memory

Main idea of flash attention

- we split the inputs $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ into blocks [**Tiling**]
 - enables us to implement our algorithm in one CUDA kernel
- load them from slow HBM to fast SRAM
- compute the attention output with respect to those blocks [**Recomputation**]
- By scaling the output of each block by the right normalization factor before adding them up, we get the correct result at the end.

Algorithm 1 FLASHATTENTION

Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ 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 $\mathbf{O} = (0)_{N \times d} \in \mathbb{R}^{N \times d}$, $\ell = (0)_N \in \mathbb{R}^N$, $m = (-\infty)_N \in \mathbb{R}^N$ in HBM.
 - 3: Divide \mathbf{Q} into $T_r = \lceil \frac{N}{B_r} \rceil$ blocks $\mathbf{Q}_1, \dots, \mathbf{Q}_{T_r}$ of size $B_r \times d$ each, and divide \mathbf{K}, \mathbf{V} into $T_c = \lceil \frac{N}{B_c} \rceil$ 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.
 - 4: Divide \mathbf{O} into T_r blocks $\mathbf{O}_1, \dots, \mathbf{O}_{T_r}$ 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 $\mathbf{K}_j, \mathbf{V}_j$ from HBM to on-chip SRAM.
 - 7: **for** $1 \leq i \leq T_r$ **do**
 - 8: Load $\mathbf{Q}_i, \mathbf{O}_i, \ell_i, m_i$ from HBM to on-chip SRAM.
 - 9: On chip, compute $\mathbf{S}_{ij} = \mathbf{Q}_i \mathbf{K}_j^T \in \mathbb{R}^{B_r \times B_c}$.
 - 10: On chip, compute $\tilde{m}_{ij} = \text{rowmax}(\mathbf{S}_{ij}) \in \mathbb{R}^{B_r}$, $\tilde{\mathbf{P}}_{ij} = \exp(\mathbf{S}_{ij} - \tilde{m}_{ij}) \in \mathbb{R}^{B_r \times B_c}$ (pointwise), $\tilde{\ell}_{ij} = \text{rowsum}(\tilde{\mathbf{P}}_{ij}) \in \mathbb{R}^{B_r}$.
 - 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 $\mathbf{O}_i \leftarrow \text{diag}(\ell_i^{\text{new}})^{-1} (\text{diag}(\ell_i) e^{m_i - m_i^{\text{new}}} \mathbf{O}_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\mathbf{P}}_{ij} \mathbf{V}_j)$ 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 \mathbf{O} .
-

We show FLASHATTENTION's correctness, runtime, and memory requirement (proof in Appendix C).

Theorem 1. *Algorithm 1 returns $\mathbf{O} = \text{softmax}(\mathbf{Q}\mathbf{K}^T)\mathbf{V}$ with $O(N^2d)$ FLOPs and requires $O(N)$ additional memory beyond inputs and output.*

Tiling

$$m(x) := \max_i x_i, \quad f(x) := [e^{x_1 - m(x)} \quad \dots \quad e^{x_B - m(x)}], \quad \ell(x) := \sum_i f(x)_i, \quad \text{softmax}(x) := \frac{f(x)}{\ell(x)}.$$

For vectors $x^{(1)}, x^{(2)} \in \mathbb{R}^B$, we can decompose the softmax of the concatenated $x = [x^{(1)} \ x^{(2)}] \in \mathbb{R}^{2B}$ as:

$$\begin{aligned} m(x) &= m([x^{(1)} \ x^{(2)}]) = \max(m(x^{(1)}), m(x^{(2)})), \quad f(x) = \begin{bmatrix} e^{m(x^{(1)}) - m(x)} f(x^{(1)}) & e^{m(x^{(2)}) - m(x)} f(x^{(2)}) \end{bmatrix}, \\ \ell(x) &= \ell([x^{(1)} \ x^{(2)}]) = e^{m(x^{(1)}) - m(x)} \ell(x^{(1)}) + e^{m(x^{(2)}) - m(x)} \ell(x^{(2)}), \quad \text{softmax}(x) = \frac{f(x)}{\ell(x)}. \end{aligned}$$

- But what is x^1 and x^2

Recomputation

- Incrementally calculates the soft max instead of storing intermediate matrices
- Uses kernel fusion



How they prove their faster algorithm

- Proved asymptotically
- We discussed this extensively last week

Key idea

- Breaks matrix apart in order to be able to put smaller matrices in SRAM

Useful youtube videos

-  Flash Attention 2.0 with Tri Dao (author)! | Discord server talks
- https://www.youtube.com/live/gMOAud7hZg4?si=1KZU1AV3IMm6_4lt
-  MedAI #54: FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awa...

Code

- <https://github.com/Dao-AILab/flash-attention>

Code notes

-

Meeting Oct 8, 2023

- Theorem 2: if you assume a certain component is $1/n^2$ then it cancels out the original runtime complexity and it makes it linear
- Fun fact: big bird is flash sparse attention
- Why is runtime connected to sequence length?

- Sequence length is the n variable
- Why can't they also tile S and P ?
 - Imo i think they can do this but they also might need to have the matrices together for softmax (because it
- What is gradient checkpointing
 - Can find this in other paper