



# FlashBias: Fast Computation of Attention with Bias

NEURAL INFORMATION PROCESSING SYSTEMS

Haixu Wu, Minghao Guo, Yuezhou Ma, Yuanxu Sun, Jianmin Wang, Wojciech Matusik, Mingsheng Long#

Attention with Bias  $\mathbf{o} = \operatorname{softmax}(\frac{\mathbf{q}\mathbf{k}^{\top}}{\sqrt{C}} + \mathbf{b})\mathbf{v}$ .



Flex/FlashAttention Fails, Try FlashBias!

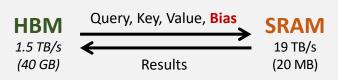
1.5x speedup for Pairformer in Alphafold 3

2x speedup for Swin Transformer v2

Key Challenge: Attention Bias is not sparse.



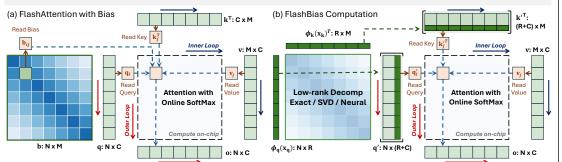
#### Have to load NxN bias matrix from HBM



Get FlashBias kernel at https://github.com/thuml/FlashBias wuhaixu98@gmail.com



# FlashBias: Reach the Theoretical Efficiency Upper Bound



## Why FlashAttention is fast? Underlying low rank assumption

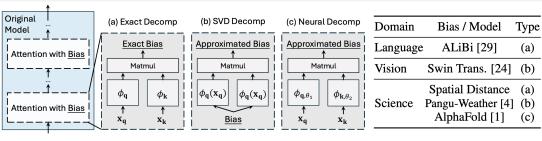
Given Sequence len N, Channel dim C, SRAM size S and  $\underline{C=\alpha N}$ ,  $\underline{S=\beta NC}$ 

- 1) FlashAttention IO Complexity is  $\Theta\left(\left(1+\frac{1}{\alpha}\right)\beta\right)$  smaller than standard attention
- 2) Suppose attention weight is of rank R,  $\alpha \ge \frac{R}{N}$  (determines the optimal speedup)

# > FlashBias based on low-rank compressed sensing theory

- 1) Low-rank Decomp  $\mathbf{b} = f(\mathbf{x_q}, \mathbf{x_k}) = \phi_{\mathbf{q}}(\mathbf{x_q}) \phi_{\mathbf{k}}(\mathbf{x_k})^{\top}, \ \phi_{\mathbf{q}}, \phi_{\mathbf{k}} : \mathbb{R}^{C'} \to \mathbb{R}^{R}.$
- 2) Fast compute  $\operatorname{softmax}(\frac{\mathbf{q}\mathbf{k}^{\top}}{\sqrt{C}} + \mathbf{b})\mathbf{v} = \operatorname{softmax}(\frac{\left[\mathbf{q}|\sqrt{C}\phi_{\mathbf{q}}(\mathbf{x}_{\mathbf{q}})\right]\left[\mathbf{k}|\phi_{\mathbf{k}}(\mathbf{x}_{\mathbf{k}})\right]^{\top}}{\sqrt{C}})\mathbf{v}$

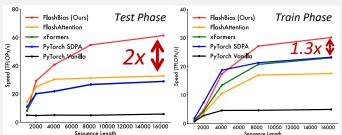
Theorem: FlashBias achieves optimal efficiency guaranteed by compressed sensing theory



Three concrete instantiations for decomposition: Exact, SVD and Neural Decomp

## Speed up without any loss of accuracy

Surpass FlashAttention, PyTorch SDPA, xFormers



#### **GPT-2 with ALiBi** (R=2, Exact decomp)

$$f(\mathbf{x}_{q,i},\mathbf{x}_{k,j})=i-j$$
 ,  $\phi_{\mathbf{q}}(\mathbf{x}_{q,i})=[1,i]$  and  $\phi_{\mathbf{k}}(\mathbf{x}_{k,j})=[-j,1]$ 

Inference: 52.5→91.6TFLOPs/s; Train: 42.7→51.8TFLOPs/s

#### SwinV2-B (R=16, SVD decomp) 0.473s→0.190s

| /lethod                              | Acc@1   | Acc@5              | Time(s)               | Mem(MB)       | ٩                 | 15.5         | 15    | 26       | 15<br>14       |
|--------------------------------------|---------|--------------------|-----------------------|---------------|-------------------|--------------|-------|----------|----------------|
| Official Code<br>Pure FlashAttention |         | 98.232%<br>19.234% |                       | 12829<br>3957 | SVD               | 15.0<br>14.5 | 13    | 13       | 13<br>12<br>11 |
| lashAttention with Bias              |         |                    |                       | 11448         |                   | 14.0         | 10    |          | 11             |
| lexAttention [11]                    | 87.142% | 98.232%            | 2.885                 | 25986         | iginal<br>Matrix  | 15.5         | 15    | 15<br>14 | 15             |
| NT8 PTQ                              | 86.46%  | Arou               | ıd 22% s <sub>i</sub> | peed up       | Origin<br>Bias Ma | 19.0         | 13    | 23       | 13             |
| TashBias (Ours)                      | 87.186% | 98.220%            | 0.190                 | 9429          | 8                 | 14.0         | 11 10 | 11       | 11<br>10       |

## AlphaFold 3 (R=94, Neural decomp) 26.8s→18.2s

|                         | PD     | B ID 7 | wux     | Par l                      |                             |
|-------------------------|--------|--------|---------|----------------------------|-----------------------------|
| Method                  | рТМ↑Т  | ime(s) | Mem(GB) |                            |                             |
| Open-sourced Code       | 0.9500 | 26.85  | 13.62   |                            | ( E) ( E                    |
| FlashAttention w/o Bias | 0.1713 | 8.27   | 12.89   |                            | 12 ES                       |
| FlashAttention w/ Bias  | 0.9500 | 20.39  | 13.62   | C. S.                      | Con A                       |
| FlashBias (Ours)        | 0.9498 | 18.19  | 13.62   | PDB ID: 7wux  Official Pre | PDB ID: 7r6r<br>diction Sne |

Ack: AlphaFold 3 is based on the Protenix from ByteDance

#### Easy to use API: Try FlashBias!

>> from flash\_bias\_triton import flash\_bias\_func

>> output = flash\_bias\_func(q, k, v, q\_bias, k\_bias, mask=None, causal=False, softmax\_scale=1/math.sqrt(headdim))