



FlashBias: Fast Computation of Attention with Bias

Haixu Wu¹, Minghao Guo², Yuezhou Ma¹, Yuanxu Sun¹, Jianmin Wang¹, Wojciech Matusik², Mingsheng Long^{1⊠} ¹School of Software, Tsinghua University, ²MIT CSAIL



Haixu Wu



Minghao Guo



Yuezhou Ma



Yuanxu Sun







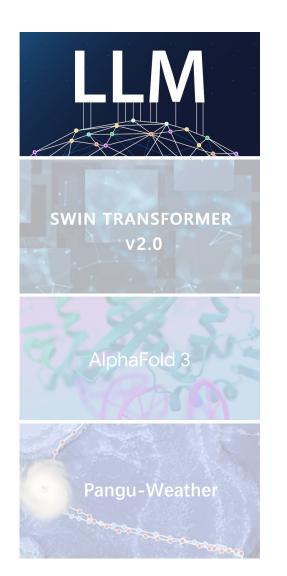
Jianmin Wang Wojciech Matusik Mingsheng Long



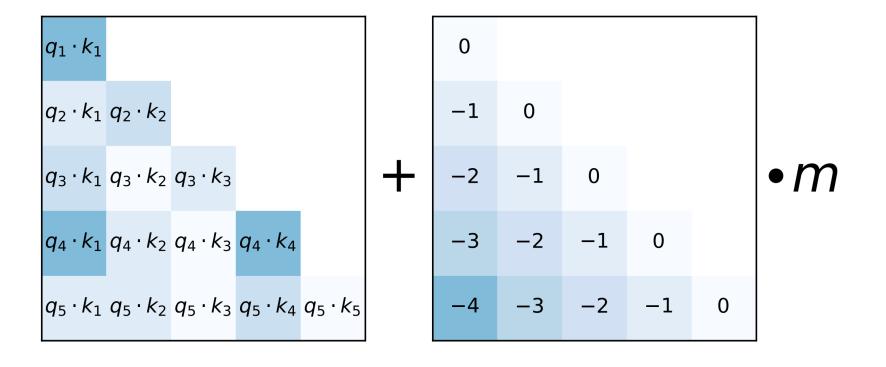
Code Link: https://github.com/thuml/FlashBias

1.5x Speedup for Pairformer in AlphaFold 3; 2x Speedup for Swin Transformer v2. Try FlashBias!

Attention in Advanced Language Models

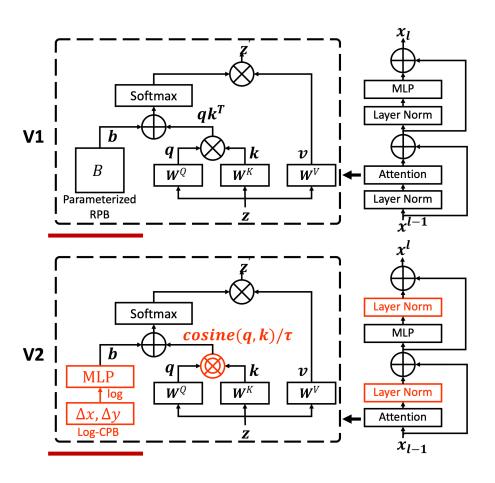


$$\operatorname{softmax}(\mathbf{q}_i\mathbf{K}^\top + m \cdot [-(i-1),...,-2,-1,0])$$



Attention in Advanced Vision Models





$$\operatorname{SoftMax}(QK^T/\sqrt{d} + \underline{B})V,$$

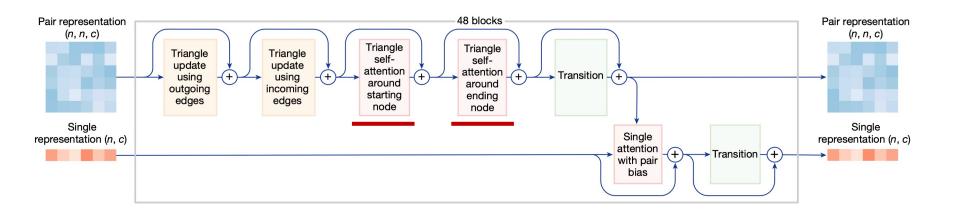
Relative Position Bias

$$\cos(\mathbf{q}_i, \mathbf{k}_j)/\tau + B_{ij}$$

Relative Position Bias

Attention in Advanced Scientific Models





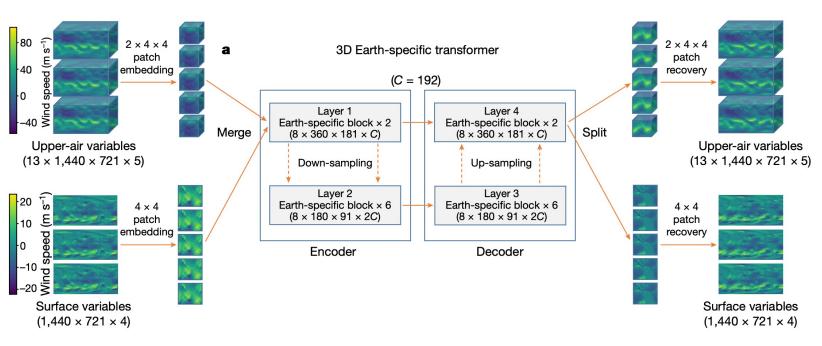
Attention

5:
$$a_{ijk}^h = \operatorname{softmax}_k \left(\frac{1}{\sqrt{c}} \ \mathbf{q}_{ij}^{h^{\top}} \mathbf{k}_{ik}^h + b_{jk}^h \right)$$

6:
$$\mathbf{o}_{ij}^h = \mathbf{g}_{ij}^h \odot \sum_k a_{ijk}^h \mathbf{v}_{ik}^h$$
 Pair Representation Bias

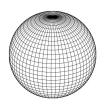
Attention in Advanced Scientific Models





Attention(
$$\mathbf{Q}, \mathbf{K}, \mathbf{V}$$
) = SoftMax($\mathbf{Q}\mathbf{K}^{\top}/\sqrt{D} + \mathbf{B}$) \mathbf{V}

Earth-specific Positional Bias



Attention with Bias

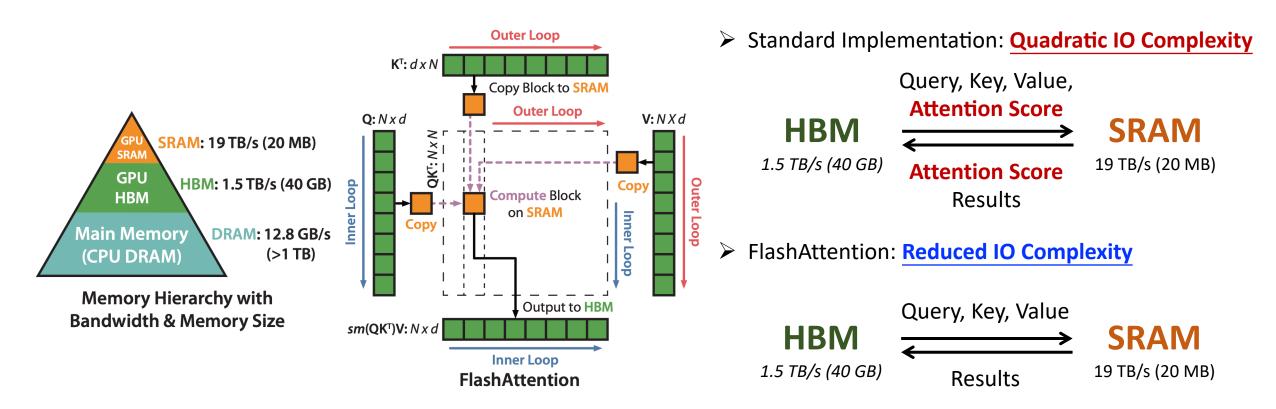
$$\mathbf{o} = \operatorname{softmax}(\frac{\mathbf{q}\mathbf{k}^{\top}}{\sqrt{C}} + \mathbf{b})\mathbf{v}.$$

queries $\mathbf{q} \in \mathbb{R}^{N \times C}$, keys $\mathbf{k} \in \mathbb{R}^{M \times C}$ and values $\mathbf{v} \in \mathbb{R}^{M \times C}$, bias $\mathbf{b} \in \mathbb{R}^{N \times M}$

Introduce prior knowledge to guide attention learning

Vanilla FlashAttention

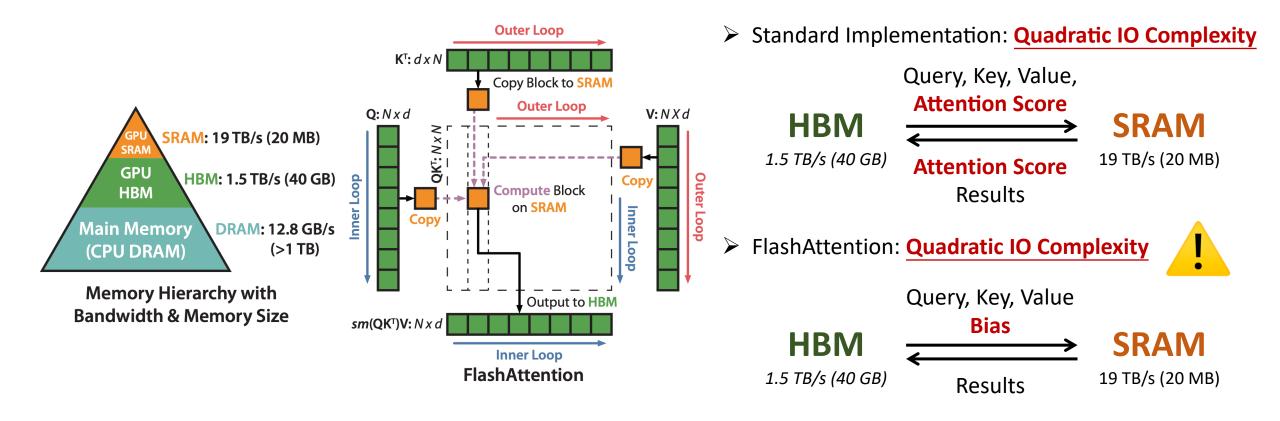
$$\mathbf{o} = \operatorname{softmax}(\frac{\mathbf{q}\mathbf{k}^{\top}}{\sqrt{C}} + \mathbf{b})\mathbf{v}.$$



Dao et al., FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness, NeurIPS 2022 Dao et al., FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning, ICLR 2024

Vanilla FlashAttention Fails

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



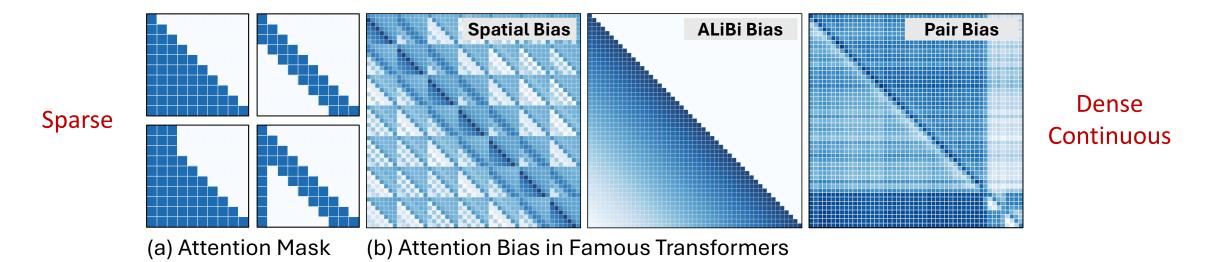
Dao et al., FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness, NeurIPS 2022 Dao et al., FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning, ICLR 2024

Challenge in Optimizing Attention with Bias

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

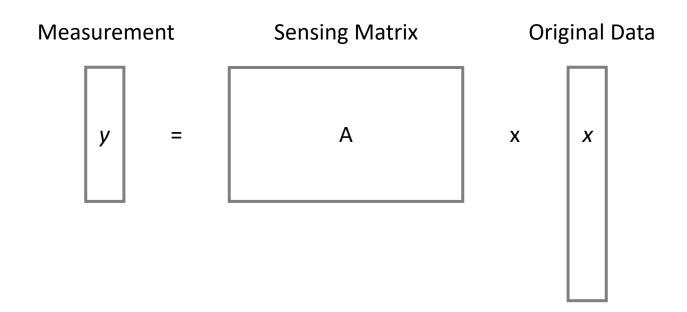
queries $\mathbf{q} \in \mathbb{R}^{N \times C}$, keys $\mathbf{k} \in \mathbb{R}^{M \times C}$ and values $\mathbf{v} \in \mathbb{R}^{M \times C}$, bias $\mathbf{b} \in \mathbb{R}^{N \times M}$

Introduce prior knowledge to guide attention learning



Inevitable IO complexity for loading the dense bias matrix

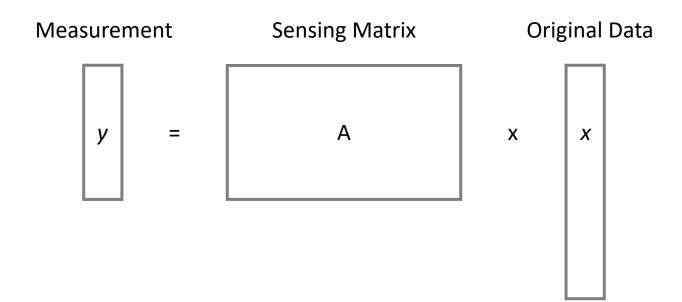
A Typical Compressed Sensing Problem



- > Compressed Sensing: "measurement" (storage) is expensive, but the computation is cheap
- > Attention Computation: IO is slow, but on-chip computation is fast

If we can compress the original data (Bias Matrix), we can reduce the IO complexity.

A Typical Compressed Sensing Problem



Theorem (from Emmanuel J. Candès and Terence Tao) For an NxN dense Matrix with rank-R, the smallest measurement is $\Theta(NR)$.

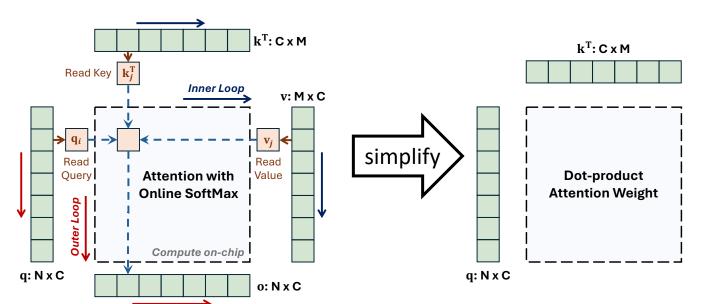
Determined by the rank of matrix

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 dot-product attention weight $\mathbf{s} = \mathbf{q}\mathbf{k}^{\mathrm{T}}$ is of rank R, $\alpha \geq \frac{R}{N}$

The speedup ratio of FlashAttention $\propto \frac{1}{\alpha}$ and $\propto \beta$. β is usually fixed. α determines performance.



Reverse thinking \(\begin{align*} \text{ } \end{align*}

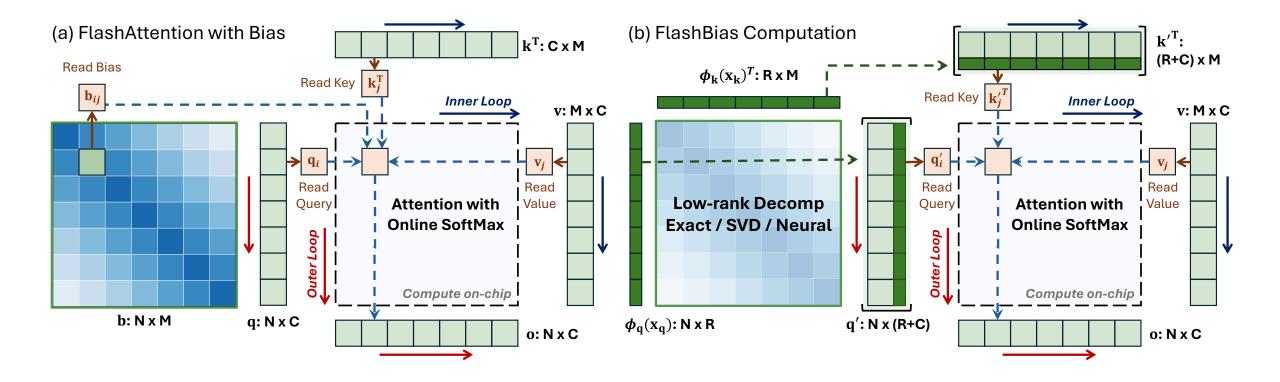


Query and key are from low-rank decomposition of attention score.

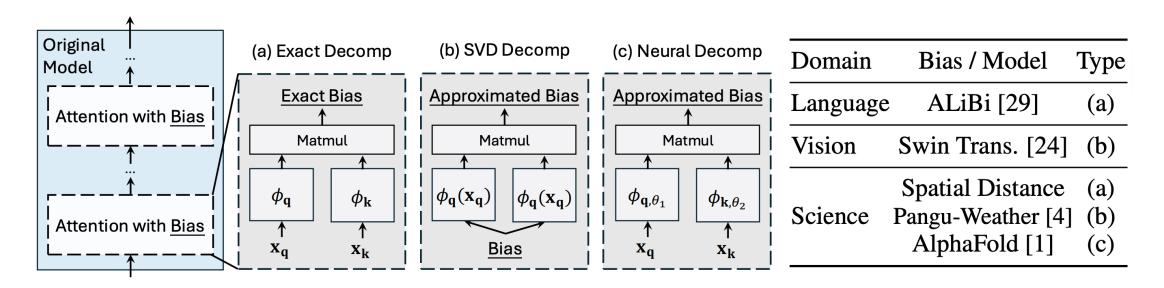
FlashBias: Achieving theoretically optimal complexity

1) Low-rank Decomp
$$\mathbf{b} = f(\mathbf{x}_{\mathbf{q}}, \mathbf{x}_{\mathbf{k}}) = \phi_{\mathbf{q}}(\mathbf{x}_{\mathbf{q}})\phi_{\mathbf{k}}(\mathbf{x}_{\mathbf{k}})^{\top}, \ \phi_{\mathbf{q}}, \phi_{\mathbf{k}} : \mathbb{R}^{C'} \to \mathbb{R}^{R}.$$

2) Fast computation
$$\mathbf{o} = \operatorname{softmax}(\frac{\mathbf{q}\mathbf{k}^{\top}}{\sqrt{C}} + \mathbf{b})\mathbf{v} = \operatorname{softmax}(\frac{\left[\mathbf{q}|\sqrt{C}\phi_{\mathbf{q}}(\mathbf{x}_{\mathbf{q}})\right]\left[\mathbf{k}|\phi_{\mathbf{k}}(\mathbf{x}_{\mathbf{k}})\right]^{\top}}{\sqrt{C}})\mathbf{v}$$



FlashBias: Three concrete instantiations for decomposition



1) Exact Decomp: for some representative bias, such as ALiBi or spatial distance bias.

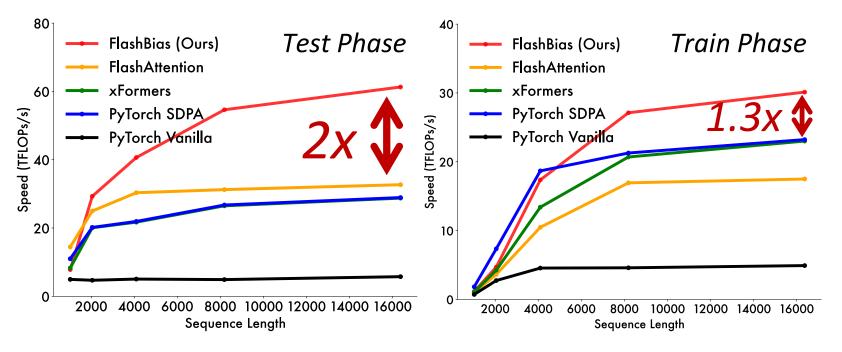
$$f(\mathbf{x_{q,i}},\mathbf{x_{k,j}})=i-j$$
 , $\phi_{\mathbf{q}}(\mathbf{x_{q,i}})=[1,i]$ and $\phi_{\mathbf{k}}(\mathbf{x_{k,j}})=[-j,1]$

- 2) SVD Decomp: when the bias term is learnable model parameters
- 3) Neural Decomp: when the bias term is data dependent

$$\min_{\theta_1, \theta_2} \mathcal{L}(\mathbf{x_q}, \mathbf{x_k}) = \|\widehat{\phi}_{\mathbf{q}, \theta_1}(\mathbf{x_q})\widehat{\phi}_{\mathbf{k}, \theta_2}(\mathbf{x_k})^{\top} - f(\mathbf{x_q}, \mathbf{x_k})\|_2^2.$$
Relative information

Usage and Comparison

- >> 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))

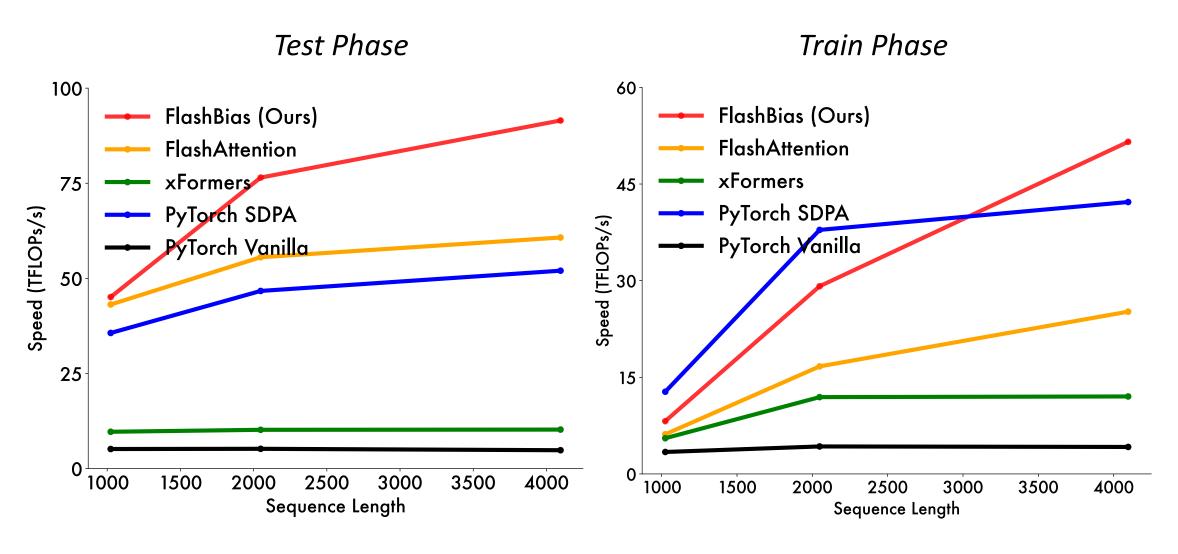


Surpass FlashAttention, PyTorch SDPA, xFormers (bs2-head4-headdim32-noncausal-rank8)

https://github.com/Dao-AlLab/flash-attentionhttps://github.com/facebookresearch/xformers

https://docs.pytorch.org/docs/stable/generated/torch.nn.functional.scaled_dot_product_attention.html

Case 1: GPT-2 with ALiBi Bias (Exact Decomp, Causal Mask, R=2)



batchsize1-head50-headdim32-causal-rank2

Case 2: Swin Transformer V2 (SVD Decomp, R=16)

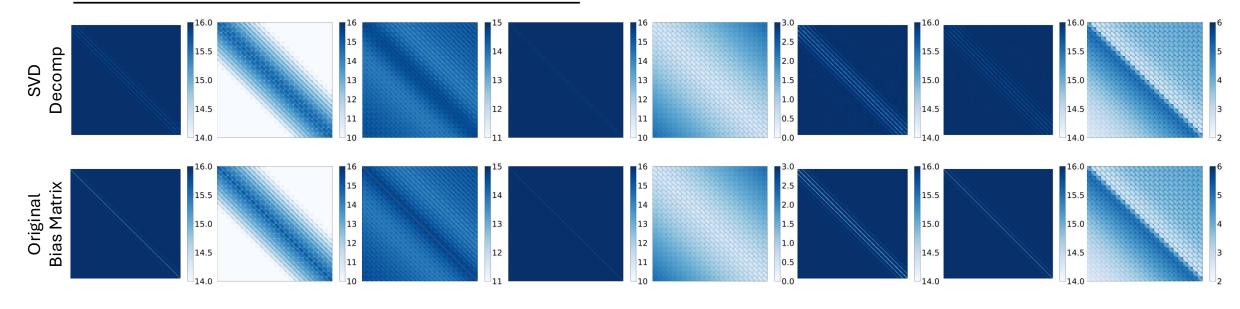
Table 4: Experiment of SwinV2-B on ImageNet-1K. #Time and #Mem correspond to inference efficiency on A100 per batch. Offline calculation of SVD for all biases takes 4.79s.

Method	Acc@1	Acc@5	Time(s)	Mem(MB)
Official Code Pure FlashAttention	87.144% 9.376%	98.232% 19.234%	0.473 0.180	12829 3957
FlashAttention with Bias FlexAttention [11] INT8 PTQ	87.142%	98.232% 98.232% <i>Arour</i>	2.885	11448 25986 peed up
FlashBias (Ours)	87.186%	98.220%	0.190	9429

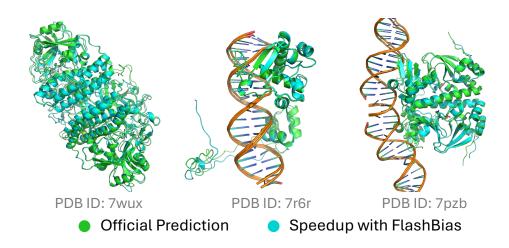
Inference time: $0.473s \rightarrow 0.190s$ (60% reduction)

GPU memory: 12829MB → 9429MB (27% reduction)

2x speedup without any loss of accuracy



Case 3: AlphaFold 3 (Neural Decomp, R=96)

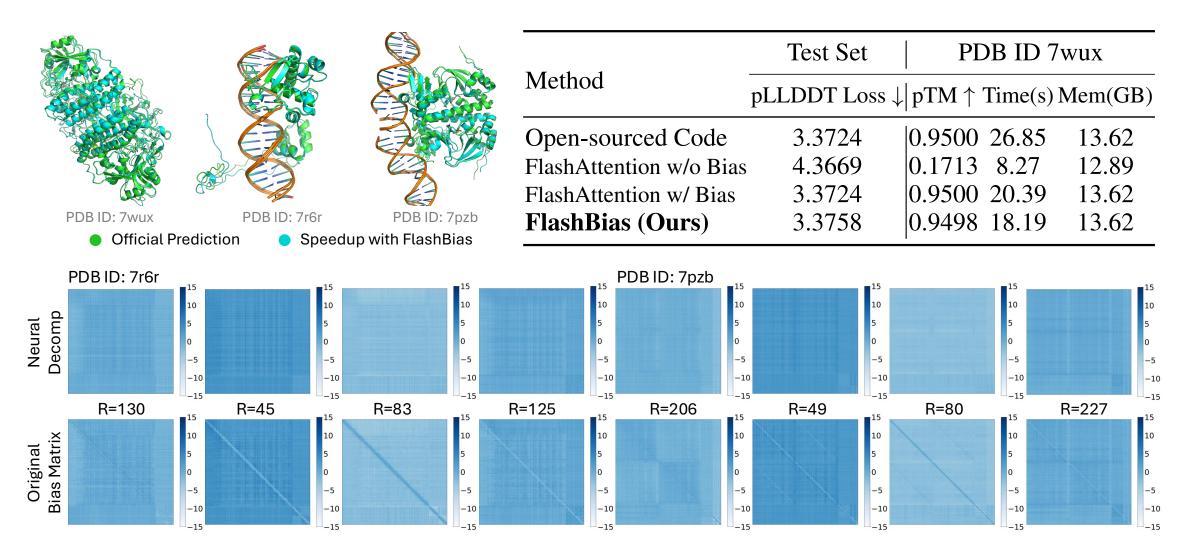


	Test Set	PDB ID 7wux		
Method	pLLDDT Loss ↓	$pTM \uparrow Time(s)$	Mem(GB)	
Open-sourced Code	3.3724	0.9500 26.85	13.62	
FlashAttention w/o Bias	4.3669	0.1713 8.27	12.89	
FlashAttention w/ Bias	3.3724	0.9500 20.39	13.62	
FlashBias (Ours)	3.3758	0.9498 18.19	13.62	

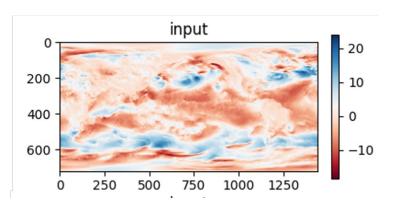
Inference time: 26.85s → 18.19s (32% reduction)

1.5x speedup without any loss of accuracy

Case 3: AlphaFold 3 (Neural Decomp, R=96)



Case 4: Pangu-Weather (SVD Decomp, R=56)



Method	Output Difference	Time(s/100iters)	Mem(MB)
Open-sourced Code	-	98.022	26552
FlashAttention w/o bias	0.0128	74.089	12141
FlashAttention w/ bias	-	79.649	13186
FlashBias (Ours)	0.0003	76.779	12222

Inference time: 98s → 77s (21% reduction)

GPU memory: 26552MB → 12222MB (54% reduction)

speedup without any loss of accuracy

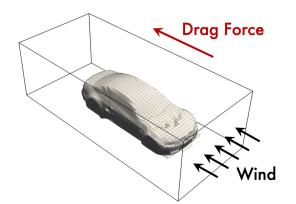
Case 5: Neural Solver with Spatial Dist Bias (Exact Decomp, R=5)

For the spatial distance $f(\mathbf{x}_{\mathbf{q},i},\mathbf{x}_{\mathbf{k},j}) = \|\mathbf{x}_{\mathbf{q},i} - \mathbf{x}_{\mathbf{k},j}\|_2^2$, it can be exactly decomposed as:

$$\phi_{\mathbf{q}}(\mathbf{x}_{\mathbf{q},i}) = [\mathbf{x}_{\mathbf{q},i,0}^{2}, 1, -2\mathbf{x}_{\mathbf{q},i,0}, \mathbf{x}_{\mathbf{q},i,1}^{2}, 1, -2\mathbf{x}_{\mathbf{q},i,1}, \mathbf{x}_{\mathbf{q},i,2}^{2}, 1, -2\mathbf{x}_{\mathbf{q},i,2}],$$

$$\phi_{\mathbf{k}}(\mathbf{x}_{\mathbf{k},j}) = [1, \mathbf{x}_{\mathbf{k},j,0}^{2}, \mathbf{x}_{\mathbf{k},j,0}, 1, \mathbf{x}_{\mathbf{k},j,1}^{2}, \mathbf{x}_{\mathbf{k},j,1}, 1, \mathbf{x}_{\mathbf{k},j,2}^{2}, \mathbf{x}_{\mathbf{k},j,2}].$$
(4)

Based on spherical coordinates, we can further reduce R to 5.



Adaptive distance bias for better geometry information encoding:

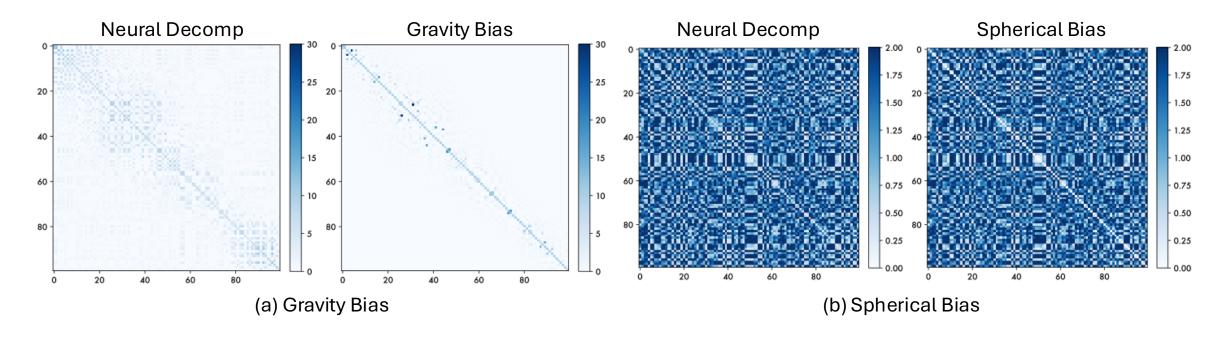
$$f(\mathbf{x}_{\mathbf{q},i}, \mathbf{x}_{\mathbf{k},j}) = \alpha_i \|\mathbf{x}_{\mathbf{q},i} - \mathbf{x}_{\mathbf{k},j}\|_2^2$$

Method (Learnable Bias)	Training Phase			Inference Phase		
	8192	16384	32186	8192	16384	32186
FlashAttention FlexAttention	12.8 / 15.4 Not suppo	OOM orted in curre	OOM nt version	4.54 / 5.46 21.9 / 184.0	15.3 / 21.2 OOM	OOM OOM
FlashBias (Ours)	1.46 / 4.54	2.02 / 14.7	2.97 / 51.1	0.98 / 1.22	1.03 / 3.48	1.13 / 12.7

Generalization to Diverse Bias (Neural Decomp, R=32)

Gravity Bias:
$$f(\mathbf{x}_{\mathbf{q},i}, \mathbf{x}_{\mathbf{k},j}) = \frac{1}{\|\mathbf{x}_{\mathbf{q},i} - \mathbf{x}_{\mathbf{k},j}\|_2^2}$$

$$\textbf{Spherical Bias:} \ f(\mathbf{x}_{\mathbf{q},i},\mathbf{x}_{\mathbf{k},j}) = 2 \cdot \arcsin \left(\sqrt{\sin^2(\frac{\mathbf{x}_{\mathbf{q},i,0} - \mathbf{x}_{\mathbf{k},j,0}}{2}) + \cos \mathbf{x}_{\mathbf{q},i,0} \cos \mathbf{x}_{\mathbf{k},j,0} \sin^2(\frac{\mathbf{x}_{\mathbf{q},i,1} - \mathbf{x}_{\mathbf{k},j,1}}{2})}{2} \right)$$



Extension to Multiplicative Bias

$$\mathbf{o} = \operatorname{softmax}(\frac{\mathbf{q}\mathbf{k}^{\top}}{\sqrt{C}} \odot \mathbf{b})\mathbf{v}$$

queries $\mathbf{q} \in \mathbb{R}^{N \times C}$, keys $\mathbf{k} \in \mathbb{R}^{M \times C}$ and values $\mathbf{v} \in \mathbb{R}^{M \times C}$, bias $\mathbf{b} \in \mathbb{R}^{N \times M}$

Introduce prior knowledge to guide attention learning

FlashBias' extension to multiplicative bias:

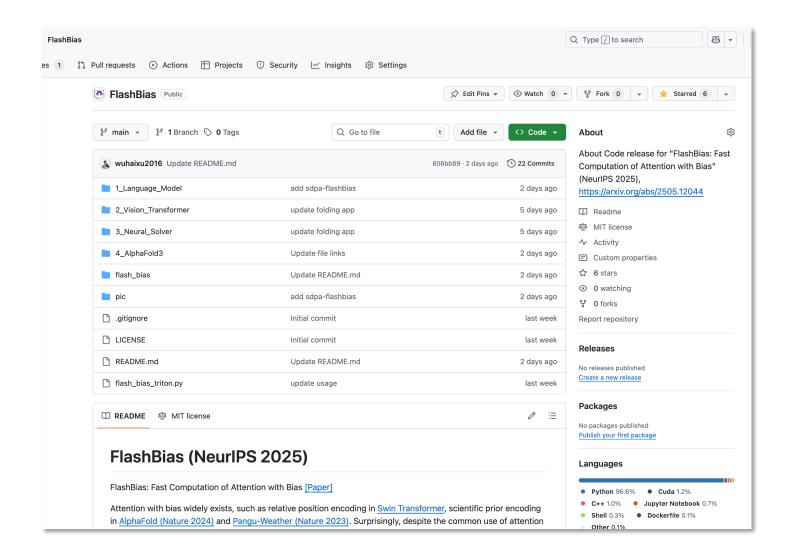
$$\mathbf{o} = \operatorname{softmax}(\frac{\mathbf{q}\mathbf{k}^{\top}}{\sqrt{C}} \odot \mathbf{b})\mathbf{v} = \operatorname{softmax}(\frac{\mathbf{q}'\mathbf{k}'^{\top}}{\sqrt{C}})\mathbf{v},$$

where
$$\mathbf{q}' = [\mathbf{q} \odot \phi_{\mathbf{q},1}, \cdots, \mathbf{q} \odot \phi_{\mathbf{q},R}] \in \mathbb{R}^{N \times CR}, \ \mathbf{k}' = [\mathbf{k} \odot \phi_{\mathbf{k},1}, \cdots, \mathbf{k} \odot \phi_{\mathbf{k},R}] \in \mathbb{R}^{N \times CR}.$$

Example:
$$\mathbf{b}_{ij} = \cos(i-j)$$

$$\phi_{\mathbf{q}}(\mathbf{x}_{\mathbf{q},i}) = [\cos(i), \sin(i)] \in \mathbb{R}^{1 \times 2}, \ \phi_{\mathbf{k}}(\mathbf{x}_{\mathbf{k},j}) = [\cos(j), \sin(j)] \in \mathbb{R}^{1 \times 2}.$$

Open Source





Code is available at

https://github.com/thuml/FlashBias





Thank You!

wuhaixu98@gmail.com

Code Link: https://github.com/thuml/FlashBias

1.5x Speedup for Pairformer in AlphaFold 3; 2x Speedup for Swin Transformer v2.



Try FlashBias!