

Linear Attention for Efficient Bidirectional Sequence Modeling

Arshia Afzal

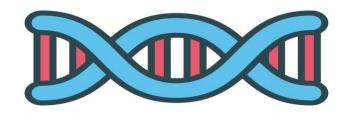
EPFL



Not all Sequences are Causal

Sometimes all tokens are available in tasks like:

- Image classification
- Masked Language Modeling (MLM)
- Diffusion LLMs
- DNA modeling



DNA Modeling

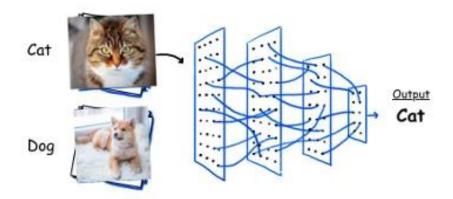
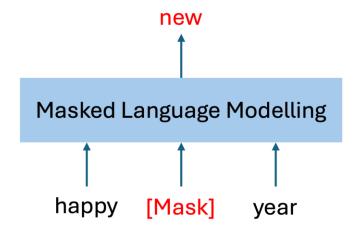


Image classification

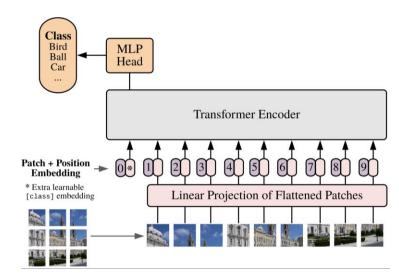


Masked Language Modeling

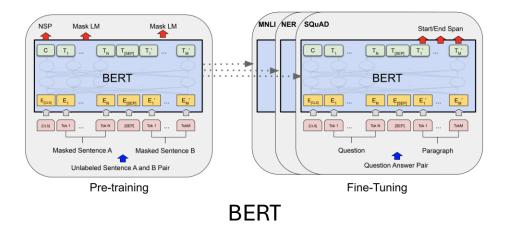
Bi-directional Transformers

Examples of Bi-directional Transformers:

- Vanilla Transformer [Machine Translation]
- Vision Transformers (ViT) [Vision]
- Bidirectional Encoder Representations from Transformers (BERT) [MLM]
- Large Language Diffusion Models (LLaDA)



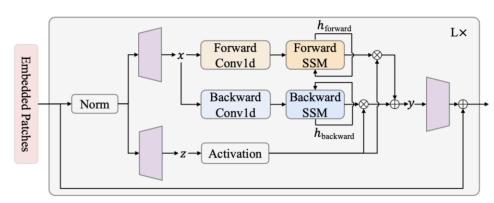
Vision Transformers (ViT)



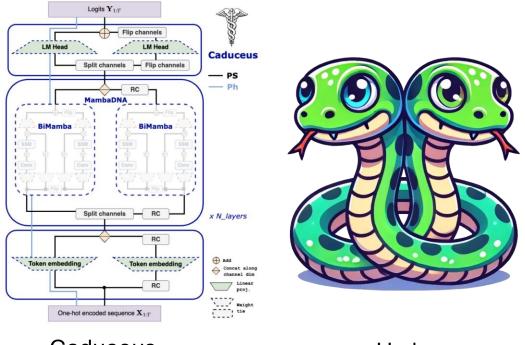
Bi-directional SSMs

Examples of Bi-directional SSMs:

- Vision Mamba (Vim) [Vision]
- Hydra [Vision, MLM]
- Caduceus [DNA]
- Lyra [DNA]



Vision Mamba (Vim)



Caduceus

Hydra

Model	Training Speed	Inference Efficiency
Transformer		×
SSM	×	✓

Model	Training Speed	Inference Efficiency
Transformer	✓	×
SSM	×	✓

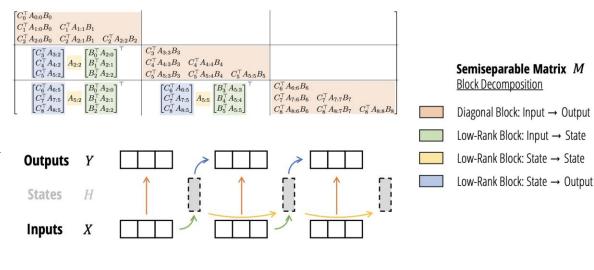
\mathbf{Model}	Top-1	Acc. (%)	Train Time (\downarrow)		
	Small	Base	Small	Base	
$\overline{ m ViT}$	72.2	77.9	×1.00	×1.00	
DeiT	79.8	81.8	$\times 1.00$	$\times 1.00$	
Hydra	78.6	81.0	$\times 2.50$	$\times 2.51$	
Vim	80.3	81.9	$\times 14.95$	$\times 10.86$	

Training Speed on ImageNet-1K Classification

Model	Training Speed	Inference Efficiency
Transformer	✓	×
SSM	×	✓

Model	Top-1	Acc. (%)	Train Time (\downarrow)		
	Small	Base	Small	Base	
$\overline{ m ViT}$	72.2	77.9	×1.00	×1.00	
DeiT	79.8	81.8	×1.00	$\times 1.00$	
Hydra	78.6	81.0	$\times 2.50$	imes 2.51 -	
Vim	80.3	81.9	$\times 14.95$	$\times 10.86$	

Training Speed on ImageNet-1K Classification



Model	Training Speed	Inference Efficiency
Transformer		×
SSM	×	

 $\begin{bmatrix} C_0^{ op} A_{0:0} B_0 \end{bmatrix}$

						$C_1^{\top} A_{1:0} B_0 C_1^{\top} A_{1:0}$				
Model	Top-1	Acc. (%)	Train T	Time (\downarrow)	_	$\begin{bmatrix} C_2^\top A_{2:0} B_0 & C_2^\top A_{2:0} \\ C_3^\top A_{3:2} \\ C_3^\top A_{4:2} \end{bmatrix}_A$	$egin{array}{cccc} A_{2:1}B_1 & C_2^ op A_{2:2}B_1 \\ \hline & B_0^ op A_{2:0} \\ B_1^ op A_{2:1} \end{bmatrix}^ op \end{array}$	$egin{array}{c c} B_2 & & & & & \\ \hline & C_3^ op A_{3:3} B_3 & & & & \\ \hline & C_4^ op A_{4:3} B_3 & C_4^ op A_{4:4} B_4 & & & \\ \hline \end{array}$		+
	Small	Base	Small	Base		$\left\lfloor C_5^{ op} A_{5:2} ight floor$	$\lfloor B_2^{ op} A_{2:2} floor$	$C_5^{\top} A_{5:3} B_3 C_5^{\top} A_{5:4} B_4 C_5^{\top} A_{5:5} B_5$	$C_6^{ op} A_{6:6} B_6$	
ViT	72.2	77.9	×1.00	×1.00	_	$\begin{bmatrix} C_6 & A_{6:5} \\ C_7^{T} A_{7:5} \\ C_8^{T} A_{8:5} \end{bmatrix} A$	$\begin{bmatrix} B_0^\top A_{2:0} \\ B_1^\top A_{2:1} \\ B_2^\top A_{2:2} \end{bmatrix}^\top$	$\begin{bmatrix} C_6^{\top} A_{6:5} \\ C_7^{\top} A_{7:5} \\ C_8^{\top} A_{8:5} \end{bmatrix} \boldsymbol{A_{5:5}} \begin{bmatrix} \boldsymbol{B_3^{\top}} A_{5:3} \\ \boldsymbol{B_4^{\top}} A_{5:4} \\ \boldsymbol{B_5^{\top}} A_{5:5} \end{bmatrix}^{\top}$	$C_7^{\top} A_{7:6} B_6 C_7^{\top} A_{7:7} B_7 $ $C_8^{\top} A_{8:6} B_6 C_8^{\top} A_{8:7} B_7 C_8^{\top} A_{8:8} B_7$	38
DeiT	79.8	81.8	$\times 1.00$	$\times 1.00$						
Hydra	78.6	81.0	$\times 2.50$	$\times 2.51$		Outputs	v F		\ \	
Vim	80.3	81.9	$\times 14.95$	$\times 10.86$		Outputs	1	1		I I
	<u> </u>				_	States	H			

Training Speed on ImageNet-1K Classification

2 SSD Algorithm

XX Flash Linear Attention: Fast Chunkwise Parallel Training

Semiseparable Matrix *M*Block Decomposition

Diagonal Block: Input → Output

Low-Rank Block: Input → State

Low-Rank Block: State → State

Low-Rank Block: State → Output

Why Bi-directional Linear Attention?

Can Bidirectional Linear Attention ...

- Train as fast as Transformers in fully parallel form?
- Be as efficient as SSMs during inference?
- Support chunkwise processing for balanced memory–speed trade-offs?

$$(\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i) = (\mathbf{x}_i \mathbf{W}_{\mathbf{q}}, \ \mathbf{x}_i \mathbf{W}_{\mathbf{k}}, \ \mathbf{x}_i \mathbf{W}_{\mathbf{v}})$$

$$\mathbf{Y} = \operatorname{softmax}(\mathbf{Q}\mathbf{K}^{\top} + \mathbf{M}^{C})\mathbf{V}, \quad \mathbf{y}_{i} = \sum_{j=1}^{i} \frac{\exp(\mathbf{q}_{i}^{\top}\mathbf{k}_{j})}{\sum_{p=1}^{i} \exp(\mathbf{q}_{i}^{\top}\mathbf{k}_{p})} \mathbf{v}_{j}, \quad \mathbf{M}^{C} \in \{-\infty, 0\}^{L \times L}$$

$$(\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i) = (\mathbf{x}_i \mathbf{W}_{\mathbf{q}}, \ \mathbf{x}_i \mathbf{W}_{\mathbf{k}}, \ \mathbf{x}_i \mathbf{W}_{\mathbf{v}})$$

$$\mathbf{Y} = \operatorname{softmax}(\mathbf{Q}\mathbf{K}^{\top} + \mathbf{M}^{C})\mathbf{V}, \quad \mathbf{y}_{i} = \sum_{j=1}^{i} \frac{\exp(\mathbf{q}_{i}^{\top}\mathbf{k}_{j})}{\sum_{p=1}^{i} \exp(\mathbf{q}_{i}^{\top}\mathbf{k}_{p})} \mathbf{v}_{j}, \quad \mathbf{M}^{C} \in \{-\infty, 0\}^{L \times L}$$

$$\mathbf{Y} = \text{SCALE}\left(\mathbf{Q}\mathbf{K}^{\top} \odot \mathbf{M}^{C}\right) \mathbf{V}, \quad \mathbf{S}_{i} = \mathbf{S}_{i-1} + \mathbf{k}_{i} \mathbf{v}_{i}^{\top}, \quad \mathbf{z}_{i} = \mathbf{z}_{i-1} + \mathbf{k}_{i}, \quad \mathbf{y}_{i} = \frac{\mathbf{q}_{i}^{\top} \mathbf{S}_{i}}{\mathbf{q}_{i}^{\top} \mathbf{z}_{i}}, \quad \mathbf{M}^{C} \in \{1, 0\}^{L \times L}$$

$$(\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i) = (\mathbf{x}_i \mathbf{W}_{\mathbf{q}}, \ \mathbf{x}_i \mathbf{W}_{\mathbf{k}}, \ \mathbf{x}_i \mathbf{W}_{\mathbf{v}})$$

$$\mathbf{Y} = \operatorname{softmax}(\mathbf{Q}\mathbf{K}^{\top} + \mathbf{M}^{C})\mathbf{V}, \quad \mathbf{y}_{i} = \sum_{j=1}^{i} \frac{\exp(\mathbf{q}_{i}^{\top}\mathbf{k}_{j})}{\sum_{p=1}^{i} \exp(\mathbf{q}_{i}^{\top}\mathbf{k}_{p})} \mathbf{v}_{j}, \quad \mathbf{M}^{C} \in \{-\infty, 0\}^{L \times L}$$

$$\mathbf{Y} = \text{SCALE}\left(\mathbf{Q}\mathbf{K}^{\top} \odot \mathbf{M}^{C}\right) \mathbf{V}, \quad \mathbf{S}_{i} = \mathbf{S}_{i-1} + \mathbf{k}_{i} \mathbf{v}_{i}^{\top}, \quad \mathbf{z}_{i} = \mathbf{z}_{i-1} + \mathbf{k}_{i}, \quad \mathbf{y}_{i} = \frac{\mathbf{q}_{i}^{\top} \mathbf{S}_{i}}{\mathbf{q}_{i}^{\top} \mathbf{z}_{i}}, \quad \mathbf{M}^{C} \in \{1, 0\}^{L \times L}$$

$$\mathbf{Y} = \text{SCALE}\left(\mathbf{Q}\mathbf{K}^{\top} \odot \mathbf{M}^{C}\right) \mathbf{V}, \quad \mathbf{S}_{i} = \lambda_{i} \mathbf{S}_{i-1} + \mathbf{k}_{i} \mathbf{v}_{i}^{\top}, \quad \mathbf{M}_{ij}^{C} = \begin{cases} \Pi_{k=j+1}^{i} \lambda_{k}, & i \geq j; \\ 0, & i < j. \end{cases}$$

Model	Recurrence	Memory read-out
Linear Attention [48, 47]	$\mathbf{S}_t = \mathbf{S}_{t-1} + oldsymbol{v}_t oldsymbol{k}_t^{^{T}}$	$oldsymbol{o}_t = \mathbf{S}_t oldsymbol{q}_t$
+ Kernel	$\mathbf{S}_t = \mathbf{S}_{t-1} + oldsymbol{v}_t \phi(oldsymbol{k}_t)^{^{T}}$	$oldsymbol{o}_t = \mathbf{S}_t \phi(oldsymbol{q}_t)$
+ Normalization	$\mathbf{S}_t = \mathbf{S}_{t-1} + oldsymbol{v}_t \phi(oldsymbol{k}_t)^{^{\! op}}, \ \ oldsymbol{z}_t = oldsymbol{z}_{t-1} + \phi(oldsymbol{k}_t)$	$oldsymbol{o}_t = \mathbf{S}_t \phi(oldsymbol{q}_t)/(oldsymbol{z}_t^{^T} \phi(oldsymbol{q}_t))$
DeltaNet [101]	$\mathbf{S}_t = \mathbf{S}_{t-1}(\mathbf{I} - eta_t oldsymbol{k}_t oldsymbol{k}_t^{^{T}}) + eta_t oldsymbol{v}_t oldsymbol{k}_t^{^{T}}$	$\boldsymbol{o}_t = \mathbf{S}_t \boldsymbol{q}_t$
Gated RFA [81]	$\mathbf{S}_t = g_t \mathbf{S}_{t-1} + (1-g_t) oldsymbol{v}_t oldsymbol{k}_t^{T}, \ oldsymbol{z}_t = g_t oldsymbol{z}_{t-1} + (1-g_t) oldsymbol{k}_t$	$oldsymbol{o}_t = \mathbf{S}_t oldsymbol{q}_t / (oldsymbol{z}_t^{^{T}} oldsymbol{q}_t)$
S4 [32, 106]	$\mathbf{S}_t = \mathbf{S}_{t-1} \odot \exp(-(\boldsymbol{\alpha} 1^{^T}) \odot \exp(\boldsymbol{A})) + \boldsymbol{B} \odot (\boldsymbol{v}_t 1^{^T})$	$oldsymbol{o}_t = (\mathbf{S}_t \odot oldsymbol{C}) 1 + oldsymbol{d} \odot oldsymbol{v}_t$
ABC [82]	$\mathbf{S}_t^{oldsymbol{k}} = \mathbf{S}_{t-1}^{oldsymbol{k}} + oldsymbol{k}_t oldsymbol{\phi}_t^{T}, \ \ \mathbf{S}_t^{oldsymbol{v}} = \mathbf{S}_{t-1}^{oldsymbol{v}} + oldsymbol{v}_t oldsymbol{\phi}_t^{T}$	$oldsymbol{o}_t = \mathbf{S}_t^{oldsymbol{v}} \operatorname{softmax}\left(\mathbf{S}_t^{oldsymbol{k}} oldsymbol{q}_t ight)$
DFW [63]	$\mathbf{S}_t = \mathbf{S}_{t-1} \odot (oldsymbol{eta}_t oldsymbol{lpha}_t^{^{T}}) + oldsymbol{v}_t oldsymbol{k}_t^{^{T}}$	$oldsymbol{o}_t = \mathbf{S}_t oldsymbol{q}_t$
RetNet [108]	$\mathbf{S}_t = \gamma \mathbf{S}_{t-1} + oldsymbol{v}_t oldsymbol{k}_t^{^T}$	$\boldsymbol{o}_t = \mathbf{S}_t \boldsymbol{q}_t$
Mamba [31]	$\mathbf{S}_t = \mathbf{S}_{t-1} \odot \exp(-(oldsymbol{lpha}_t 1^{^{T}}) \odot \exp(oldsymbol{A})) + (oldsymbol{lpha}_t \odot oldsymbol{v}_t) oldsymbol{k}_t^{^{T}}$	$oldsymbol{o}_t = \mathbf{S}_t oldsymbol{q}_t + oldsymbol{d} \odot oldsymbol{v}_t$
GLA [124]	$\mathbf{S}_t = \mathbf{S}_{t-1} \odot (1oldsymbol{lpha}_t^{^{T}}) + oldsymbol{v}_t oldsymbol{k}_t^{^{T}} = \mathbf{S}_{t-1} \mathrm{Diag}(oldsymbol{lpha}_t) + oldsymbol{v}_t oldsymbol{k}_t^{^{T}}$	$\boldsymbol{o}_t = \mathbf{S}_t \boldsymbol{q}_t$
RWKV-6 [79]	$\mathbf{S}_t = \mathbf{S}_{t-1} \mathrm{Diag}(oldsymbol{lpha}_t) + oldsymbol{v}_t oldsymbol{k}_t^{^{T}}$	$oldsymbol{o}_t = (\mathbf{S}_{t-1} + (oldsymbol{d} \odot oldsymbol{v}_t) oldsymbol{k}_t^{^{T}}) oldsymbol{q}_t$
HGRN-2 [92]	$\mathbf{S}_t = \mathbf{S}_{t-1} \mathrm{Diag}(oldsymbol{lpha}_t) + oldsymbol{v}_t (1 - oldsymbol{lpha}_t)^{^{T}}$	$\boldsymbol{o}_t = \mathbf{S}_t \boldsymbol{q}_t$
mLSTM [9]	$\mathbf{S}_t = f_t \mathbf{S}_{t-1} + i_t oldsymbol{v}_t oldsymbol{k}_t^{^{\intercal}}, \ oldsymbol{z}_t = f_t oldsymbol{z}_{t-1} + i_t oldsymbol{k}_t$	$oldsymbol{o}_t = \mathbf{S}_t oldsymbol{q}_t / \max\{1, oldsymbol{z}_t^{^{T}} oldsymbol{q}_t \}$
Mamba-2 [19]	$\mathbf{S}_t = \gamma_t \mathbf{S}_{t-1} + oldsymbol{v}_t oldsymbol{k}_t^{^{T}}$	$\boldsymbol{o}_t = \mathbf{S}_t \boldsymbol{q}_t$
GSA [131]	$\mathbf{S}_t^{m{k}} = \mathbf{S}_{t-1}^{m{k}} \operatorname{Diag}(m{lpha}_t) + m{k}_t m{\phi}_t^{T}, \ \ \mathbf{S}_t^{m{v}} = \mathbf{S}_{t-1}^{m{v}} \operatorname{Diag}(m{lpha}_t) + m{v}_t m{\phi}_t^{T}$	$oldsymbol{o}_t = \mathbf{S}_t^{oldsymbol{v}} ext{ softmax} \left(\mathbf{S}_t^{oldsymbol{k}} oldsymbol{q}_t ight)$
Gated DeltaNet [125]	$\mathbf{S}_t = \mathbf{S}_{t-1} \left(lpha_t (\mathbf{I} - eta_t oldsymbol{k}_t oldsymbol{k}_t^{T}) ight) + eta_t oldsymbol{v}_t oldsymbol{k}_t^{T}$	$oldsymbol{o}_t = \mathbf{S}_t oldsymbol{q}_t$

[?] Can have bi-directional framework for above linear transformers.

Key Components of Bidirectional Linear Transformers

- Full Linear Attention
- Equal Bi-directional RNN
- Chunkwise Parallel Representation

Key Representation of Bidirectional Linear Transformers

- Full Linear Attention
- Equal Bi-directional RNN
- Chunkwise Parallel Representation

$$\mathbf{Y} = \operatorname{softmax}(\mathbf{Q}\mathbf{K}^{ op})\,\mathbf{V}$$
 Full Softmax Attention

$$\mathbf{Y} = \operatorname{softmax}(\mathbf{Q}\mathbf{K}^{ op})\,\mathbf{V}$$
 Full Softmax Attention

$$\mathbf{Y} = \operatorname{SCALE}\left(\mathbf{Q}\mathbf{K}^{ op}\odot\mathbf{M}\right)\mathbf{V}$$
 Full Linear Attention

Full Linear Attention

$$\mathbf{Y} = \operatorname{softmax}(\mathbf{Q}\mathbf{K}^{\top})\mathbf{V}$$

Full Softmax Attention

$$\mathbf{Y} = \operatorname{Scale}\left(\mathbf{Q}\mathbf{K}^{\top} \odot \mathbf{M}\right) \mathbf{V}$$

$$\mathbf{M}_{ij} = \begin{cases} \Pi_{k=j}^{i-1} \lambda_k, & i > j \\ 1 & i = j \\ \Pi_{k=i+1}^{j} \lambda_k, & i < j. \end{cases} \qquad \mathbf{M}_{ij} = \lambda^{|i-j|}, \qquad \mathbf{M}_{ij} = 1.$$

$$|\mathbf{M}_{ij}| = \lambda^{|i-j|}, \quad |\mathbf{M}_{ij}| = 1$$

LION-Full Linear Attention

$$\mathbf{Y} = \operatorname{softmax}(\mathbf{Q}\mathbf{K}^{\top})\mathbf{V}$$

Full Softmax Attention

$$\mathbf{Y} = \operatorname{Scale}\left(\mathbf{Q}\mathbf{K}^{\top} \odot \mathbf{M}\right) \mathbf{V}$$

$$\mathbf{M}_{ij} = \begin{cases} \Pi_{k=j}^{i-1} \lambda_k, & i > j \\ 1 & i = j \\ \Pi_{k=i+1}^{j} \lambda_k, & i < j. \end{cases} \qquad \mathbf{M}_{ij} = \lambda^{|i-j|}, \qquad \mathbf{M}_{ij} = 1.$$

$$|\mathbf{M}_{ij}| = \lambda^{|i-j|}, \quad |\mathbf{M}_{ij}| = 1$$

$$\begin{pmatrix} \begin{pmatrix} \mathbf{q}_{1}^{\top}\mathbf{k}_{1} & \mathbf{q}_{1}^{\top}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\top}\mathbf{k}_{L} \\ \mathbf{q}_{2}^{\top}\mathbf{k}_{1} & \mathbf{q}_{2}^{\top}\mathbf{k}_{2} & \cdots & \mathbf{q}_{2}^{\top}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\top}\mathbf{k}_{1} & \mathbf{q}_{L}^{\top}\mathbf{k}_{2} & \cdots & \mathbf{q}_{L}^{\top}\mathbf{k}_{L} \end{pmatrix} \underbrace{\begin{pmatrix} \mathbf{1} & \boldsymbol{\lambda}_{2} & \boldsymbol{\lambda}_{2}\boldsymbol{\lambda}_{3} & \cdots & \boldsymbol{\lambda}_{2}\cdots\boldsymbol{\lambda}_{L} \\ \boldsymbol{\lambda}_{1} & \mathbf{1} & \boldsymbol{\lambda}_{3} & \cdots & \boldsymbol{\lambda}_{3}\cdots\boldsymbol{\lambda}_{L} \\ \boldsymbol{\lambda}_{1}\boldsymbol{\lambda}_{2} & \boldsymbol{\lambda}_{2} & \mathbf{1} & \cdots & \boldsymbol{\lambda}_{4}\cdots\boldsymbol{\lambda}_{L} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{\lambda}_{L-1}\cdots\boldsymbol{\lambda}_{1} & \boldsymbol{\lambda}_{L-1}\cdots\boldsymbol{\lambda}_{2} & \boldsymbol{\lambda}_{L-1}\cdots\boldsymbol{\lambda}_{3} & \cdots & \mathbf{1} \end{pmatrix}} \begin{pmatrix} \mathbf{v}_{1}^{\top} \\ \mathbf{v}_{2}^{\top} \\ \mathbf{v}_{3}^{\top} \\ \vdots \\ \mathbf{v}_{L}^{\top} \end{pmatrix} \qquad \mathbf{Causal}$$

$$\mathbf{A} = \mathbf{Q}\mathbf{K}^{\top}$$

$$\mathbf{M}$$

Key Representation of Bidirectional Linear Transformers

- Full Linear Attention
- Equal Bi-directional RNN
- Chunkwise Parallel Representation

Equivalent Bi-directional RNN for Linear Attention

$$\underbrace{\begin{pmatrix} \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{= \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{1} & \frac{1}{2}\mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{2} & \cdots \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \frac{1}{2}\mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{+ \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \frac{1}{2}\mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{+ \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \frac{1}{2}\mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{+ \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \frac{1}{2}\mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{+ \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \frac{1}{2}\mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{+ \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \frac{1}{2}\mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{+ \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{+ \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{+ \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}$$

Equivalent Bi-directional RNN for Linear Attention

Equivalent Bi-directional RNN for Linear Attention

$$\frac{\begin{pmatrix} \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}{\mathbf{A}^{\mathsf{T}}\mathbf{k}_{1}^{\mathsf{T}}\mathbf{k}_{2}^{\mathsf{T}} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}} + \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{\mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1}} + \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{\mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2}^{\mathsf{T}}\mathbf{k}_{2}^{\mathsf{T}}\mathbf{k}_{2}^{\mathsf{T}}\mathbf{k}_{2} \end{pmatrix}} + \underbrace{\begin{pmatrix} \frac{1}{2}\mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{1}^{\mathsf{T}}\mathbf{k}_{L} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{L} \\ \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{1} & \mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{2} & \cdots & \mathbf{q}_{2}^{\mathsf{T}}\mathbf{k}_{L} \end{pmatrix}}_{\mathbf{q}_{L}^{\mathsf{T}}\mathbf{k}_{L}} \\ \mathbf{A}^{B}$$

$$\mathbf{A}^{B}$$

Parallel Form for Diagonal Decay

$$\mathbf{Y} = (\mathbf{C}^{F} + \mathbf{C}^{B})^{-1} \left(\mathbf{A}^{F} \odot \mathbf{M}^{F} + \underbrace{\mathbf{A}^{F} \odot \mathbf{M}^{B}}_{\mathbf{A}^{F} \odot \mathbf{I}} + \underbrace{\mathbf{A}^{B} \odot \mathbf{M}^{F}}_{\mathbf{A}^{B} \odot \mathbf{I}} + \mathbf{A}^{B} \odot \mathbf{M}^{B} - \mathbf{A}^{F} \odot \mathbf{I} - \mathbf{A}^{B} \odot \mathbf{I} \right) \mathbf{V}$$

$$= (\mathbf{C}^{F} + \mathbf{C}^{B})^{-1} \left(\underbrace{(\mathbf{A}^{F} \odot \mathbf{M}^{F}) \mathbf{V}}_{\text{FORWARD}} + \underbrace{(\mathbf{A}^{B} \odot \mathbf{M}^{B}) \mathbf{V}}_{\text{BACKWARD}} \right).$$

The backward (anti-causal) component of full linear attention is equivalent to a recurrent neural network operating in the reverse direction.

$$\underbrace{\begin{pmatrix} \frac{1}{2}\frac{\mathbf{q}_{L}^{\top}\mathbf{k}_{L}}{\mathbf{q}_{L}^{\top}\mathbf{z}_{L}} & \frac{1}{2}\frac{\mathbf{q}_{L-1}^{\top}\mathbf{k}_{L-1}}{\mathbf{q}_{L}^{\top}\mathbf{z}_{L}} & \\ \vdots & \vdots & \ddots & \\ \frac{\mathbf{q}_{L}^{\top}\mathbf{k}_{L}}{\mathbf{q}_{L}^{\top}\mathbf{z}_{L}} & \frac{\mathbf{q}_{L}^{\top}\mathbf{k}_{L-1}}{\mathbf{q}_{L}^{\top}\mathbf{z}_{L}} & \cdots & \frac{1}{2}\frac{\mathbf{q}_{L}^{\top}\mathbf{k}_{L}}{\mathbf{q}_{L}^{\top}\mathbf{z}_{L}} \end{pmatrix}} \underbrace{\begin{pmatrix} \mathbf{1} & & & & \\ \boldsymbol{\lambda}_{L} & \mathbf{1} & & & & \\ \boldsymbol{\lambda}_{L} \boldsymbol{\lambda}_{L-1} & \boldsymbol{\lambda}_{L-1} & \mathbf{1} & & & \\ \vdots & \vdots & \vdots & \ddots & & \\ \boldsymbol{\lambda}_{L} \cdots \boldsymbol{\lambda}_{2} & \boldsymbol{\lambda}_{L} \cdots \boldsymbol{\lambda}_{3} & \boldsymbol{\lambda}_{L} \cdots \boldsymbol{\lambda}_{4} & \cdots & \mathbf{1} \end{pmatrix}}_{F(\mathbf{M}^{B})} \underbrace{\begin{pmatrix} \mathbf{v}_{L}^{\top} \\ \mathbf{v}_{L-1}^{\top} \\ \mathbf{v}_{L-2}^{\top} \\ \vdots \\ \mathbf{v}_{1}^{\top} \end{pmatrix}}_{F(\mathbf{M}^{B})}$$

$$\mathbf{Y} = (\mathbf{C}^F + \mathbf{C}^B)^{-1}(\mathbf{Y}^F + \mathbf{Y}^B), \text{ where}$$

$$\mathbf{Y}^F = (\mathbf{A}^F \odot \mathbf{M}^F)\mathbf{V}, \quad \mathbf{Y}^B = (\mathbf{A}^B \odot \mathbf{M}^B)\mathbf{V} = \text{FLIP}\Big(\big(F(\mathbf{A}^B) \odot F(\mathbf{M}^B)\big)\text{FLIP}(\mathbf{V})\Big).$$

$$\mathbf{S}_{i}^{F/B} = \lambda_{i} \mathbf{S}_{i-1}^{F/B} + \mathbf{k}_{i} \mathbf{v}_{i}^{\mathsf{T}}, \qquad (19a) \qquad \mathbf{y}_{i}^{F/B} = \mathbf{q}_{i}^{\mathsf{T}} \mathbf{S}_{i}^{F/B} - \frac{\mathbf{q}_{i}^{\mathsf{T}} \mathbf{k}_{i}}{2} \mathbf{v}_{i}, \qquad (20a)$$

$$\mathbf{z}_{i}^{F/B} = \lambda_{i} \mathbf{z}_{i-1}^{F/B} + \mathbf{k}_{i}, \qquad (19b)$$

$$c_{i}^{F/B} = \mathbf{q}_{i}^{\mathsf{T}} \mathbf{z}_{i}^{F/B} - \frac{\mathbf{q}_{i}^{\mathsf{T}} \mathbf{k}_{i}}{2}, \qquad (19c)$$

$$OUTPUT: \mathbf{y}_{i} = \frac{\mathbf{y}_{i}^{F} + \mathbf{y}_{i}^{B}}{c_{i}^{F} + c_{i}^{B}} \qquad (20b)$$

$$\mathbf{S}_{i}^{F/B} = \lambda_{i} \mathbf{S}_{i-1}^{F/B} + \mathbf{k}_{i} \mathbf{v}_{i}^{\mathsf{T}}, \qquad (19a) \qquad \mathbf{y}_{i}^{F/B} = \mathbf{q}_{i}^{\mathsf{T}} \mathbf{S}_{i}^{F/B} - \frac{\mathbf{q}_{i}^{\mathsf{T}} \mathbf{k}_{i}}{2} \mathbf{v}_{i}, \qquad (20a)$$

$$\mathbf{z}_{i}^{F/B} = \lambda_{i} \mathbf{z}_{i-1}^{F/B} + \mathbf{k}_{i}, \qquad (19b)$$

$$c_{i}^{F/B} = \mathbf{q}_{i}^{\mathsf{T}} \mathbf{z}_{i}^{F/B} - \frac{\mathbf{q}_{i}^{\mathsf{T}} \mathbf{k}_{i}}{2}, \qquad (19c)$$

$$OUTPUT: \mathbf{y}_{i} = \frac{\mathbf{y}_{i}^{F} + \mathbf{y}_{i}^{B}}{c_{i}^{F} + c_{i}^{B}} \qquad (20b)$$

$$\mathbf{S}_{i}^{F/B} = \lambda_{i} \mathbf{S}_{i-1}^{F/B} + \mathbf{k}_{i} \mathbf{v}_{i}^{\top}, \qquad (19a) \qquad \mathbf{y}_{i}^{F/B} = \mathbf{q}_{i}^{\top} \mathbf{S}_{i}^{F/B} - \frac{\mathbf{q}_{i}^{\top} \mathbf{k}_{i}}{2} \mathbf{v}_{i}, \qquad (20a)$$

$$\mathbf{z}_{i}^{F/B} = \lambda_{i} \mathbf{z}_{i-1}^{F/B} + \mathbf{k}_{i}, \qquad (19b)$$

$$c_{i}^{F/B} = \mathbf{q}_{i}^{\top} \mathbf{z}_{i}^{F/B} - \frac{\mathbf{q}_{i}^{\top} \mathbf{k}_{i}}{2}, \qquad (19c)$$

$$OUTPUT: \mathbf{y}_{i} = \frac{\mathbf{y}_{i}^{F} + \mathbf{y}_{i}^{B}}{c_{i}^{F} + c_{i}^{B}} \qquad (20b)$$

- LION-LIT for $\lambda_i = 1$ which is bi-directional form of Vanilla Linear Transformer [25].
- LION-D for $\lambda_i = \lambda$ fixed decay, and bi-directional form of RetNet (with scaling) [44].
- LION-S, where $\lambda_i = \sigma(\mathbf{W}\mathbf{x}_i + b)$ is the bidirectional extension of GRFA [34] (with shifted SiLU activation) and also inspired by the selectivity of Mamba2.

$$\mathbf{S}_{i}^{F/B} = \lambda_{i} \mathbf{S}_{i-1}^{F/B} + \mathbf{k}_{i} \mathbf{v}_{i}^{\mathsf{T}}, \qquad (19a) \qquad \mathbf{y}_{i}^{F/B} = \mathbf{q}_{i}^{\mathsf{T}} \mathbf{S}_{i}^{F/B} - \frac{\mathbf{q}_{i}^{\mathsf{T}} \mathbf{k}_{i}}{2} \mathbf{v}_{i}, \qquad (20a)$$

$$\mathbf{z}_{i}^{F/B} = \lambda_{i} \mathbf{z}_{i-1}^{F/B} + \mathbf{k}_{i}, \qquad (19b)$$

$$c_{i}^{F/B} = \mathbf{q}_{i}^{\mathsf{T}} \mathbf{z}_{i}^{F/B} - \frac{\mathbf{q}_{i}^{\mathsf{T}} \mathbf{k}_{i}}{2}, \qquad (19c)$$

$$OUTPUT: \mathbf{y}_{i} = \frac{\mathbf{y}_{i}^{F} + \mathbf{y}_{i}^{B}}{c_{i}^{F} + c_{i}^{B}} \qquad (20b)$$

- LION-LIT for $\lambda_i = 1$ which is bi-directional form of Vanilla Linear Transformer [25].
- LION-D for $\lambda_i = \lambda$ fixed decay, and bi-directional form of RetNet (with scaling) [44].
- LION-S, where $\lambda_i = \sigma(\mathbf{W}\mathbf{x}_i + b)$ is the bidirectional extension of GRFA [34] (with shifted SiLU activation) and also inspired by the selectivity of Mamba2.

For diagonal forget gate: $\mathbf{S}_i = \mathbf{\Lambda}_i \mathbf{S}_{i-1} + \mathbf{k}_i \mathbf{v}_i^{ op}$

$$\mathbf{Y} = \mathbf{M} \mathbf{V}$$

$$\mathbf{M}_{ij} = \mathbf{Q}_{i}^{\top} \mathbf{\Lambda}_{0:i}^{\times} (\mathbf{\Lambda}_{0:j}^{\times})^{-1} \mathbf{K}_{j} = \mathbf{C}_{i}^{\top} (\mathbf{\lambda}_{0} \odot \mathbf{\lambda}_{1} \odot \mathbf{\lambda}_{2} \ldots \mathbf{\lambda}_{i}) (\mathbf{\lambda}_{0} \odot \mathbf{\lambda}_{1} \odot \mathbf{\lambda}_{2} \ldots \mathbf{\lambda}_{j})^{-1} \mathbf{K}_{j} = \mathbf{L}_{i} = \mathbf{\lambda}_{0} \odot \mathbf{\lambda}_{1} \odot \cdots \odot \mathbf{\lambda}_{i}$$

$$(\mathbf{Q}_{i} \odot \mathbf{\lambda}_{0} \odot \mathbf{\lambda}_{1} \odot \mathbf{\lambda}_{2} \ldots \mathbf{\lambda}_{i})^{\top} (\mathbf{\lambda}_{0} \odot \mathbf{\lambda}_{1} \odot \mathbf{\lambda}_{2} \ldots \mathbf{\lambda}_{j})^{-1} \odot \mathbf{K}_{j}$$

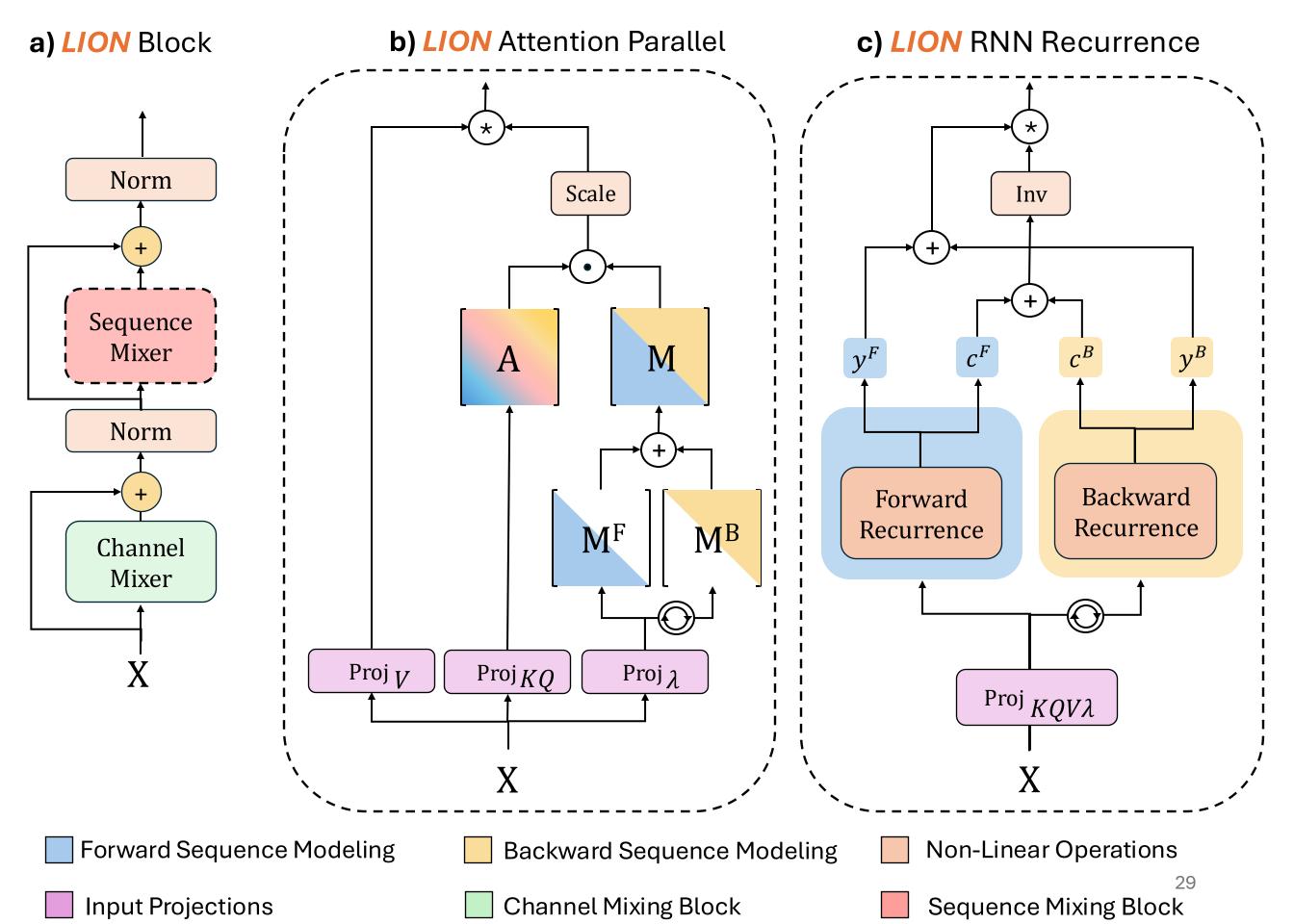
$$\mathbf{U}_{i} = (\mathbf{\lambda}_{0} \odot \mathbf{\lambda}_{1} \odot \cdots \odot \mathbf{\lambda}_{i})^{-1}$$

$$\mathbf{L} = \operatorname{cumprod}(\mathbf{D}), \quad \text{with } \mathbf{D}_{i} = \operatorname{DIAG}(\mathbf{\Lambda}_{i})$$

$$\mathbf{Y} = (\operatorname{TRIL}\left[(\mathbf{Q} \odot \mathbf{L}^{F}) \ (\mathbf{K} \odot \operatorname{Inv}(\mathbf{L}^{F}))^{\top}\right] + \operatorname{TRIU}\left[(\mathbf{Q} \odot \mathbf{L}^{B}) \ (\mathbf{K} \odot \operatorname{Inv}(\mathbf{L}^{B}))^{\top}\right]) \mathbf{V}$$
(81)

$$\mathbf{L}^{F} = \operatorname{cumprod}(\mathbf{D}), \quad \text{with } \mathbf{D}_{i} = \operatorname{DIAG}(\boldsymbol{\Lambda}_{i})$$
(82)

$$\mathbf{L}^{B} = \operatorname{cumprod}(\operatorname{Flip}(\mathbf{D})), \quad \text{with } \mathbf{D}_{i} = \operatorname{DIAG}(\boldsymbol{\Lambda}_{i})$$
(83)



Bi-directional Linear Transformers

Model	Causal Recurrence	LION Bi-directional RNN	LION Full Linear Attention
LinAtt [25]	$egin{aligned} \mathbf{S}_i &= \mathbf{S}_{i-1} + \mathbf{k}_i \mathbf{v}_i^{ op} \ \mathbf{y}_i &= \mathbf{q}_i^{ op} \mathbf{S}_i \end{aligned}$	$egin{aligned} \mathbf{S}_i^{F/B} &= \mathbf{S}_{i-1}^{F/B} + \mathbf{k}_i \mathbf{v}_i^{ op} \ \mathbf{y}_i^{F/B} &= \mathbf{q}_i^{ op} (\mathbf{S}_i^{F/B} - rac{1}{2} \mathbf{k}_i \mathbf{v}_i), \mathbf{y}_i = \mathbf{y}_i^F + \mathbf{y}_i^B \end{aligned}$	$\mathbf{Y} = \mathbf{Q} \mathbf{K}^{ op} \mathbf{V}$
+ Scaling	$egin{aligned} \mathbf{z}_i &= \mathbf{z}_{i-1} + \mathbf{k}_i \ \mathbf{y}_i &= \mathbf{q}_i^ op \mathbf{S}_i/\mathbf{q}_i^ op \mathbf{z}_i \end{aligned}$	$\mathbf{z}_i^{F/B} = \mathbf{z}_{i-1}^{F/B} + \mathbf{k}_i, c_i^{F/B} = \mathbf{q}_i^{\top} (\mathbf{z}_i^{F/B} - \frac{1}{2}\mathbf{k}_i)$ $\mathbf{y}_i^{F/B} = \mathbf{q}_i^{\top} (\mathbf{S}_i^{F/B} - \frac{1}{2}\mathbf{k}_i\mathbf{v}_i), \mathbf{y}_i = \frac{\mathbf{y}_i^F + \mathbf{y}_i^B}{c_i^F + c_i^B} \text{ LION-LIT}$	$\mathbf{Y} = \mathrm{SCALE}(\mathbf{Q}\mathbf{K}^{\top})\mathbf{V}$
RetNet [44]	$egin{aligned} \mathbf{S}_i &= \lambda \mathbf{S}_{i-1} + \mathbf{k}_i \mathbf{v}_i^ op \ \mathbf{y}_i &= \mathbf{q}_i^ op \mathbf{S}_i \end{aligned}$	$egin{aligned} \mathbf{S}_i^{F/B} &= \lambda \mathbf{S}_{i-1}^{F/B} + \mathbf{k}_i \mathbf{v}_i^{ op} \ \mathbf{y}_i^{F/B} &= \mathbf{q}_i^{ op} (\mathbf{S}_i^{F/B} - rac{1}{2} \mathbf{k}_i \mathbf{v}_i), \mathbf{y}_i &= \mathbf{y}_i^F + \mathbf{y}_i^B \end{aligned}$	$egin{aligned} \mathbf{Y} &= (\mathbf{Q}\mathbf{K}^{ op} \odot \mathbf{M}) \mathbf{V}, \ \mathbf{M}_{ij} &= \lambda^{ i-j } \end{aligned}$
+ Scaling	$egin{aligned} \mathbf{z}_i &= \pmb{\lambda} \mathbf{z}_{i-1} + \mathbf{k}_i \ \mathbf{y}_i &= \mathbf{q}_i^ op \mathbf{S}_i / \mathbf{q}_i^ op \mathbf{z}_i \end{aligned}$	$\mathbf{z}_i^{F/B} = \lambda \mathbf{z}_{i-1}^{F/B} + \mathbf{k}_i, c_i^{F/B} = \mathbf{q}_i^{\top} (\mathbf{z}_i^{F/B} - \frac{1}{2}\mathbf{k}_i)$ $\mathbf{y}_i^{F/B} = \mathbf{q}_i^{\top} (\mathbf{S}_i^{F/B} - \frac{1}{2}\mathbf{k}_i \mathbf{v}_i), \mathbf{y}_i = \frac{\mathbf{y}_i^F + \mathbf{y}_i^B}{c_i^F + c_i^B}$ LION-D	$egin{aligned} \mathbf{Y} &= ext{SCALE}(\mathbf{Q}\mathbf{K}^{ op} \odot \mathbf{M})\mathbf{V}, \ \mathbf{M}_{ij} &= \lambda^{ i-j } \end{aligned}$
Gated RFA [34]	$egin{aligned} \mathbf{S}_i &= \mathbf{\sigma}(\mathbf{W}\mathbf{x}_i)\mathbf{S}_{i-1} + \mathbf{k}_i\mathbf{v}_i^{ op} \ \mathbf{z}_i &= \mathbf{\sigma}(\mathbf{W}\mathbf{x}_i)\mathbf{z}_{i-1} + \mathbf{k}_i \ \mathbf{y}_i &= \mathbf{q}_i^{ op}\mathbf{S}_i/\mathbf{q}_i^{ op}\mathbf{z}_i \end{aligned}$	$\begin{aligned} \mathbf{S}_{i}^{F/B} &= \sigma(\mathbf{W}\mathbf{x}_{i})\mathbf{S}_{i-1}^{F/B} + \mathbf{k}_{i}\mathbf{v}_{i}^{\top} \\ \mathbf{y}_{i} &= \mathbf{q}_{i}^{\top}\mathbf{S}_{i} \\ \mathbf{z}_{i}^{F/B} &= \sigma(\mathbf{W}\mathbf{x}_{i})\mathbf{z}_{i-1}^{F/B} + \mathbf{k}_{i}, c_{i}^{F/B} &= \mathbf{q}_{i}^{\top}(\mathbf{z}_{i}^{F/B} - \frac{1}{2}\mathbf{k}_{i}) \\ \mathbf{y}_{i}^{F/B} &= \mathbf{q}_{i}^{\top}(\mathbf{S}_{i}^{F/B} - \frac{1}{2}\mathbf{k}_{i}\mathbf{v}_{i}), \mathbf{y}_{i} &= \frac{\mathbf{y}_{i}^{F} + \mathbf{y}_{i}^{B}}{c_{i}^{F} + c_{i}^{B}} \text{LION-S} \end{aligned}$	$\mathbf{Y} = ext{SCALE}(\mathbf{Q}\mathbf{K}^{ op} \odot \mathbf{M}^*)\mathbf{V}, \ \mathbf{M}^*{}_{ij} = egin{cases} \Pi_i^{j+1}\lambda_k & ext{if } i>j, \ \Pi_{i+1}^{j}\lambda_k & ext{if } i$
Mamba-2 [6]	$egin{aligned} \mathbf{S}_i &= \lambda_i \mathbf{S}_{i-1} + \mathbf{k}_i \ \mathbf{v}_i^{ op} \ \mathbf{y}_i &= \mathbf{q}_i^{ op} \mathbf{S}_i \end{aligned}$	$egin{aligned} \mathbf{S}_i^{F/B} &= \lambda_i \mathbf{S}_{i-1}^{F/B} + \mathbf{k}_i \mathbf{v}_i^{ op} \ \mathbf{y}_i^{F/B} &= \mathbf{q}_i^{ op} (\mathbf{S}_i^{F/B} - rac{1}{2} \mathbf{k}_i \mathbf{v}_i), \mathbf{y}_i = \mathbf{y}_i^F + \mathbf{y}_i^B \end{aligned}$	$\mathbf{Y} = (\mathbf{Q}\mathbf{K}^{ op} \odot \mathbf{M})\mathbf{V} \ \mathbf{M} = \mathbf{M}^*$
RWKV-6[33] Mamba [13] GLA [51]	$egin{aligned} \mathbf{S}_i &= rac{ extsf{Diag}(\lambda_i)}{\mathbf{S}_{i-1}} + \mathbf{k}_i \ \mathbf{v}_i^ op \ \mathbf{y}_i &= \mathbf{q}_i^ op \mathbf{S}_i \end{aligned}$	$egin{aligned} \mathbf{S}_i^{F/B} &= \mathbf{Diag}(\lambda_{\mathbf{i}}) \mathbf{S}_{i-1}^{F/B} + \mathbf{k}_i \mathbf{v}_i^{ op} \ \mathbf{y}_i^{F/B} &= \mathbf{q}_i^{ op} (\mathbf{S}_i^{F/B} - rac{1}{2} \mathbf{k}_i \mathbf{v}_i), \mathbf{y}_i = \mathbf{y}_i^F + \mathbf{y}_i^B \end{aligned}$	$egin{aligned} \mathbf{Y^F} &= \mathbf{Tril}(\mathbf{Q}\odot\mathbf{L^F})(\mathbf{K}\odot{(\mathbf{L^F})}^{-1})\mathbf{V} \ \mathbf{Y^B} &= \mathbf{Triu}(\mathbf{Q}\odot\mathbf{L^B})(\mathbf{K}\odot{(\mathbf{L^B})}^{-1})\mathbf{V} \ \mathbf{Y} &= \mathbf{Y^F} + \mathbf{Y^B} \end{aligned}$
HGRN-2 [36]	$egin{aligned} \mathbf{S}_i &= rac{ extsf{Diag}(\lambda_i)\mathbf{S}_{i-1} + (1-\lambda_i)\mathbf{v}_i^ op \ \mathbf{y}_i &= \mathbf{q}_i^ op \mathbf{S}_i \end{aligned}$	$\begin{aligned} \mathbf{S}_i^{F/B} &= \mathbf{Diag}(\lambda_i) \mathbf{S}_{i-1}^{F/B} + (1 - \lambda_i) \mathbf{v}_i^\top, & \mathbf{k}_i &= (1 - \lambda_i) \\ \mathbf{y}_i^{F/B} &= \mathbf{q}_i^\top (\mathbf{S}_i^{F/B} - \frac{1}{2} \mathbf{k}_i \mathbf{v}_i), & \mathbf{y}_i &= \mathbf{y}_i^F + \mathbf{y}_i^B \end{aligned}$	$egin{aligned} \mathbf{Y^F} &= \mathbf{Tril}(\mathbf{Q}\odot\mathbf{L^F})(\mathbf{K}\odot{(\mathbf{L^F})}^{-1})\mathbf{V} \ \mathbf{Y^B} &= \mathbf{Triu}(\mathbf{Q}\odot\mathbf{L^B})(\mathbf{K}\odot{(\mathbf{L^B})}^{-1})\mathbf{V} \ \mathbf{Y} &= \mathbf{Y^F} + \mathbf{Y^B} \end{aligned}$
xLSTM [3]	$egin{aligned} \mathbf{S}_i &= f_i \mathbf{S}_{i-1} + i_i \mathbf{k}_i \mathbf{v}_i^{ op} \ \mathbf{z}_i &= f_i \mathbf{z}_{i-1} + i_i \mathbf{k}_i \ \mathbf{y}_i &= \mathbf{q}_i^{ op} \mathbf{S}_i / \mathrm{max}(\mathbf{q}_i^{ op} \mathbf{z}_i , 1) \end{aligned}$	$egin{aligned} \mathbf{S}_i^{F/B} &= f_i \mathbf{S}_{i-1}^{F/B} + i_i \mathbf{k}_i \mathbf{v}_i^{ op}, \ \mathbf{y}_i &= \mathbf{q}_i^{ op} \mathbf{S}_i \ \mathbf{z}_i^{F/B} &= f_i \mathbf{z}_{i-1}^{F/B} + i_i \mathbf{k}_i, c_i^{F/B} &= \mathbf{q}_i^{ op} (\mathbf{z}_i^{F/B} - rac{1}{2} \mathbf{k}_i) \ \mathbf{y}_i^{F/B} &= \mathbf{q}_i^{ op} (\mathbf{S}_i^{F/B} - rac{1}{2} \mathbf{k}_i \mathbf{v}_i), \mathbf{y}_i &= rac{\mathbf{y}_i^F + \mathbf{y}_i^B}{\max(c_i^F + c_i^B, 1)} \end{aligned}$	$\mathbf{Y} = ext{SCALE}_{max}(\mathbf{Q}\mathbf{U}^{ op}) \odot \mathbf{M})\mathbf{V}, \ \mathbf{M}_{ij} = egin{cases} \Pi_{i+1}^{j+1}f_k & ext{if } i > j, \ \Pi_{i+1}^{j}f_k & ext{if } i < j, \ 1 & ext{if } i = j, \end{cases}$
DeltaNet [52]	$egin{aligned} \mathbf{S}_i &= (1 - eta_i \mathbf{k}_i \mathbf{k}_i^ op) \mathbf{S}_{i-1} + eta_i \mathbf{k}_i \mathbf{v}_i^ op \ \mathbf{y}_i &= \mathbf{q}_i^ op \mathbf{S}_i \end{aligned}$	$egin{aligned} \mathbf{S}_i^{F/B} &= (1 - eta_i \mathbf{k}_i \mathbf{k}_i^ op) \mathbf{S}_{i-1}^{F/B} + eta_i \mathbf{k}_i \mathbf{v}_i^ op \ \mathbf{y}_i^{F/B} &= \mathbf{q}_i^ op (\mathbf{S}_i^{F/B} - rac{1}{2} \mathbf{k}_i \mathbf{v}_i), \mathbf{y}_i &= \mathbf{y}_i^F + \mathbf{y}_i^B \end{aligned}$	$egin{aligned} \mathbf{Y} &= (\mathbf{Q}\mathbf{K}^{ op})\mathbf{T}\mathbf{V}, \ \mathbf{T} &= rac{1}{2}(\mathbf{T}^F + \mathbf{T}^B), \ \mathbf{T}^F &= (\mathbf{I} + \mathrm{tril}(\mathrm{diag}(eta_i)\mathbf{K}_i\mathbf{K}_i^{ op}, -1))^{-1}\mathrm{diag}(eta_i), \ \mathbf{T}^B &= (\mathbf{I} + \mathrm{triu}(\mathrm{diag}(eta_i)\mathbf{K}_i\mathbf{K}_i^{ op}, -1))^{-1}\mathrm{diag}(eta_i) \end{aligned}$

Key Representation of Bidirectional Linear Transformers

- Full Linear Attention
- Equal Bi-directional RNN
- Chunkwise Parallel Representation

Chunkwise Parallel Form of LION

$q_1^T k_1$	$q_1^T k_2$	$q_1^T k_3$	$q_1^T k_4$	$\boldsymbol{q_1^Tk_5}$	$q_1^T k_6$	$q_1^T k_7$	$q_1^Tk_8$	$q_1^T k_9$
$q_2^T k_1$	$q_2^T k_2$	$q_2^T k_3$	$q_2^T k_4$	$q_2^T k_5$	$q_2^T k_6$	$q_2^T k_7$	$q_2^T k_8$	$q_2^T k_9$
$q_3^T k_1$	$q_3^T k_2$	$q_3^T k_3$	$q_3^T k_4$	$q_3^T k_5$	$q_3^T k_6$	$q_3^T k_7$	$q_3^T k_8$	$q_3^T k_9$
$q_4^T k_1$	$q_4^T k_2$	$q_4^T k_3$	$q_4^T k_4$	$q_4^T k_5$	$q_4^T k_6$	$q_4^T k_7$	$q_4^T k_8$	$_{4}^{T}k_{9}$
$q_5^T k_1$	$q_5^T k_2$	$q_5^T k_3$	$q_5^T k_4$	$q_5^T k_5$	$q_5^T k_6$	$q_5^T k_7$	$q_5^T k_8$	$q_5^T k_9$
$q_6^T k_1$	$q_6^T k_2$	$q_6^T k_3$	$q_6^T k_4$	$q_6^T k_5$	$q_6^T k_6$	$q_6^T k_7$	$q_6^T k_8$	$q_6^T k_9$
$q_7^T k_1$	$q_7^T k_2$	$q_7^T k_3$	$q_7^T k_4$	$q_7^T k_5$	$q_7^T k_6$	$q_7^T k_7$	$q_7^T k_8$	$q_7^T k_9$
$q_8^T k_1$	$q_8^T k_2$	$q_8^T k_3$	$q_8^T k_4$	$q_8^T k_5$	$q_8^T k_6$	$q_8^T k_7$	$q_8^T k_8$	$q_8^T k_9$
$q_9^T k_1$	$q_9^T k_2$	$q_9^T k_3$	$q_9^T k_4$	$q_9^T k_5$	$q_9^T k_6$	$q_9^T k_7$	$q_9^T k_8$	$q_9^T k_9$

 $Q[1]K[3]^{T}$

 $Q[3]K[1]^T$

$$A = QK^T$$

Chunkwise Parallel Form of LION

1	λ	λ^2	λ^3	λ^4	λ^5	λ^6	λ^7	λ^8	
λ	1	λ	λ^2	λ^3	λ^4	λ^5	λ^6	λ^7	$(1)^T$
λ^2	λ	1	λ	λ^2	λ^3	λ^4	λ^5	λ^6	$\int L \left(\frac{1}{L}\right)^T \lambda^{ i-1 }$
λ^3	λ^2	λ	1	λ	λ^2	λ^3	λ^4	λ^5	(-)
λ^4	λ^3	λ^2	λ	1	λ	λ^2		λ^4	
λ^5	λ^4	λ^3	λ^2	λ	1	λ	λ^2	λ^3	
λ^6	λ^5	λ^4	λ^3	λ^2	λ	1	λ	λ^2	
λ^7	λ^6	λ^5	λ^4	λ^3	λ^2	λ	1	λ	
λ^8	λ^7	λ^6	λ^5	λ^4	λ^3	λ^2	λ	1	

$$M = \lambda^{|i-j|}, \qquad L_i = \lambda^i$$

 $L\left(\frac{1}{L}\right)^T \lambda^{|i-j|}$

Chunkwise Parallel Form of LION

$$L_{F}[1] \left(\frac{1}{L_{F}[3]} \right)^{T} \begin{bmatrix} \lambda_{2} & \lambda_{2}\lambda_{3} & \lambda_{2}\lambda_{3}\lambda_{4} & \lambda_{2} \dots \lambda_{5} & \lambda_{2} \dots \lambda_{6} & \lambda_{2} \dots \lambda_{6} & \lambda_{2} \dots \lambda_{8} & \lambda_{2} \dots \lambda_{9} \\ \lambda_{1}\lambda_{2} & \lambda_{2} & 1 & \lambda_{3} & \lambda_{3}\lambda_{4}\lambda_{5} & \lambda_{3}\dots \lambda_{6} & \lambda_{3}\dots \lambda_{7} & \lambda_{3}\dots \lambda_{8} & \lambda_{3}\dots \lambda_{9} \\ \lambda_{1}\lambda_{2}\lambda_{3} & \lambda_{2}\lambda_{3} & \lambda_{3} & 1 & \lambda_{5} & \lambda_{5}\lambda_{6} & \lambda_{5}\lambda_{6} & \lambda_{5}\lambda_{6}\lambda_{7} & \lambda_{5}\dots \lambda_{8} & \lambda_{5}\dots \lambda_{9} \\ \lambda_{1}\dots \lambda_{4} & \lambda_{2}\lambda_{3}\lambda_{4} & \lambda_{3}\lambda_{4} & \lambda_{4}\lambda_{5} & \lambda_{5} & 1 & \lambda_{7} & \lambda_{6}\lambda_{7}\lambda_{8} & \lambda_{7}\lambda_{8}\lambda_{9} \\ \lambda_{1}\dots \lambda_{5} & \lambda_{2}\dots \lambda_{5} & \lambda_{3}\lambda_{4}\lambda_{5} & \lambda_{4}\lambda_{5} & \lambda_{5} & 1 & \lambda_{7} & \lambda_{7}\lambda_{8} & \lambda_{7}\lambda_{8}\lambda_{9} \\ \lambda_{1}\dots \lambda_{5} & \lambda_{2}\dots \lambda_{6} & \lambda_{3}\dots \lambda_{6} & \lambda_{4}\lambda_{5}\lambda_{6} & \lambda_{5}\lambda_{6} & \lambda_{5}\lambda_{6} & \lambda_{6}\lambda_{7} & \lambda_{7}\lambda_{8} & \lambda_{7}\lambda_{8}\lambda_{9} \\ \lambda_{1}\dots \lambda_{7} & \lambda_{2}\dots \lambda_{7} & \lambda_{3}\dots \lambda_{7} & \lambda_{4}\dots \lambda_{7} & \lambda_{5}\lambda_{6}\lambda_{7} & \lambda_{6}\lambda_{7} & \lambda_{7} & 1 & \lambda_{9} \\ \lambda_{1}\dots \lambda_{8} & \lambda_{2}\dots \lambda_{8} & \lambda_{3}\dots \lambda_{8} & \lambda_{4}\dots \lambda_{8} & \lambda_{5}\dots \lambda_{8} & \lambda_{6}\lambda_{7}\lambda_{8} & \lambda_{7}\lambda_{8} & \lambda_{8} & 1 \end{bmatrix}$$

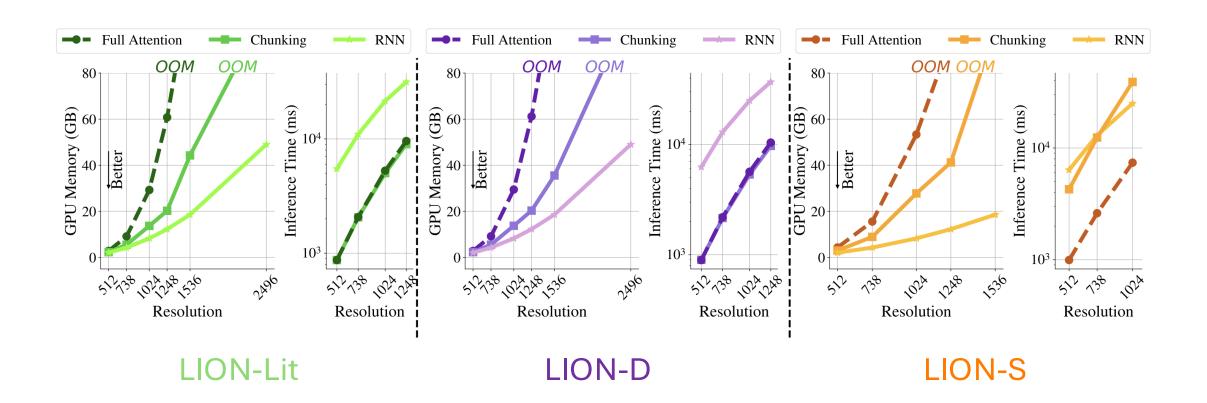
$$\mathbf{M}_{[ij]} = \begin{cases} \mathbf{L}_{[i]}^F \frac{1}{\mathbf{L}_{[i]}^F}^\top & \text{if } i > j, \\ \mathbf{L}_{[j]}^B \frac{1}{\mathbf{L}_{[i]}^B}^\top & \text{if } i < j, \end{cases} \quad \mathbf{L}_{[i]}^F = \operatorname{cumprod}(\lambda)_{iC+1:(i+1)C}, \\ \operatorname{Tril}\left(\mathbf{L}_{[i]}^F \frac{1}{\mathbf{L}_{[i]}^F}^\top\right) + \operatorname{Triu}\left(\mathbf{L}_{[i]}^B \frac{1}{\mathbf{L}_{[i]}^B}^\top\right) - \mathbf{I} \quad \text{if } i = j, \end{cases} \quad \mathbf{L}_{[i]}^F = \operatorname{cumprod}(\operatorname{Flip}(\lambda))_{iC+1:(i+1)C},$$

LION-Chunk

$$\mathbf{A}_{[ij]} = \mathbf{Q}_{[i]} \mathbf{K}_{[j]}^{\top} \odot \mathbf{M}_{[ij]}, \quad \mathbf{C}_{[ij]} = \mathbf{C}_{[ij-1]} + Sum(\mathbf{A}_{[ij]}), \quad \mathbf{S}_{[ij]} = \mathbf{S}_{[ij-1]} + \mathbf{A}_{[ij]} \mathbf{V}_{[j]}, \quad \mathbf{Y}_{[i]} = \frac{\mathbf{S}_{[iN]}}{\mathbf{C}_{[iN]}} \quad (21)$$

LION-Chunk

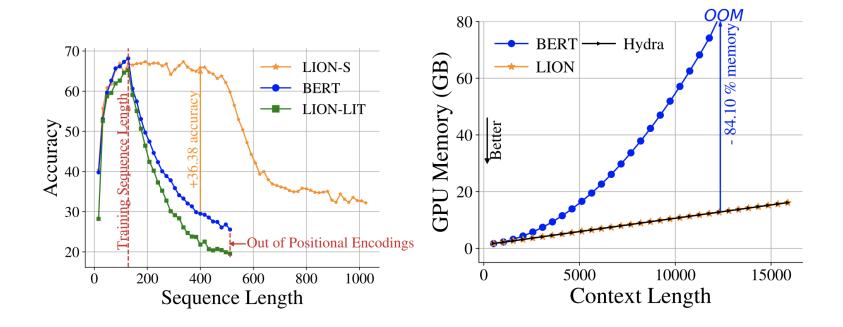
$$\mathbf{A}_{[ij]} = \mathbf{Q}_{[i]} \mathbf{K}_{[j]}^{\top} \odot \mathbf{M}_{[ij]}, \quad \mathbf{C}_{[ij]} = \mathbf{C}_{[ij-1]} + Sum(\mathbf{A}_{[ij]}), \quad \mathbf{S}_{[ij]} = \mathbf{S}_{[ij-1]} + \mathbf{A}_{[ij]} \mathbf{V}_{[j]}, \quad \mathbf{Y}_{[i]} = \frac{\mathbf{S}_{[iN]}}{\mathbf{C}_{[iN]}} \quad (21)$$



Experiments

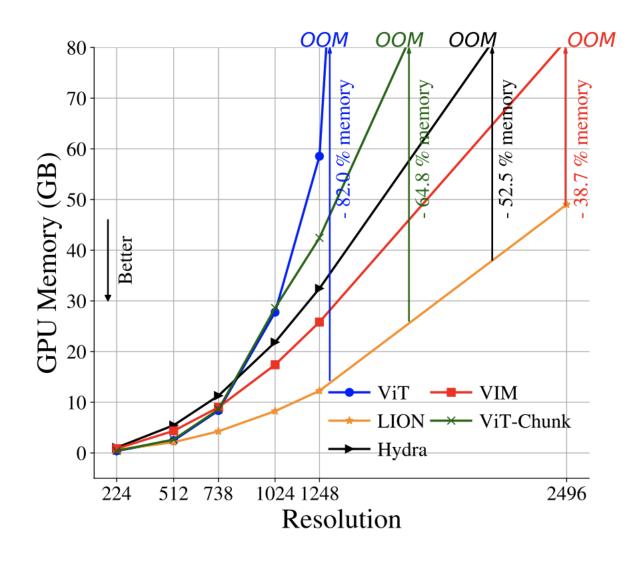
Masked Language Modeling

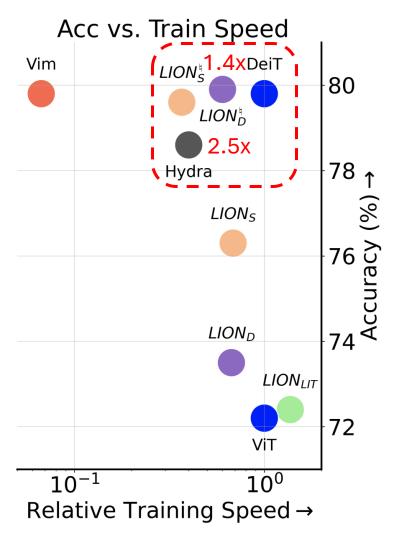
Model	MLM Acc.	GLUE	Train. time
BERT	69.88	82.95	×1
Hydra	71.18	<u>81.77</u>	$\times 3.13$
LION-LIT	67.11	80.76	× 0.95
LION-D	68.64	81.34	$\times \underline{1.10}$
LION-S	69.16	81.58	$\times 1.32$



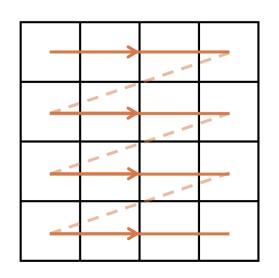
Experiments

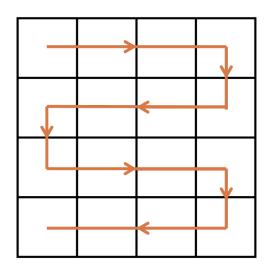
Image Classification

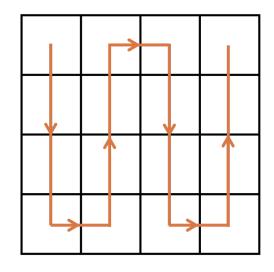




Linear Vision Transformers Require Multiple Scans







LION-LIT	22M	86M	72.4	74.7	×0.74	×0.73
LION-D	22M	86M	73.5	77.8	×1.49	$\times 1.39$
Lion-d ^{\(\beta\)}	22M	86M	79.9	80.2	×1.66	$\times 1.48$
LION-S	22M	86M	74.0	76.3	×2.03	$\times 1.46$
Lion-s ^{\bar{\bar{\bar{\bar{\bar{\bar{\bar{}	22M	86M	79.6	79.9	×2.72	×1.68

Future Directions

- Using LION Masks as 2D Positional Embeddings in Softmax Transformers
- Leveraging LION to Accelerate Hydra's Training
- Extending LION to Other Domains: DNA and Diffusion Language Models

Authors and Links









Arshia Afzal

Elias Abad Rocamora

Leyla Naz Candogan

Pol Puigdemont Plana







Francesco Tonin



Mahsa Shoaran



Volkan Cevher







