

Sequence Models

Ramin Hasani

ACDL 2024

Agenda

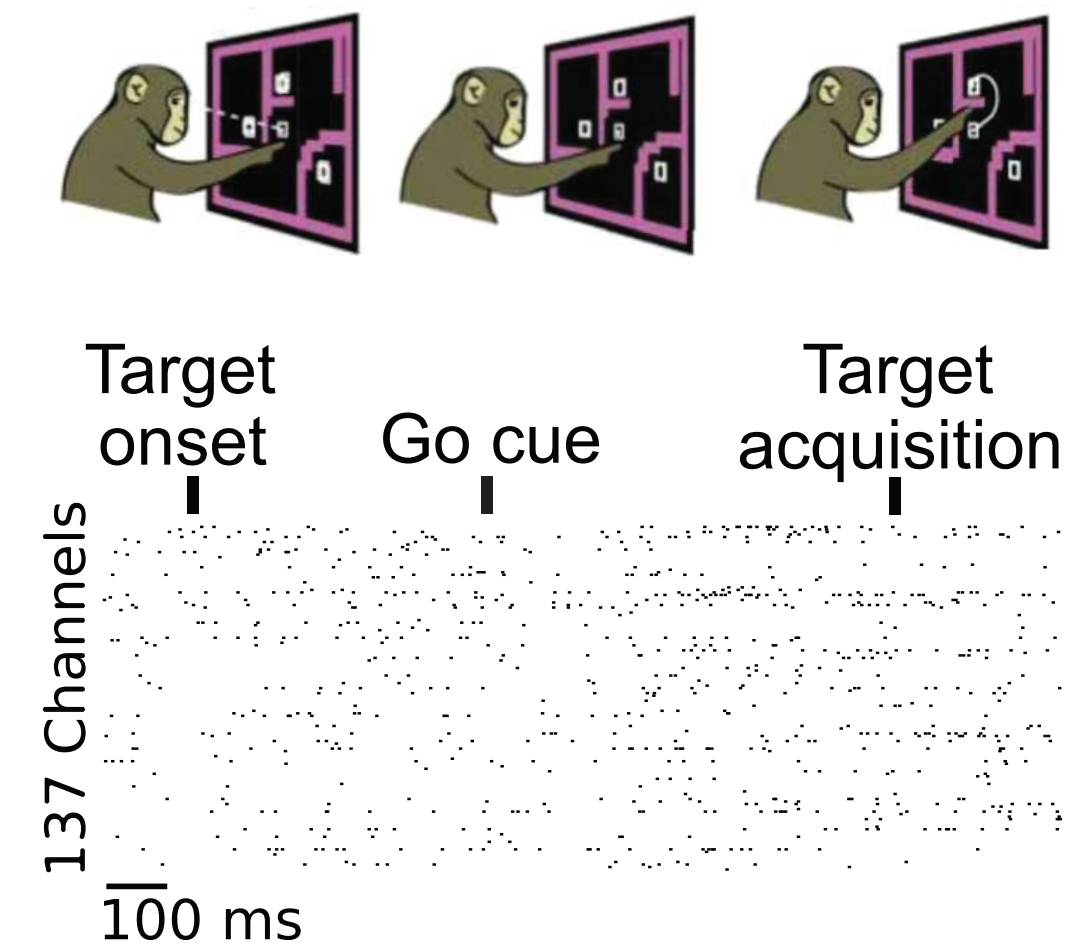
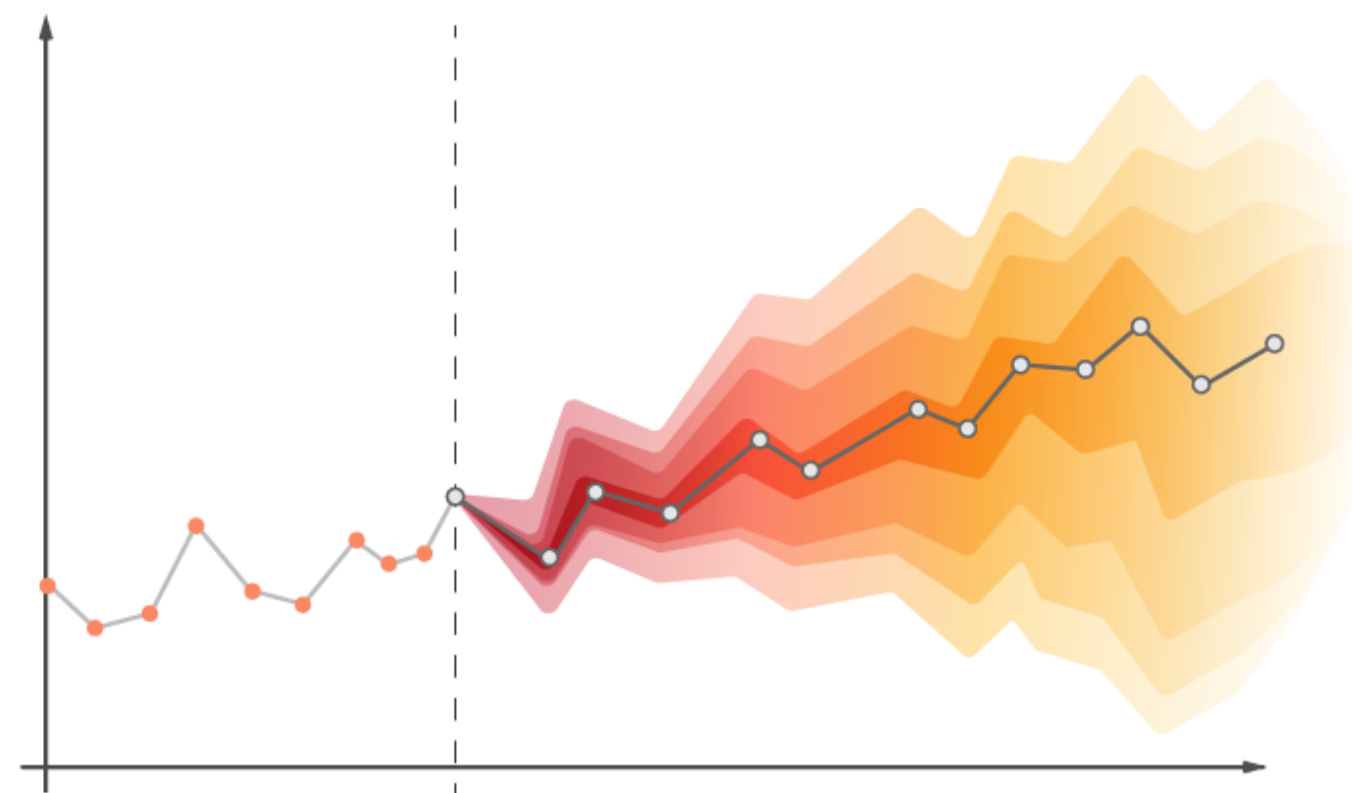
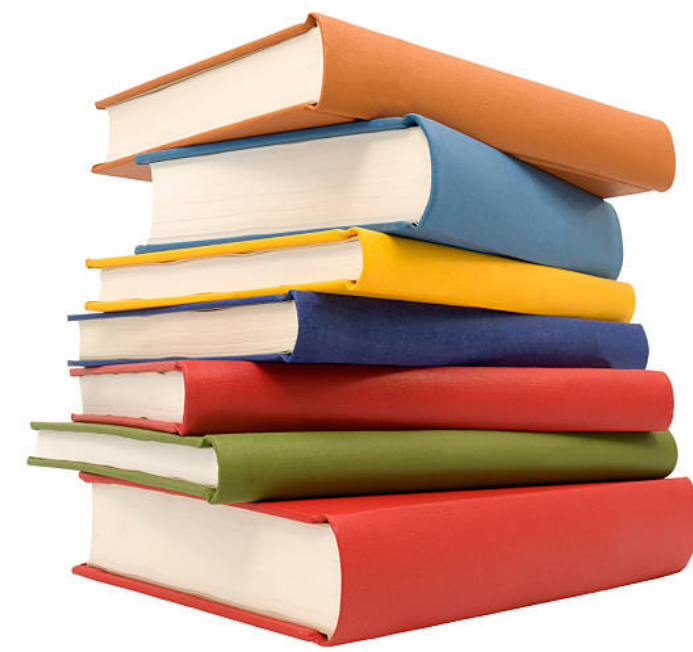
- Introduction, motivation, prior approaches
- Linear state space models (SSMs) overview
- S4, convolutions, parameterization
- S5, diagonalization, parallel scans
- S6/Mamba/Hawk/Griffin, data-dependent dynamics
- Challenges and opportunities

Agenda

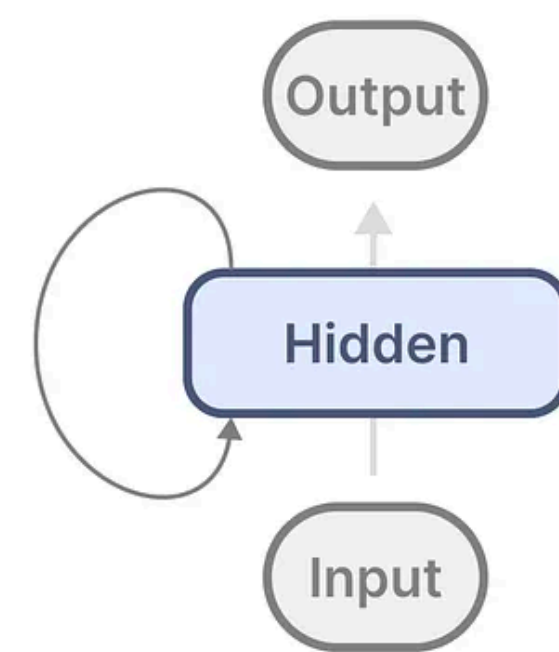
- **Introduction, motivation, prior approaches**
- Linear state space models (SSMs) overview
- S4, convolutions, parameterization
- S5, diagonalization, parallel scans
- S6/Mamba/Hawk/Griffin, data-dependent dynamics
- Challenges and opportunities

Motivation: Efficiently modeling long sequences

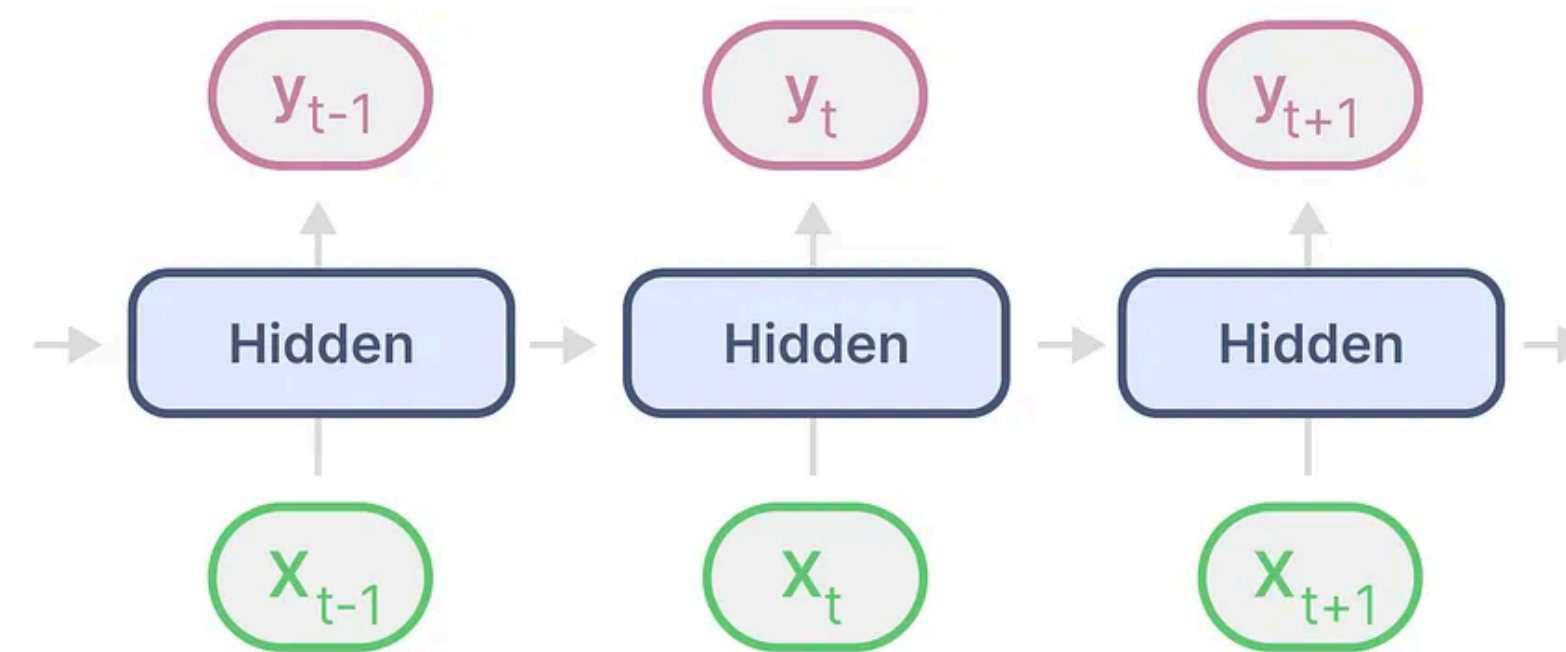
Applications: text, audio, forecasting, neuroscience, images, videos



Nonlinear Recurrent Neural Networks

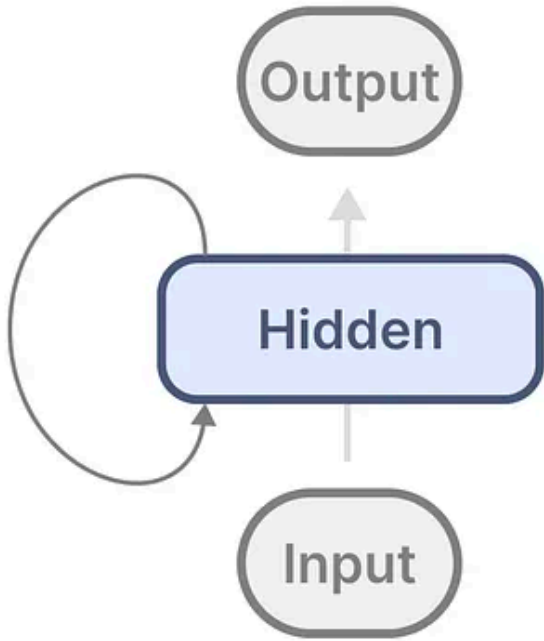


RNN

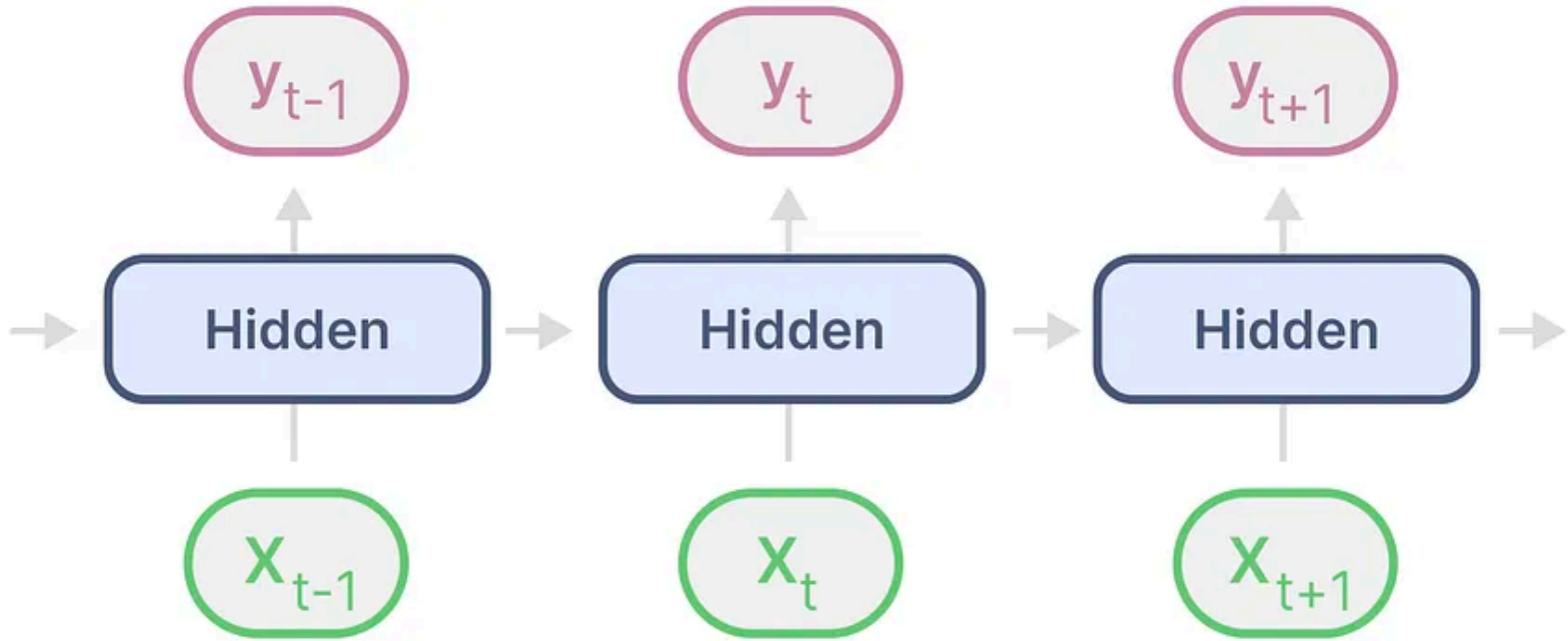


RNN
(Unfolded)

Nonlinear Recurrent Neural Networks



RNN

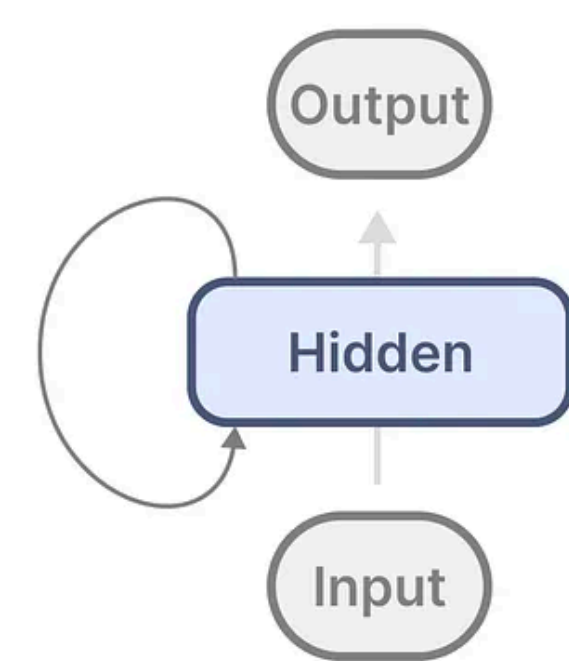


RNN
(Unfolded)

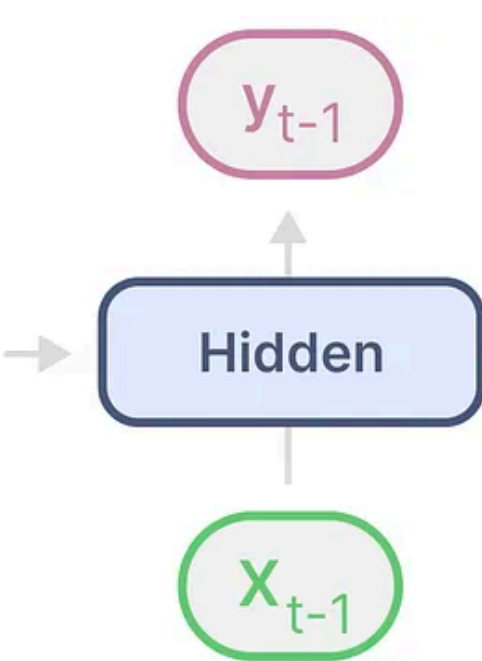
	Recurrent Neural Networks (RNNs)
Parallelizable training	X

Inherently sequential forward and backward pass

Nonlinear Recurrent Neural Networks



RNN



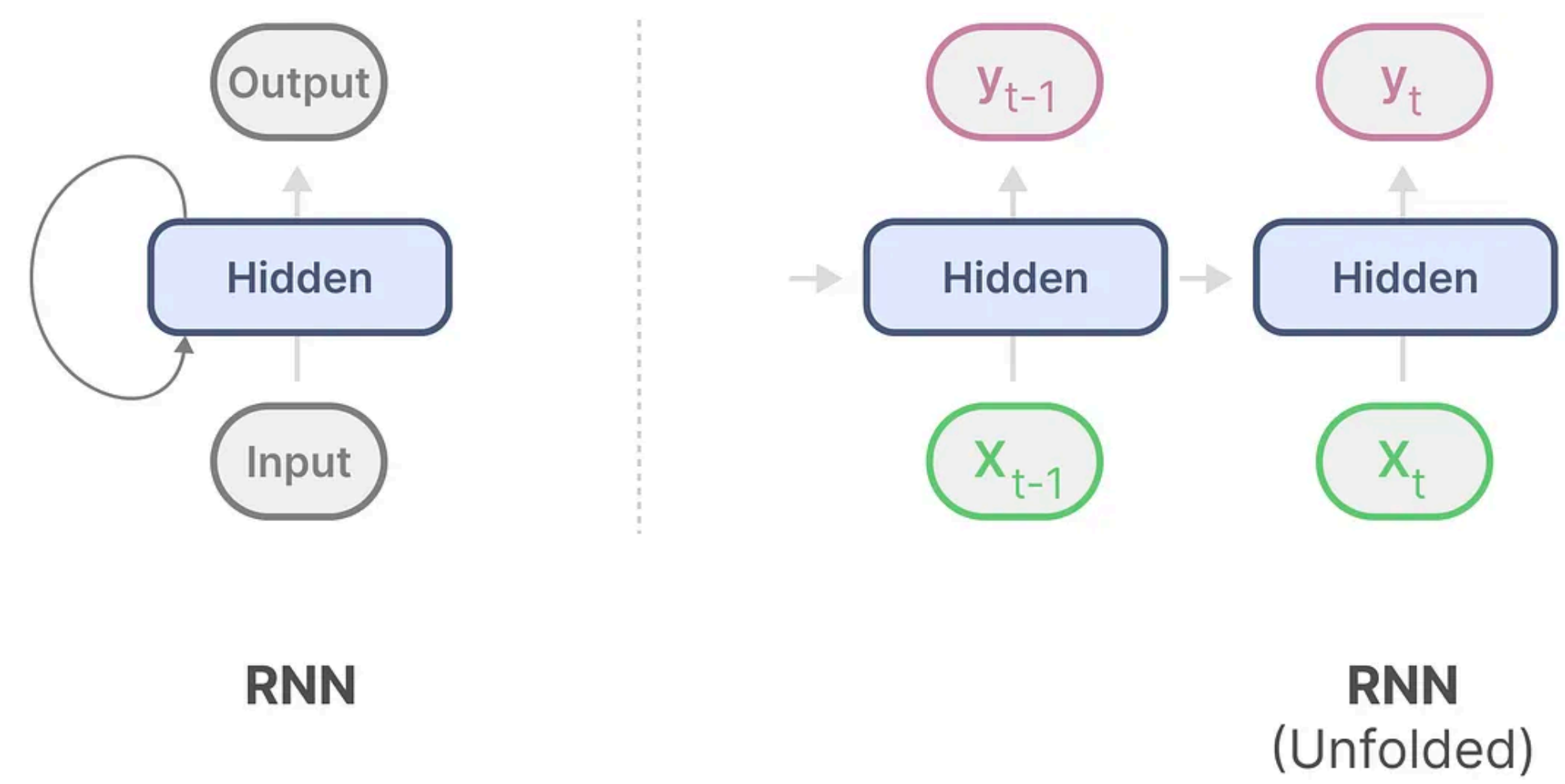
RNN
(Unfolded)

	Recurrent Neural Networks (RNNs)
Parallelizable training	X
Fast autoregressive generation	✓

Inherently sequential forward and backward pass

Constant time and space required to perform single step of generation

Nonlinear Recurrent Neural Networks

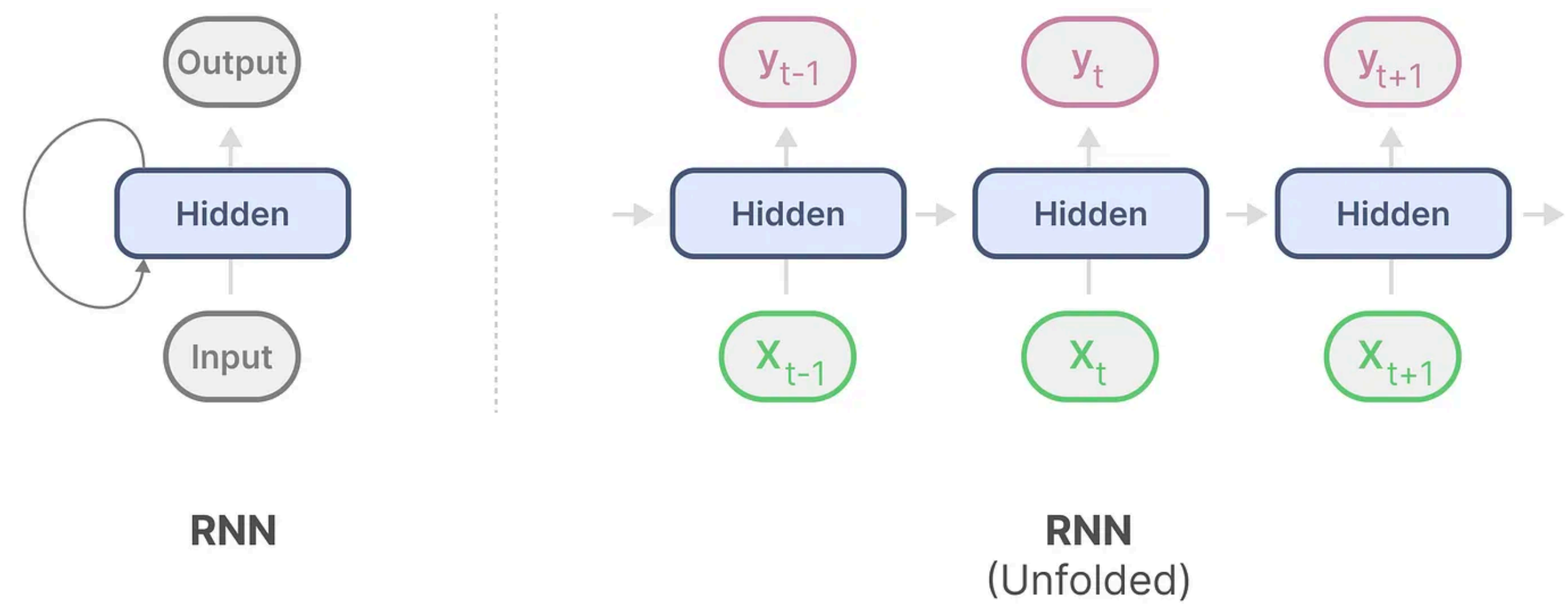


	Recurrent Neural Networks (RNNs)
Parallelizable training	X
Fast autoregressive generation	✓

Inherently sequential forward and backward pass

Constant time and space required to perform single step of generation

Nonlinear Recurrent Neural Networks

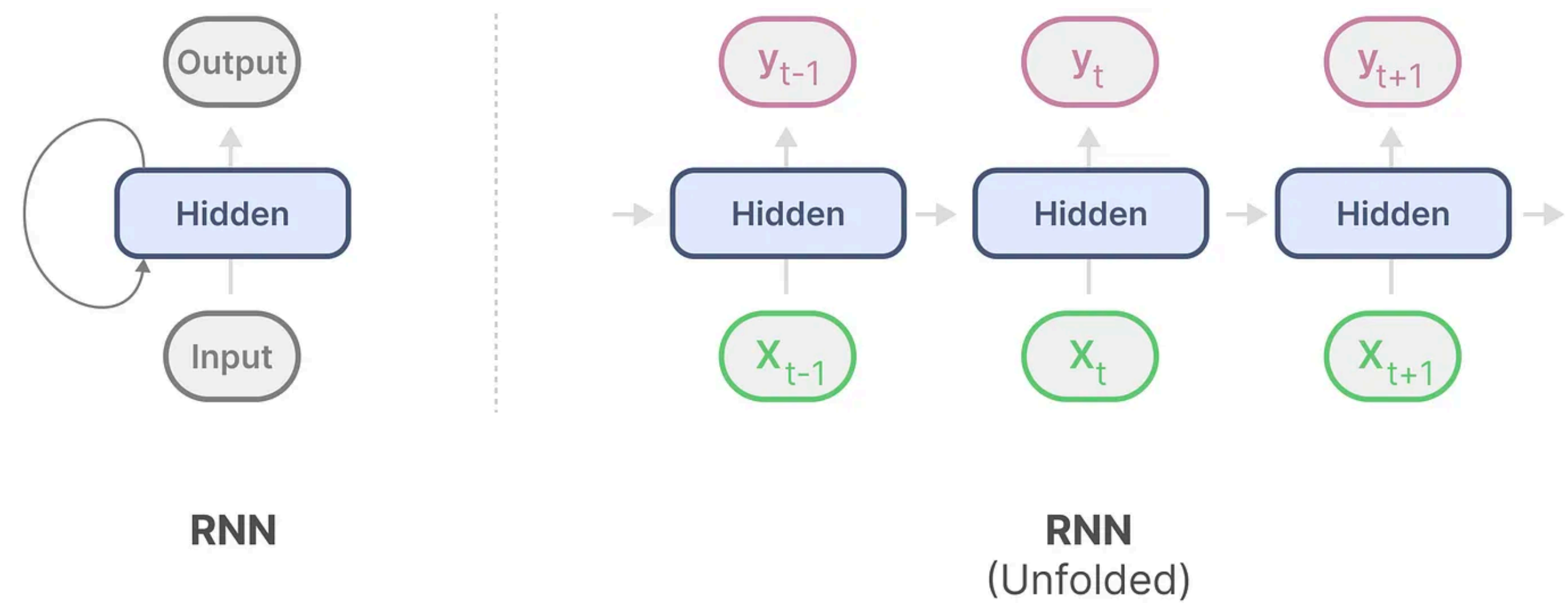


	Recurrent Neural Networks (RNNs)
Parallelizable training	X
Fast autoregressive generation	✓

Inherently sequential forward and backward pass

Constant time and space required to perform single step of generation

Nonlinear Recurrent Neural Networks



	Recurrent Neural Networks (RNNs)
Parallelizable training	✗
Fast autoregressive generation	✓
Avoid vanishing gradients	✗

- ✗ Inherently sequential forward and backward pass
- ✓ Constant time and space required to perform single step of generation
- ✗ Difficult to train to retain information from the past due to this

Attention

Training

$$\text{softmax}\left(\frac{\begin{matrix} \textcolor{violet}{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \textcolor{brown}{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \square & \square \\ \hline \end{array} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \textcolor{blue}{V} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$
$$= \begin{matrix} \textcolor{pink}{Z} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

Attention

Training

$$\text{softmax}\left(\frac{\overset{\text{Q}}{\begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix}} \times \overset{\text{K}^T}{\begin{matrix} \square & \square \\ \square & \square \\ \square & \square \end{matrix}}}{\sqrt{d_k}}\right) \overset{\text{V}}{\begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix}}$$

$$= \overset{\text{Z}}{\begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix}}$$

	Attention
Parallelizable training	✓

Matrix multiplications, modern hardware (GPUs/TPUs) is highly optimized for this (though quadratic complexity)

Attention

Training

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{matrix} \square & \square \\ \square & \square \\ \square & \square \end{matrix} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

$$= \begin{matrix} \text{Z} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

	Attention
Parallelizable training	✓
Fast autoregressive generation	✗

Matrix multiplications, modern hardware (GPUs/TPUs) is highly optimized for this (though quadratic complexity)

Again quadratic complexity, have to compare to all past keys and values each step (growing “state” size, aka KV cache)

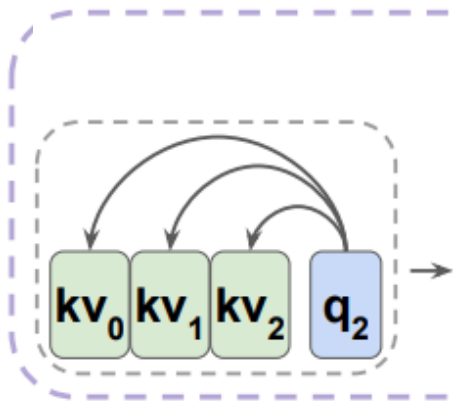
Attention

Training

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{matrix} \square & \square \\ \square & \square \\ \square & \square \end{matrix} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

$$= \begin{matrix} \text{Z} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

Autoregressive Generation



	Attention
Parallelizable training	✓
Fast autoregressive generation	✗

Matrix multiplications, modern hardware (GPUs/TPUs) is highly optimized for this (though quadratic complexity)

Again quadratic complexity, have to compare to all past keys and values each step (growing “state” size, aka KV cache)

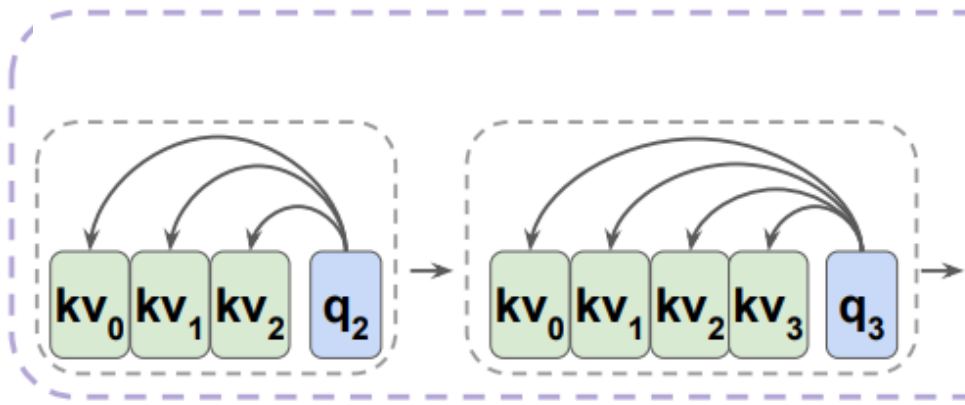
Attention

Training

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{matrix} \square & \square \\ \square & \square \\ \square & \square \end{matrix} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

$$= \begin{matrix} \text{Z} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

Autoregressive Generation



	Attention
Parallelizable training	✓
Fast autoregressive generation	✗

Matrix multiplications, modern hardware (GPUs/TPUs) is highly optimized for this (though quadratic complexity)

Again quadratic complexity, have to compare to all past keys and values each step (growing “state” size, aka KV cache)

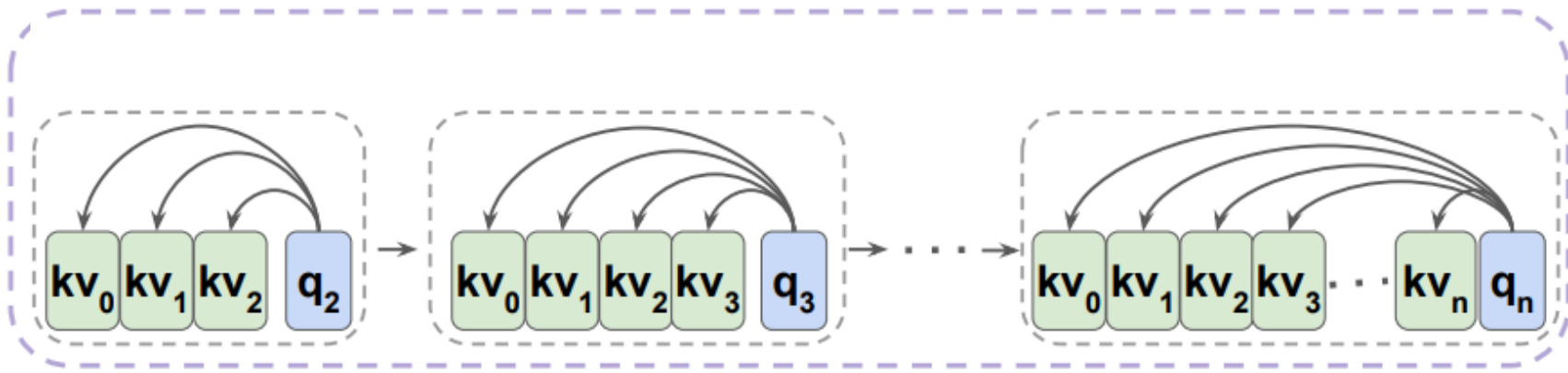
Attention

Training

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \text{2x3 grid} \end{matrix} \times \begin{matrix} \text{K}^T \\ \text{3x2 grid} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \text{3x2 grid} \end{matrix}$$

$$= \begin{matrix} \text{Z} \\ \text{2x3 grid} \end{matrix}$$

Autoregressive Generation



	Attention
Parallelizable training	✓
Fast autoregressive generation	✗

Matrix multiplications, modern hardware (GPUs/TPUs) is highly optimized for this (though quadratic complexity)

Again quadratic complexity, have to compare to all past keys and values each step (growing “state” size, aka KV cache)

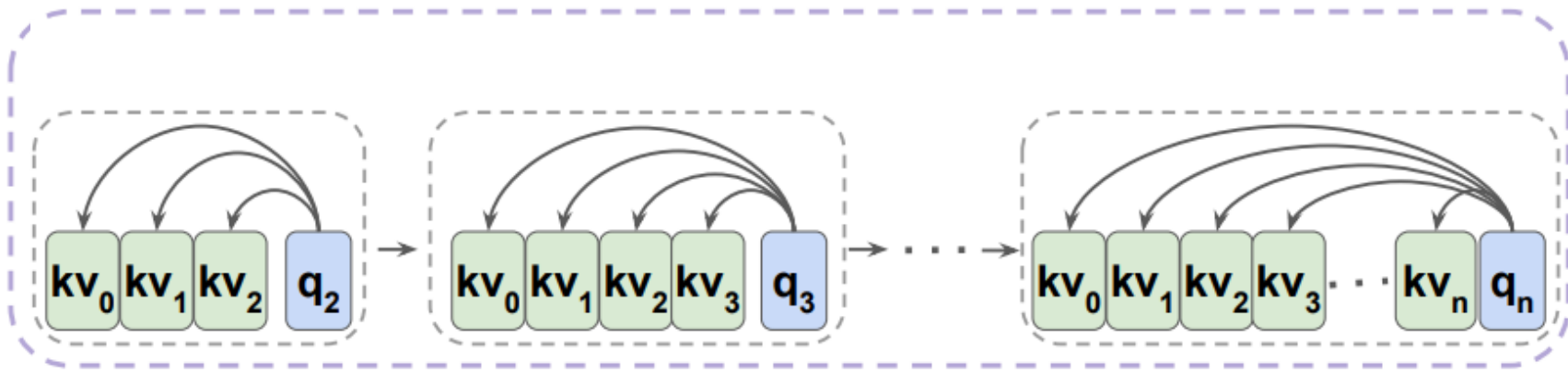
Attention

Training

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{matrix} \square & \square \\ \square & \square \\ \square & \square \end{matrix} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

$$= \begin{matrix} \text{Z} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

Autoregressive Generation



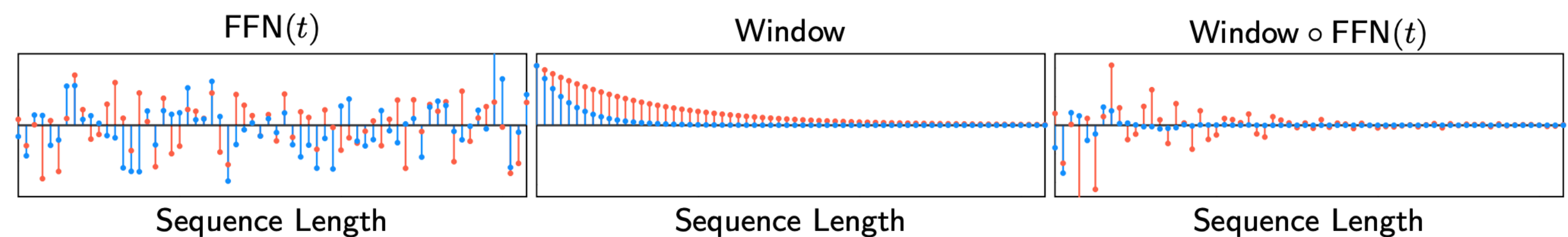
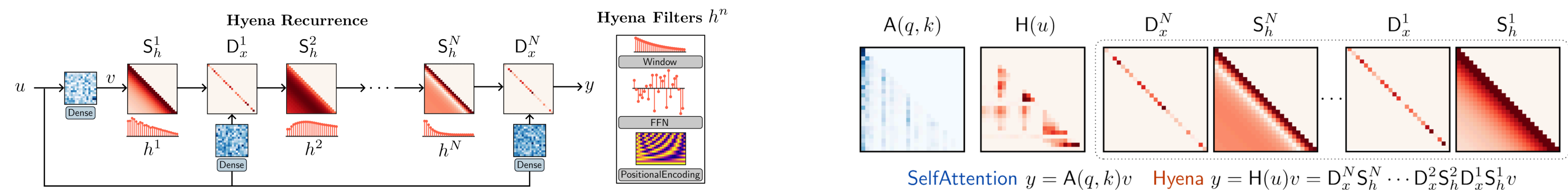
	Attention
Parallelizable training	✓
Fast autoregressive generation	✗
Avoid vanishing gradients	✓

Matrix multiplications, modern hardware (GPUs/TPUs) is highly optimized for this (though quadratic complexity)

Again quadratic complexity, have to compare to all past keys and values each step (growing “state” size, aka KV cache)

O(1) maximum path length between tokens

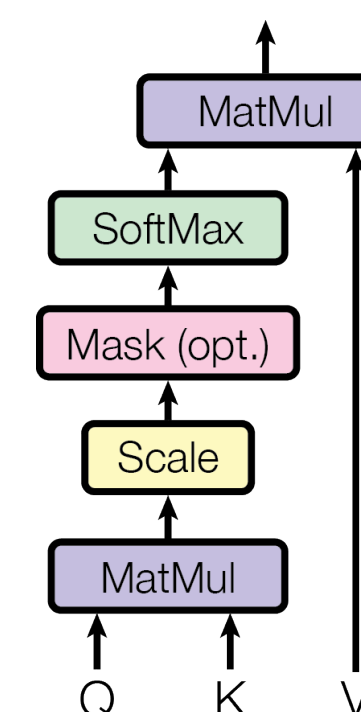
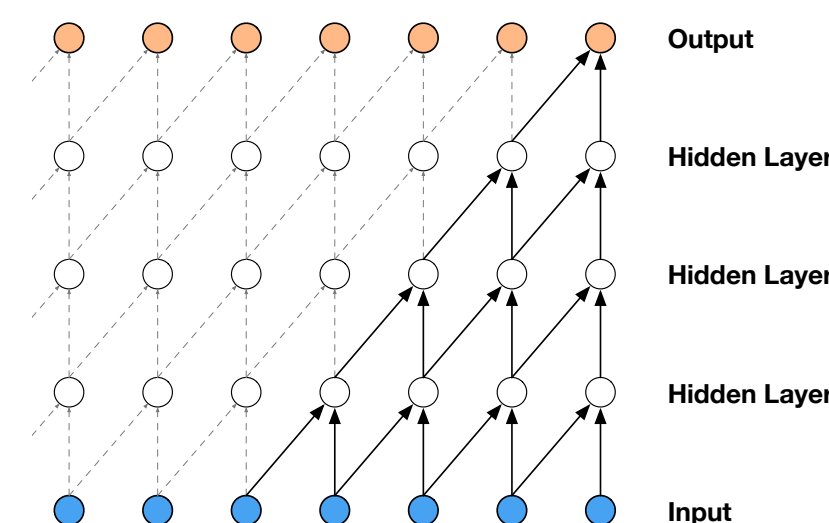
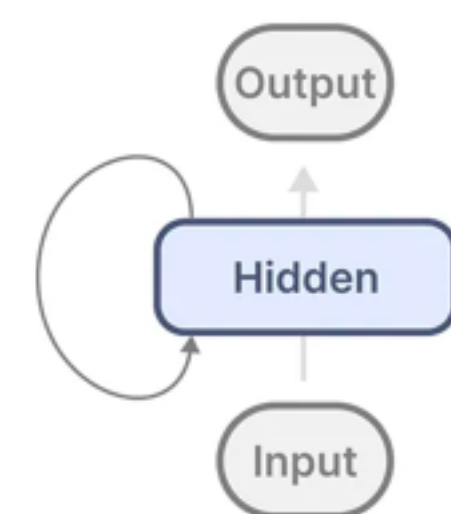
Long Convolutions



	Attention
Parallelizable training	✓
Fast autoregressive generation	✗
Avoid vanishing gradients	✓

- subquadratic complexity (FFTs)
- Quadratic complexity, but can distill into SSM post training (Laughing Hyena, Massaroli 2023)
- No recurrence to have to compute gradients through.

Prior approaches to model long sequences



	Recurrent Neural Networks (RNNs)	Convolutions	Attention
Parallelizable training	X	✓	✓ (Quadratic complexity)
Fast autoregressive generation	✓	X	X
Avoid vanishing gradients	X	✓	✓

Attention Approximations

Linearized Attention

Attention

$$\mathbf{y}_i = \sum_{j=1}^i \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j / \sqrt{d}) \mathbf{v}_j}{\sum_{n=1}^i \exp(\mathbf{q}_i^\top \mathbf{k}_n / \sqrt{d})}$$

Linear Attention

$$\mathbf{y}_i = \sum_{j=1}^i \frac{\phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_j) \mathbf{v}_j}{\sum_{n=1}^i \phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_n)}$$

1. Outer product over key, value head dims
2. Sum over sequence length
3. Dot product over query, key-value head dims

$$\mathbf{y}_i = \frac{\phi(\mathbf{q}_i)^\top \left(\sum_{j=1}^i (\phi(\mathbf{k}_j) \mathbf{v}_j^\top) \right)}{\phi(\mathbf{q}_i)^\top \sum_{n=1}^i \phi(\mathbf{k}_n)}$$

e.g. Linear Transformers (Katharopoulos 2020),
Based (Arora 2024)

Attention Approximations

Linearized Attention

Attention

$$\mathbf{y}_i = \sum_{j=1}^i \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j / \sqrt{d}) \mathbf{v}_j}{\sum_{n=1}^i \exp(\mathbf{q}_i^\top \mathbf{k}_n / \sqrt{d})}$$

Linear Attention

$$\mathbf{y}_i = \sum_{j=1}^i \frac{\phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_j) \mathbf{v}_j}{\sum_{n=1}^i \phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_n)}$$

1. Outer product over key, value head dims
2. Sum over sequence length
3. Dot product over query, key-value head dims

$$\mathbf{y}_i = \frac{\phi(\mathbf{q}_i)^\top \left(\sum_{j=1}^i (\phi(\mathbf{k}_j) \mathbf{v}_j^\top) \right)}{\phi(\mathbf{q}_i)^\top \sum_{n=1}^i \phi(\mathbf{k}_n)}$$

e.g. Linear Transformers (Katharopoulos 2020),
Based (Arora 2024)

Can also be viewed as a linear RNN at inference

Attention Approximations

Linearized Attention

Attention
Linear Attention

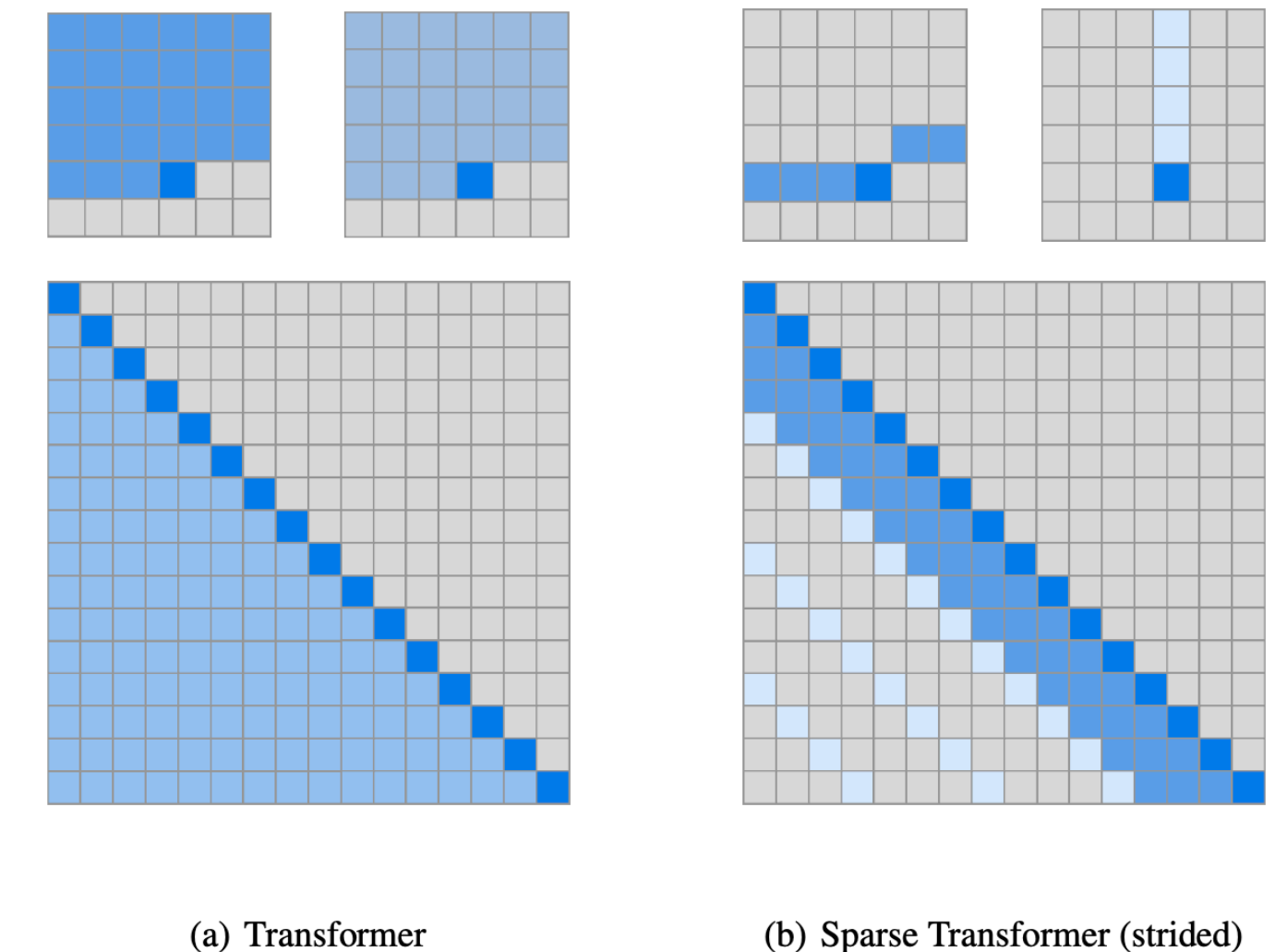
$$\mathbf{y}_i = \sum_{j=1}^i \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j / \sqrt{d}) \mathbf{v}_j}{\sum_{n=1}^i \exp(\mathbf{q}_i^\top \mathbf{k}_n / \sqrt{d})} \quad \longrightarrow \quad \mathbf{y}_i = \sum_{j=1}^i \frac{\phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_j) \mathbf{v}_j}{\sum_{n=1}^i \phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_n)}$$

1. Outer product over key, value head dims
2. Sum over sequence length
3. Dot product over query, key-value head dims

$$\mathbf{y}_i = \frac{\phi(\mathbf{q}_i)^\top \left(\sum_{j=1}^i \left(\phi(\mathbf{k}_j) \mathbf{v}_j^\top \right) \right)}{\phi(\mathbf{q}_i)^\top \sum_{n=1}^i \phi(\mathbf{k}_n)}$$

e.g. Linear Transformers (Katharopoulos 2020),
Based (Arora 2024)

Sparse Attention



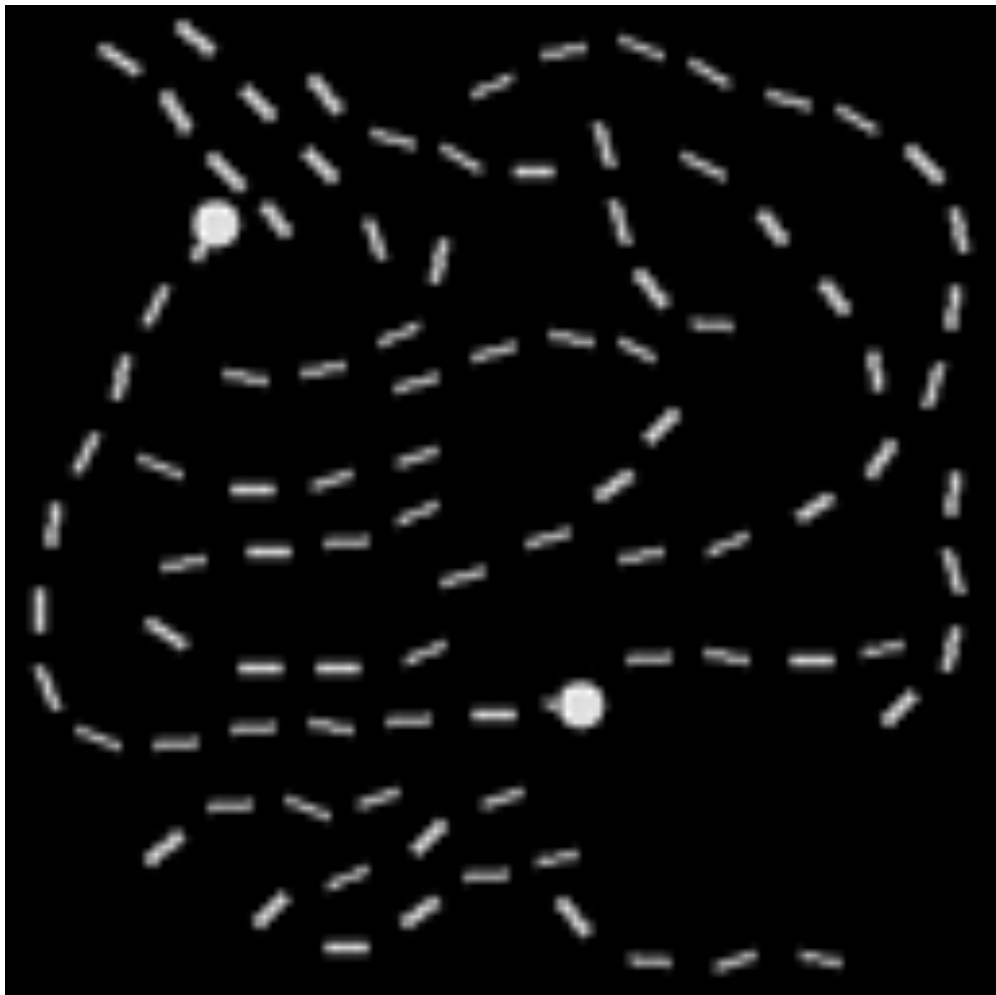
e.g. Sparse Transformers (Child 2019),
Big Bird (Zaheer 2020)

Long range benchmarks

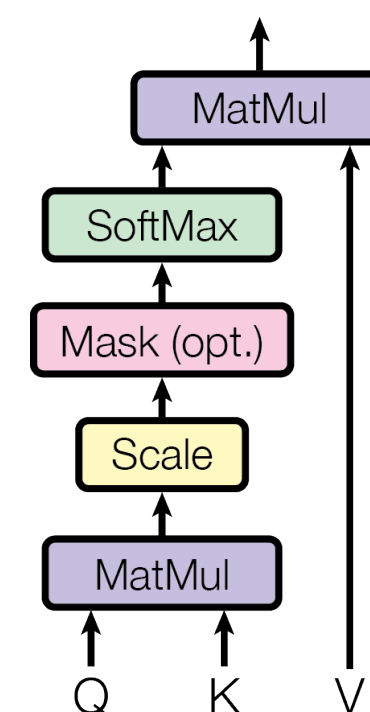
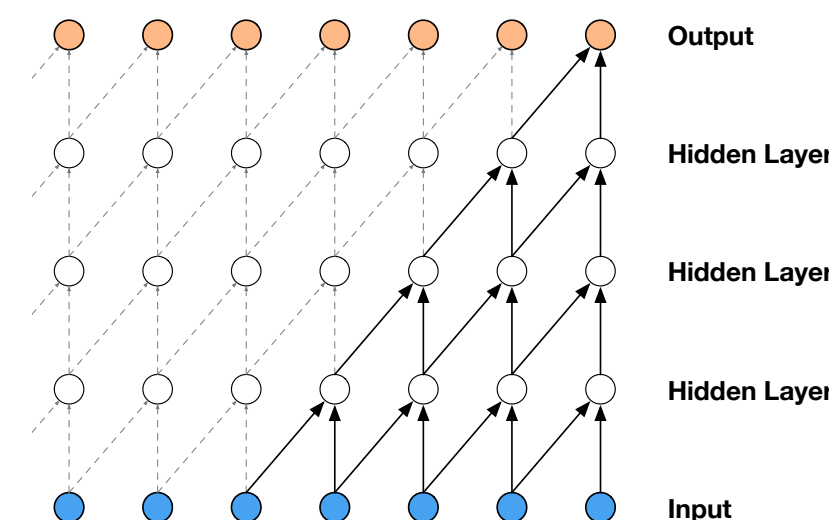
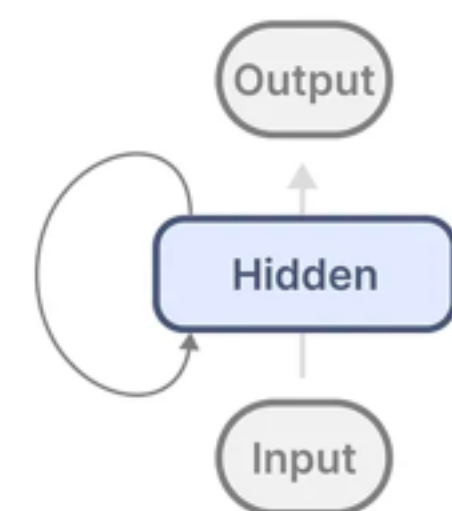
Long Range Arena

Model	LISTOPS	TEXT	RETRIEVAL	IMAGE	PATHFINDER	PATH-X	AVG
Random	10.00	50.00	50.00	10.00	50.00	50.00	36.67
Transformer	36.37	64.27	57.46	42.44	71.40	✗	53.66
Local Attention	15.82	52.98	53.39	41.46	66.63	✗	46.71
Sparse Trans.	17.07	63.58	59.59	44.24	71.71	✗	51.03
Longformer	35.63	62.85	56.89	42.22	69.71	✗	52.88
Linformer	35.70	53.94	52.27	38.56	76.34	✗	51.14
Reformer	<u>37.27</u>	56.10	53.40	38.07	68.50	✗	50.56
Sinkhorn Trans.	33.67	61.20	53.83	41.23	67.45	✗	51.23
Synthesizer	36.99	61.68	54.67	41.61	69.45	✗	52.40
BigBird	36.05	64.02	59.29	40.83	74.87	✗	54.17
Linear Trans.	16.13	<u>65.90</u>	53.09	42.34	75.30	✗	50.46
Performer	18.01	65.40	53.82	42.77	77.05	✗	51.18
FNet	35.33	65.11	59.61	38.67	<u>77.80</u>	✗	54.42
Nyströmformer	37.15	65.52	<u>79.56</u>	41.58	70.94	✗	57.46
Luna-256	37.25	64.57	79.29	<u>47.38</u>	77.72	✗	<u>59.37</u>

Path-X example:

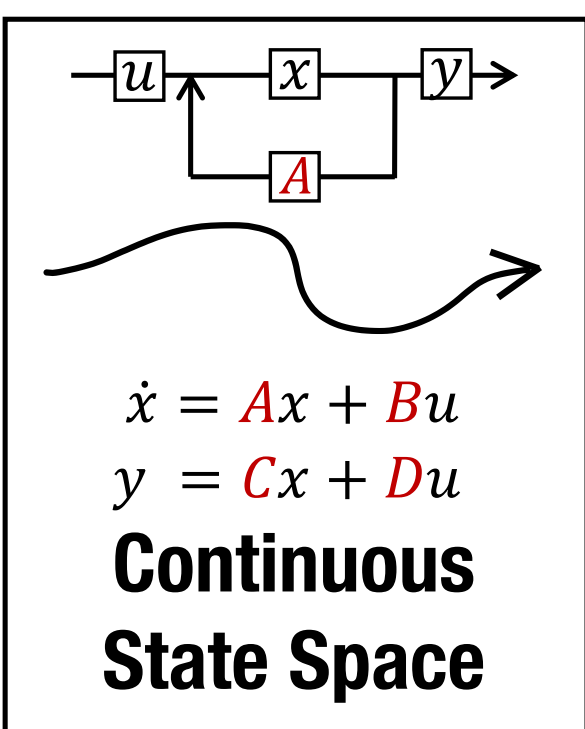
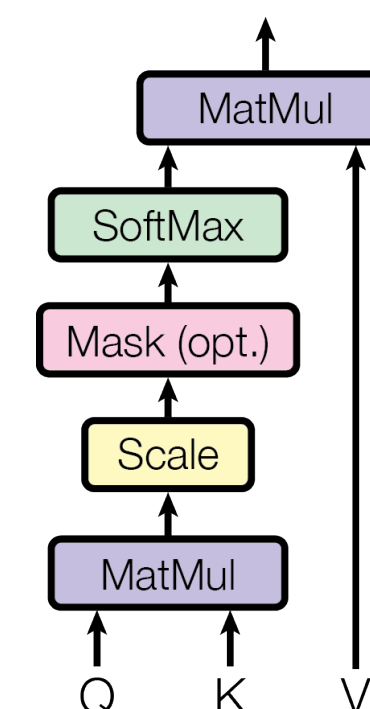
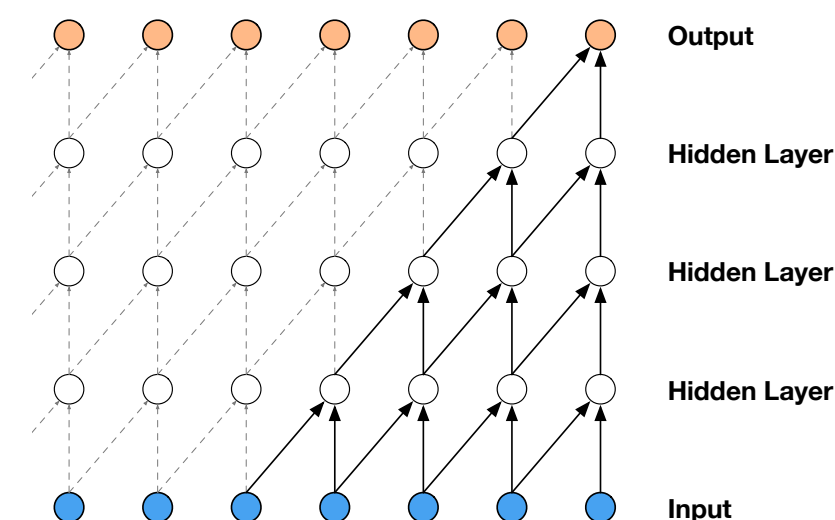
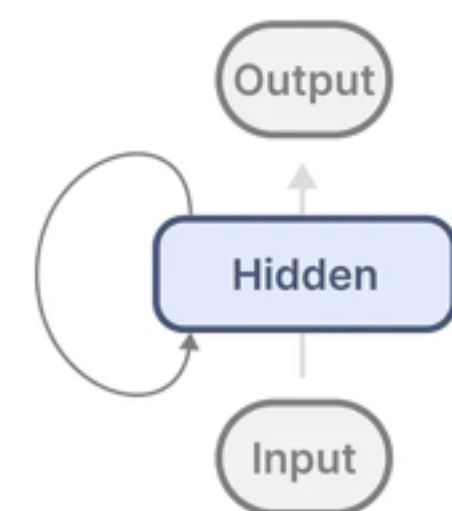


Deep SSMs



	Recurrent Neural Networks (RNNs)	Convolutions	Attention
Parallelizable training	X	✓	✓ (Quadratic complexity)
Fast autoregressive generation	✓	X	X
Avoid vanishing gradients	X	✓	✓

Deep SSMs



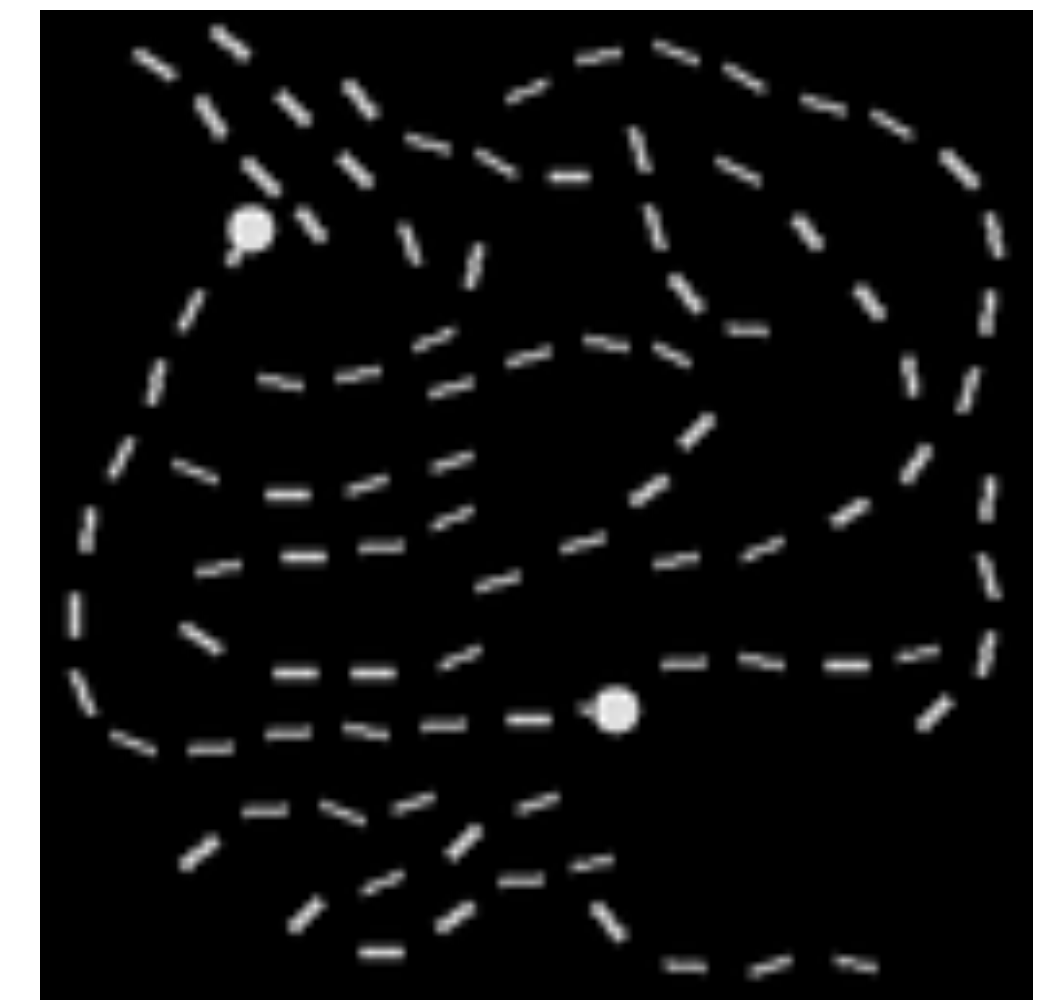
	Recurrent Neural Networks (RNNs)	Convolutions	Attention	Deep SSMs e.g. S4 (Gu et al. ICLR 2022)
Parallelizable training	X	✓	✓ (Quadratic complexity)	✓ (Subquadratic complexity)
Fast autoregressive generation	✓	X	X	✓
Avoid vanishing gradients	X	✓	✓	✓

S4 captures long-range dependencies

Long Range Arena

Model	LISTOPS	TEXT	RETRIEVAL	IMAGE	PATHFINDER	PATH-X	AVG
Random	10.00	50.00	50.00	10.00	50.00	50.00	36.67
Transformer	36.37	64.27	57.46	42.44	71.40	✗	53.66
Local Attention	15.82	52.98	53.39	41.46	66.63	✗	46.71
Sparse Trans.	17.07	63.58	59.59	44.24	71.71	✗	51.03
Longformer	35.63	62.85	56.89	42.22	69.71	✗	52.88
Linformer	35.70	53.94	52.27	38.56	76.34	✗	51.14
Reformer	<u>37.27</u>	56.10	53.40	38.07	68.50	✗	50.56
Sinkhorn Trans.	33.67	61.20	53.83	41.23	67.45	✗	51.23
Synthesizer	36.99	61.68	54.67	41.61	69.45	✗	52.40
BigBird	36.05	64.02	59.29	40.83	74.87	✗	54.17
Linear Trans.	16.13	<u>65.90</u>	53.09	42.34	75.30	✗	50.46
Performer	18.01	65.40	53.82	42.77	77.05	✗	51.18
FNet	35.33	65.11	59.61	38.67	<u>77.80</u>	✗	54.42
Nyströmformer	37.15	65.52	<u>79.56</u>	41.58	70.94	✗	57.46
Luna-256	37.25	64.57	79.29	<u>47.38</u>	77.72	✗	<u>59.37</u>
S4 (original)	58.35	76.02	87.09	87.26	86.05	88.10	80.48
S4 (updated)	59.60	86.82	90.90	88.65	94.20	96.35	86.09

Path-X example:

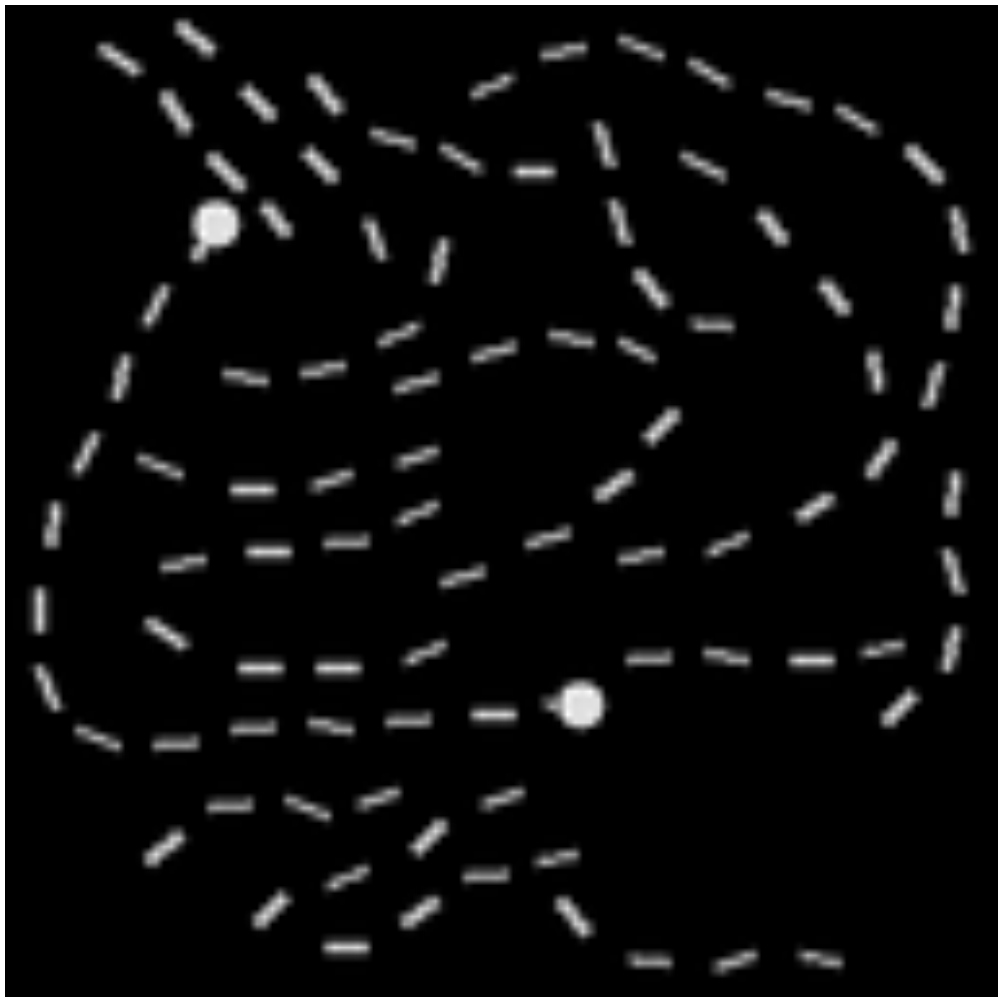


S4 captures long-range dependencies

Long Range Arena

Model	LISTOPS	TEXT	RETRIEVAL	IMAGE	PATHFINDER	PATH-X	AVG
Random	10.00	50.00	50.00	10.00	50.00	50.00	36.67
Transformer	36.37	64.27	57.46	42.44	71.40	✗	53.66
Local Attention	15.82	52.98	53.39	41.46	66.63	✗	46.71
Sparse Trans.	17.07	63.58	59.59	44.24	71.71	✗	51.03
Longformer	35.63	62.85	56.89	42.22	69.71	✗	52.88
Linformer	35.70	53.94	52.27	38.56	76.34	✗	51.14
Reformer	<u>37.27</u>	56.10	53.40	38.07	68.50	✗	50.56
Sinkhorn Trans.	33.67	61.20	53.83	41.23	67.45	✗	51.23
Synthesizer	36.99	61.68	54.67	41.61	69.45	✗	52.40
BigBird	36.05	64.02	59.29	40.83	74.87	✗	54.17
Linear Trans.	16.13	<u>65.90</u>	53.09	42.34	75.30	✗	50.46
Performer	18.01	65.40	53.82	42.77	77.05	✗	51.18
FNet	35.33	65.11	59.61	38.67	<u>77.80</u>	✗	54.42
Nyströmformer	37.15	65.52	<u>79.56</u>	41.58	70.94	✗	57.46
Luna-256	37.25	64.57	79.29	<u>47.38</u>	77.72	✗	<u>59.37</u>
S4 (original)	58.35	76.02	87.09	87.26	86.05	88.10	80.48
S4 (updated)	59.60	86.82	90.90	88.65	94.20	96.35	86.09

Path-X example:

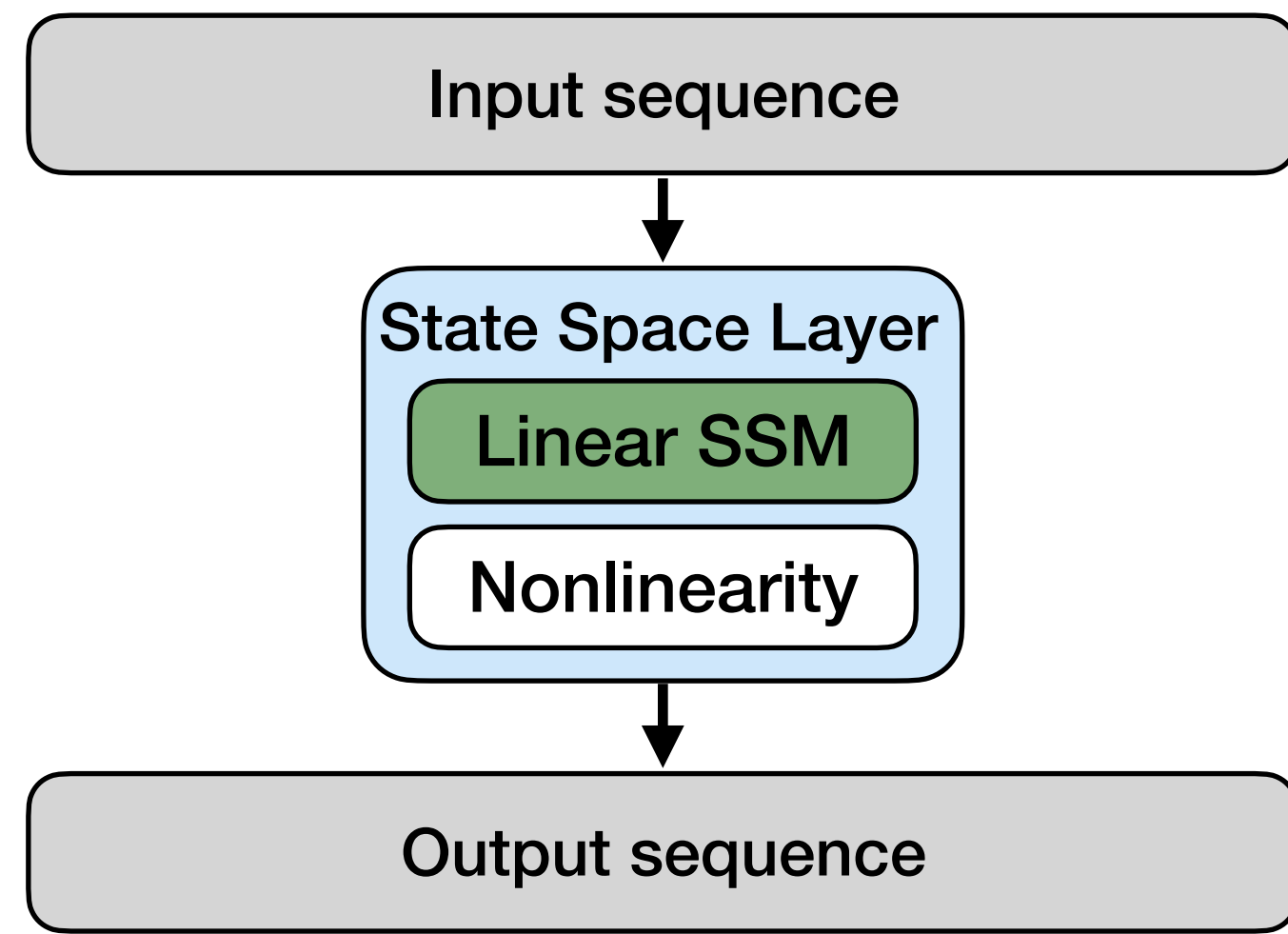


Agenda

- Introduction, motivation, prior approaches
- **Linear state space models (SSMs) overview**
- S4, convolutions, parameterization
- S5, diagonalization, parallel scans
- S6/Mamba/Hawk/Griffin, data-dependent dynamics
- Challenges and opportunities

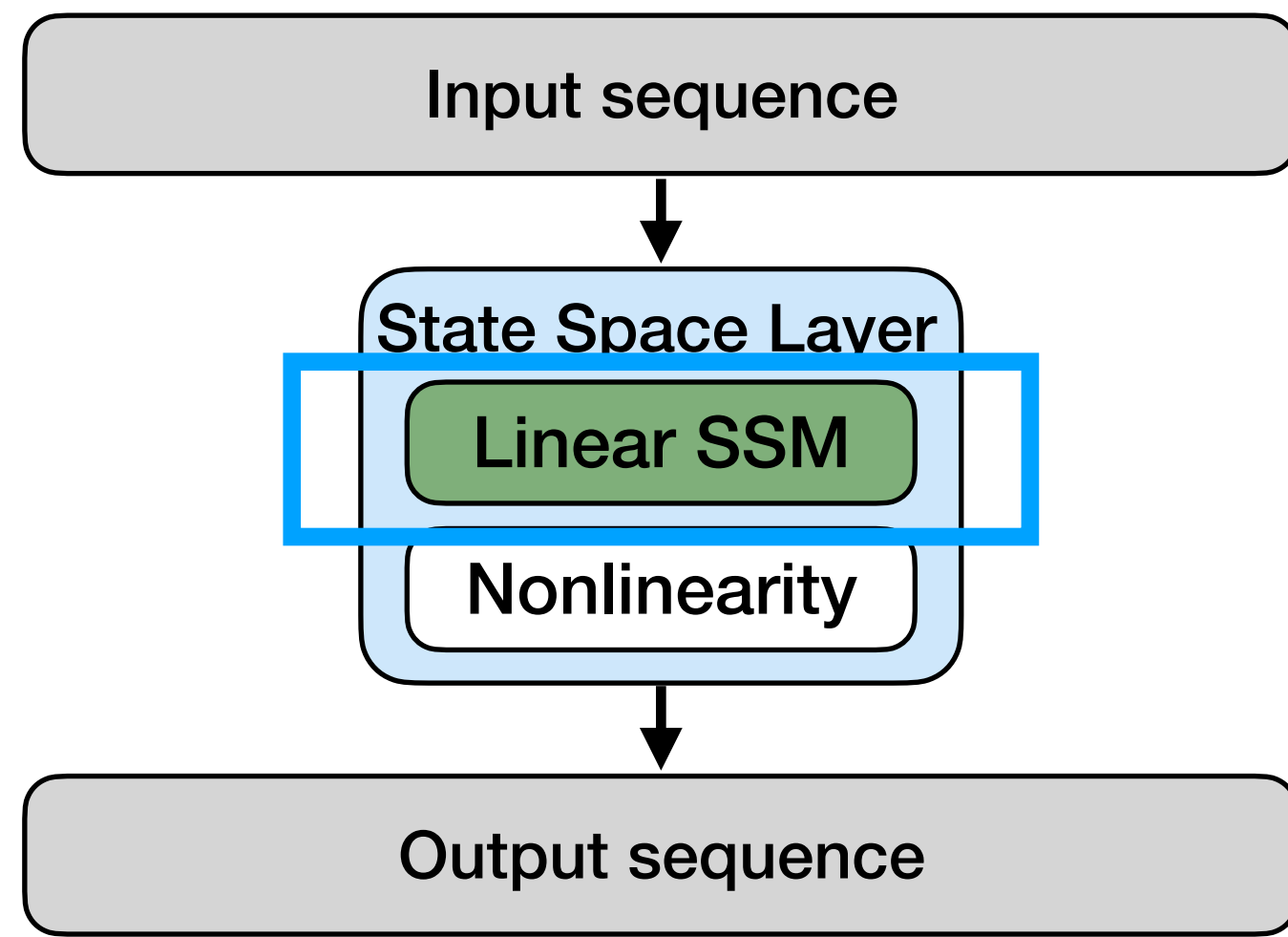
Key idea of Deep SSMs: Linear in time, nonlinear in depth

Key idea of Deep SSMs: Linear in time, nonlinear in depth



Key idea of Deep SSMs: Linear in time, nonlinear in depth

Continuous-time, linear state space model (SSM):



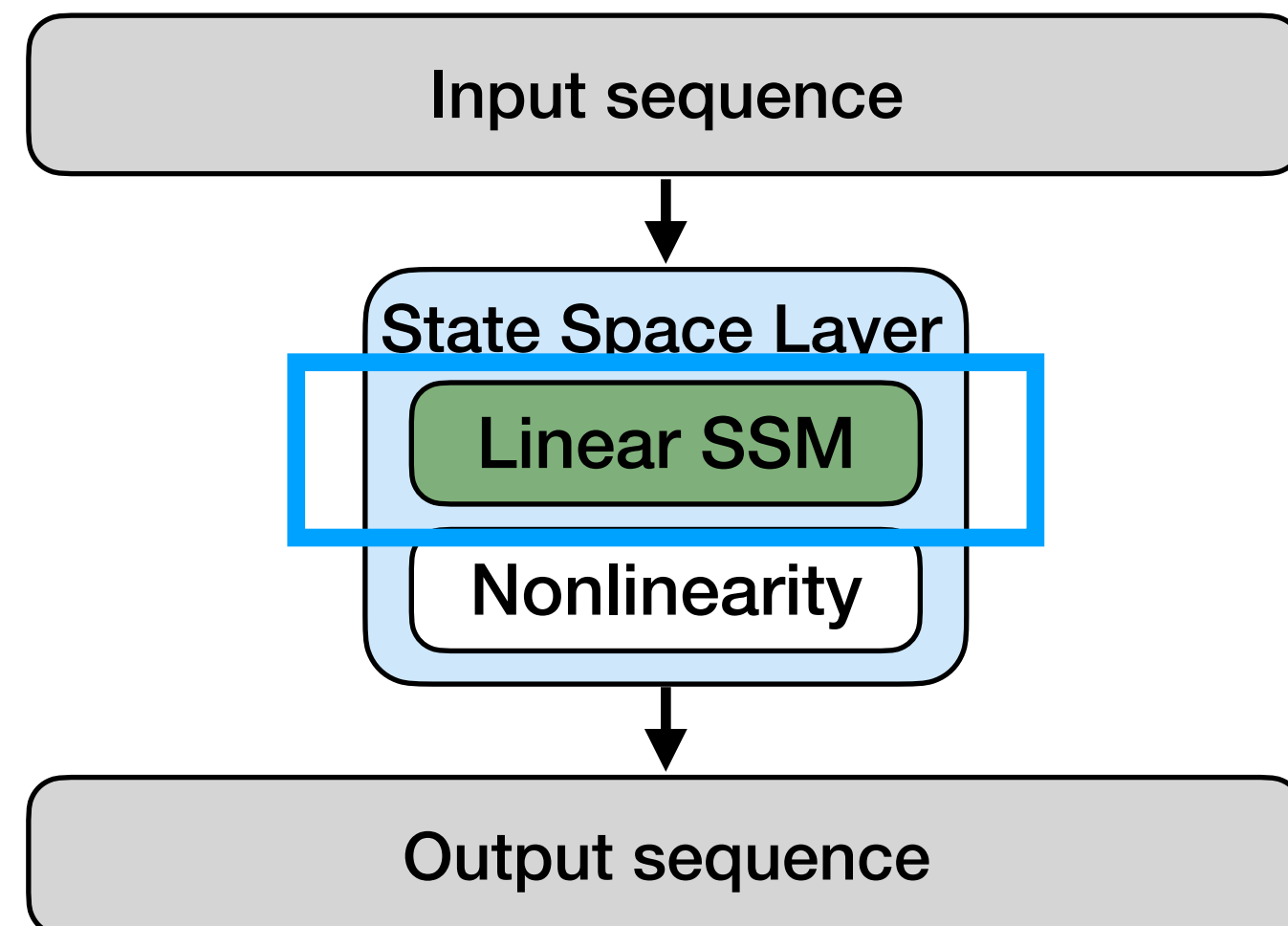
Key idea of Deep SSMs: Linear in time, nonlinear in depth

Continuous-time, linear state space model (SSM):

Input Signal: $\mathbf{u}(t) \in \mathbb{R}^U$

Hidden State: $\mathbf{x}(t) \in \mathbb{R}^N$

Output Signal: $\mathbf{y}(t) \in \mathbb{R}^M$



Key idea of Deep SSMs: Linear in time, nonlinear in depth

Continuous-time, linear state space model (SSM):

Input Signal: $\mathbf{u}(t) \in \mathbb{R}^U$

Hidden State: $\mathbf{x}(t) \in \mathbb{R}^N$

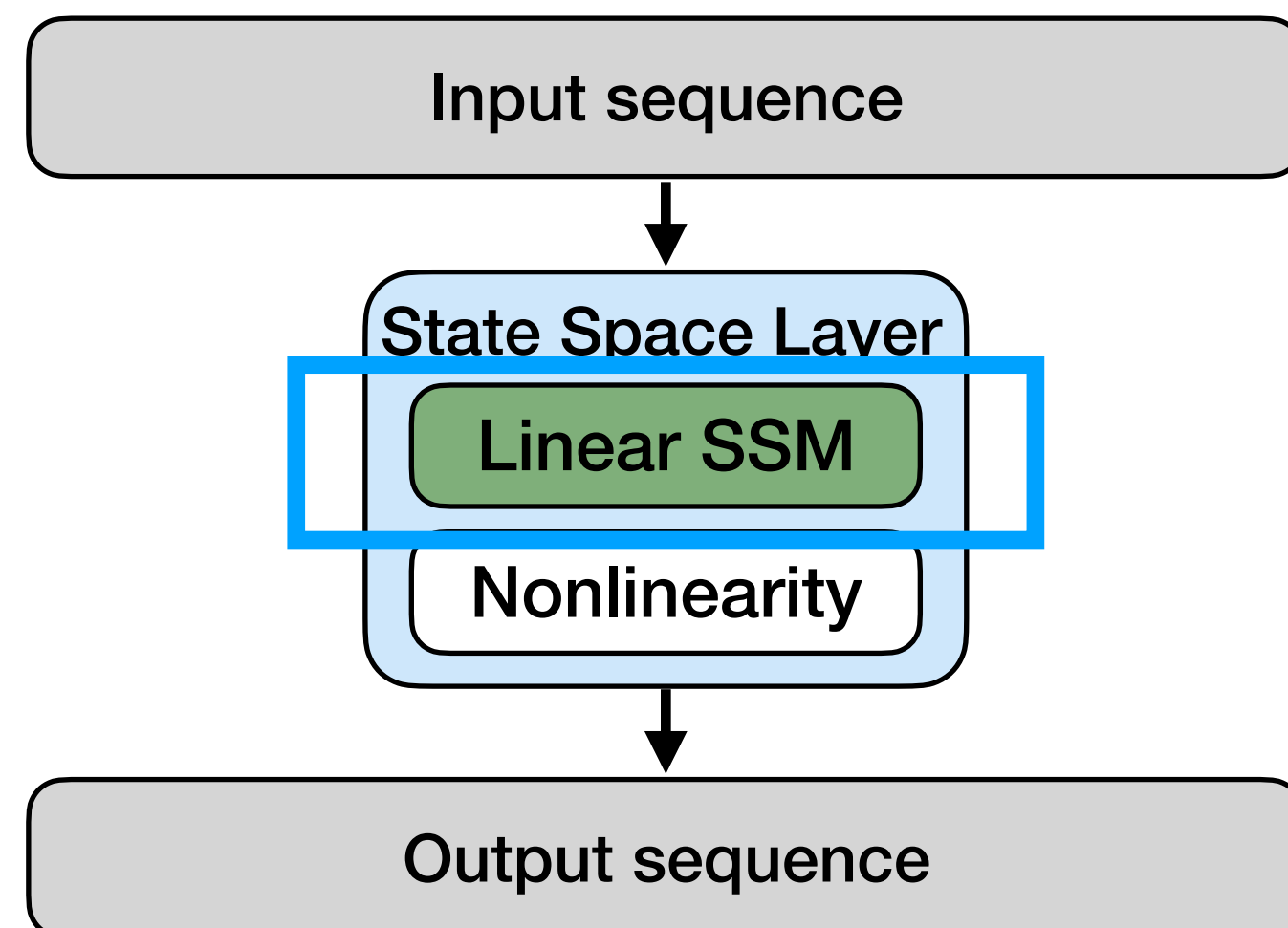
Output Signal: $\mathbf{y}(t) \in \mathbb{R}^M$

State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times U}$

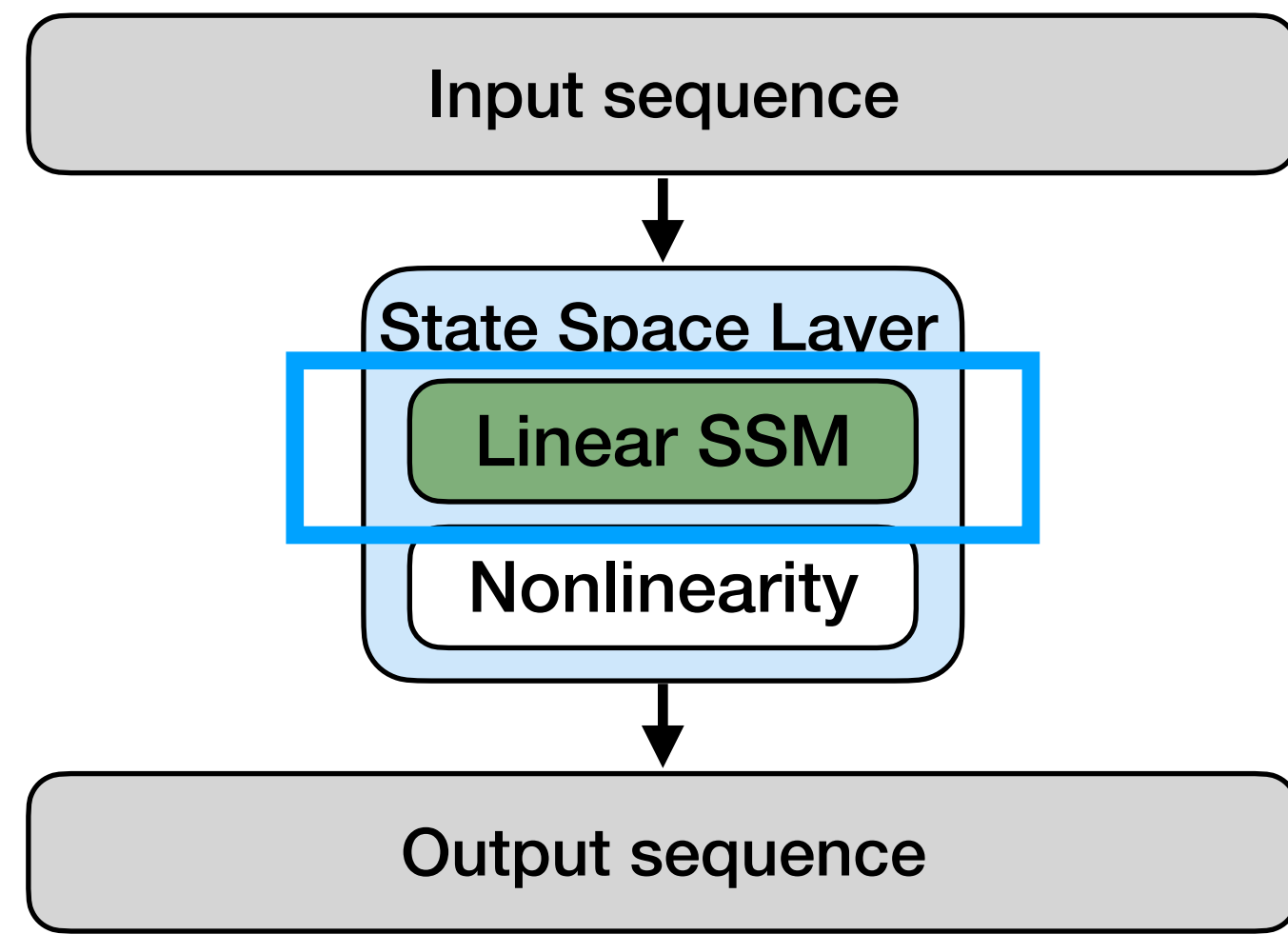
Output Matrix: $\mathbf{C} \in \mathbb{R}^{M \times N}$

Feedthrough Matrix: $\mathbf{D} \in \mathbb{R}^{M \times U}$



Key idea of Deep SSMs: Linear in time, nonlinear in depth

Continuous-time, linear state space model (SSM):



Input Signal: $\mathbf{u}(t) \in \mathbb{R}^U$

Hidden State: $\mathbf{x}(t) \in \mathbb{R}^N$

Output Signal: $\mathbf{y}(t) \in \mathbb{R}^M$

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$$

$$\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t)$$

State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

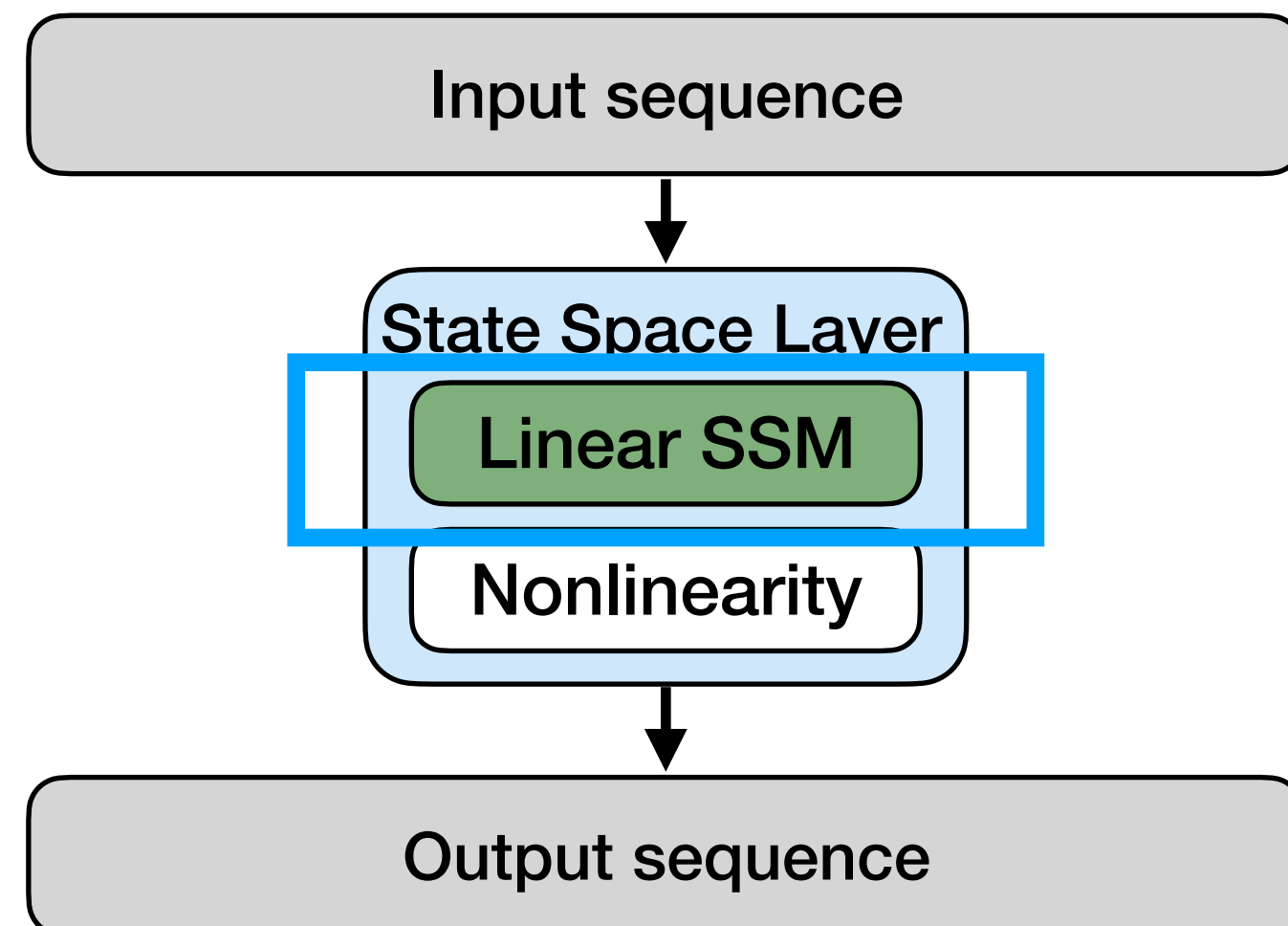
Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times U}$

Output Matrix: $\mathbf{C} \in \mathbb{R}^{M \times N}$

Feedthrough Matrix: $\mathbf{D} \in \mathbb{R}^{M \times U}$

Key idea of Deep SSMs: Linear in time, nonlinear in depth

Continuous-time, linear state space model (SSM):



Input Signal: $\mathbf{u}(t) \in \mathbb{R}^U$

Hidden State: $\mathbf{x}(t) \in \mathbb{R}^N$

Output Signal: $\mathbf{y}(t) \in \mathbb{R}^M$

State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times U}$

Output Matrix: $\mathbf{C} \in \mathbb{R}^{M \times N}$

Feedthrough Matrix: $\mathbf{D} \in \mathbb{R}^{M \times U}$

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$$

$$\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t)$$

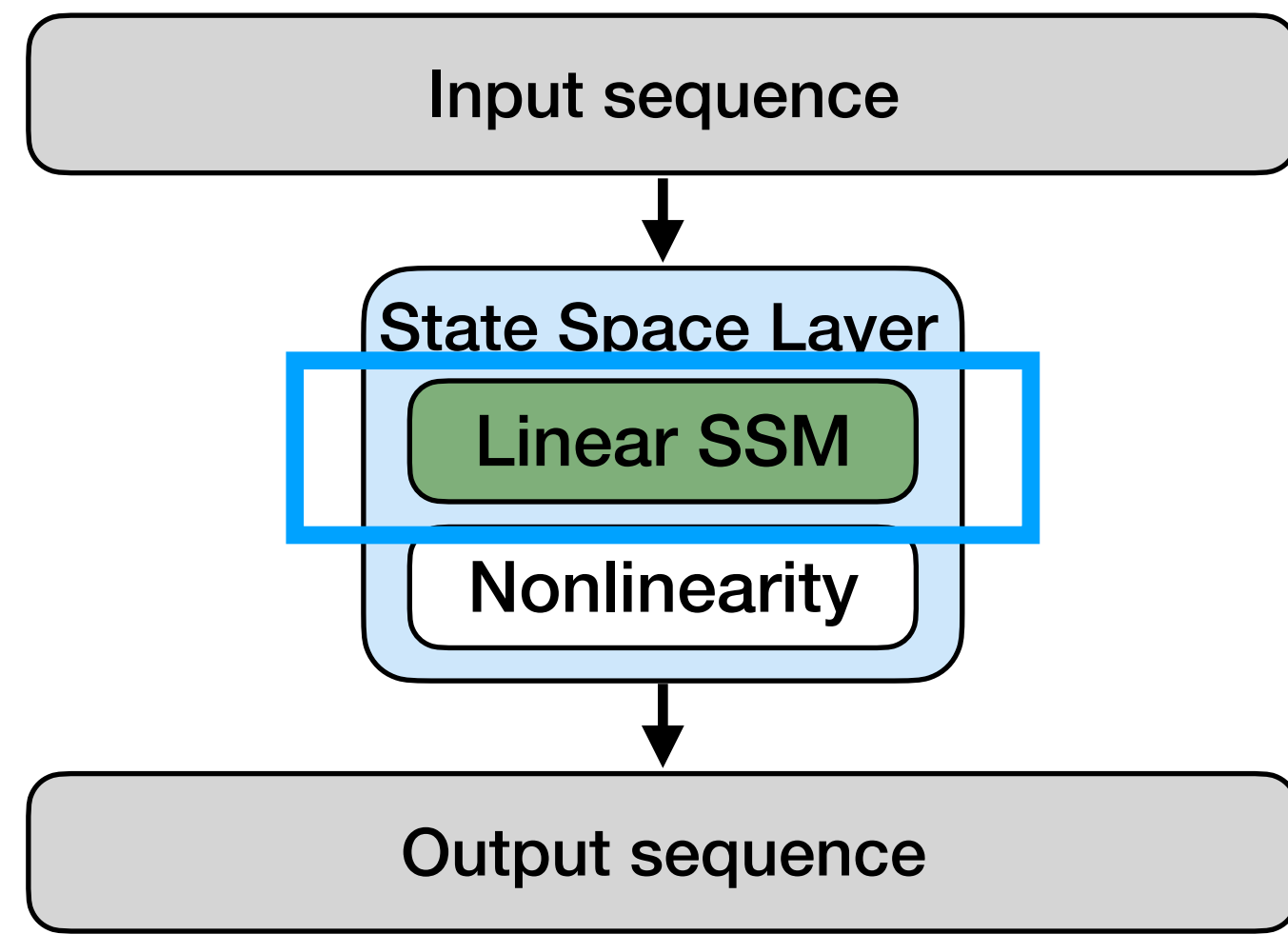
Discretized linear SSM:

$$\mathbf{x}_k = \bar{\mathbf{A}}\mathbf{x}_{k-1} + \bar{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \bar{\mathbf{C}}\mathbf{x}_k + \bar{\mathbf{D}}\mathbf{u}_k$$

Key idea of Deep SSMs: Linear in time, nonlinear in depth

Continuous-time, linear state space model (SSM):



Input Signal: $\mathbf{u}(t) \in \mathbb{R}^U$

Hidden State: $\mathbf{x}(t) \in \mathbb{R}^N$

Output Signal: $\mathbf{y}(t) \in \mathbb{R}^M$

State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times U}$

Output Matrix: $\mathbf{C} \in \mathbb{R}^{M \times N}$

Feedthrough Matrix: $\mathbf{D} \in \mathbb{R}^{M \times U}$

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$$

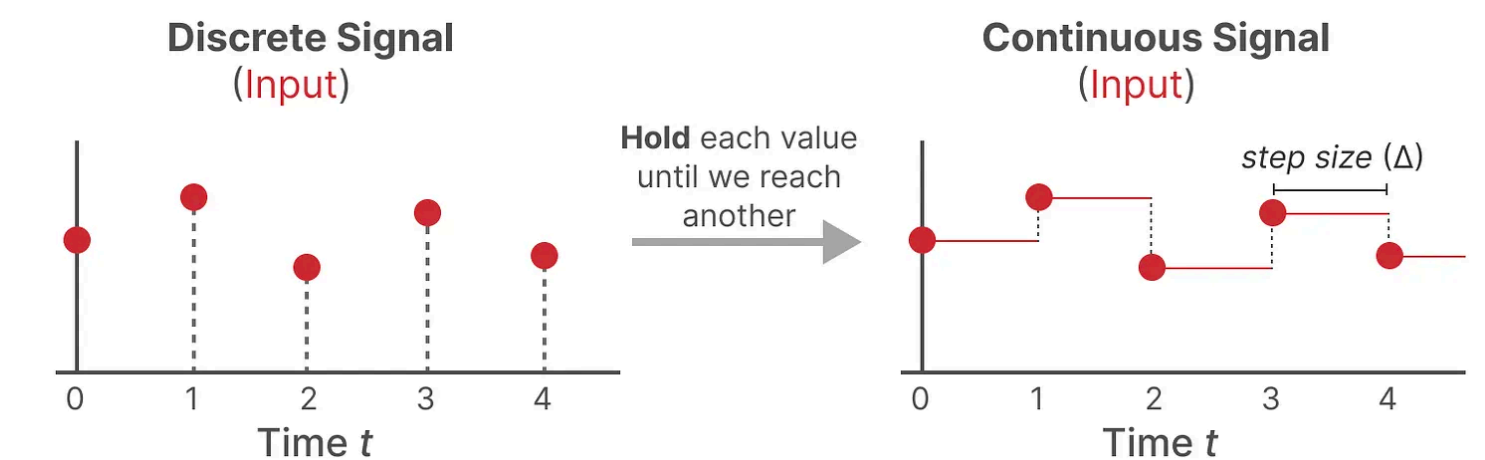
$$\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t)$$

Discretized linear SSM:

$$\mathbf{x}_k = \bar{\mathbf{A}}\mathbf{x}_{k-1} + \bar{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \bar{\mathbf{C}}\mathbf{x}_k + \bar{\mathbf{D}}\mathbf{u}_k$$

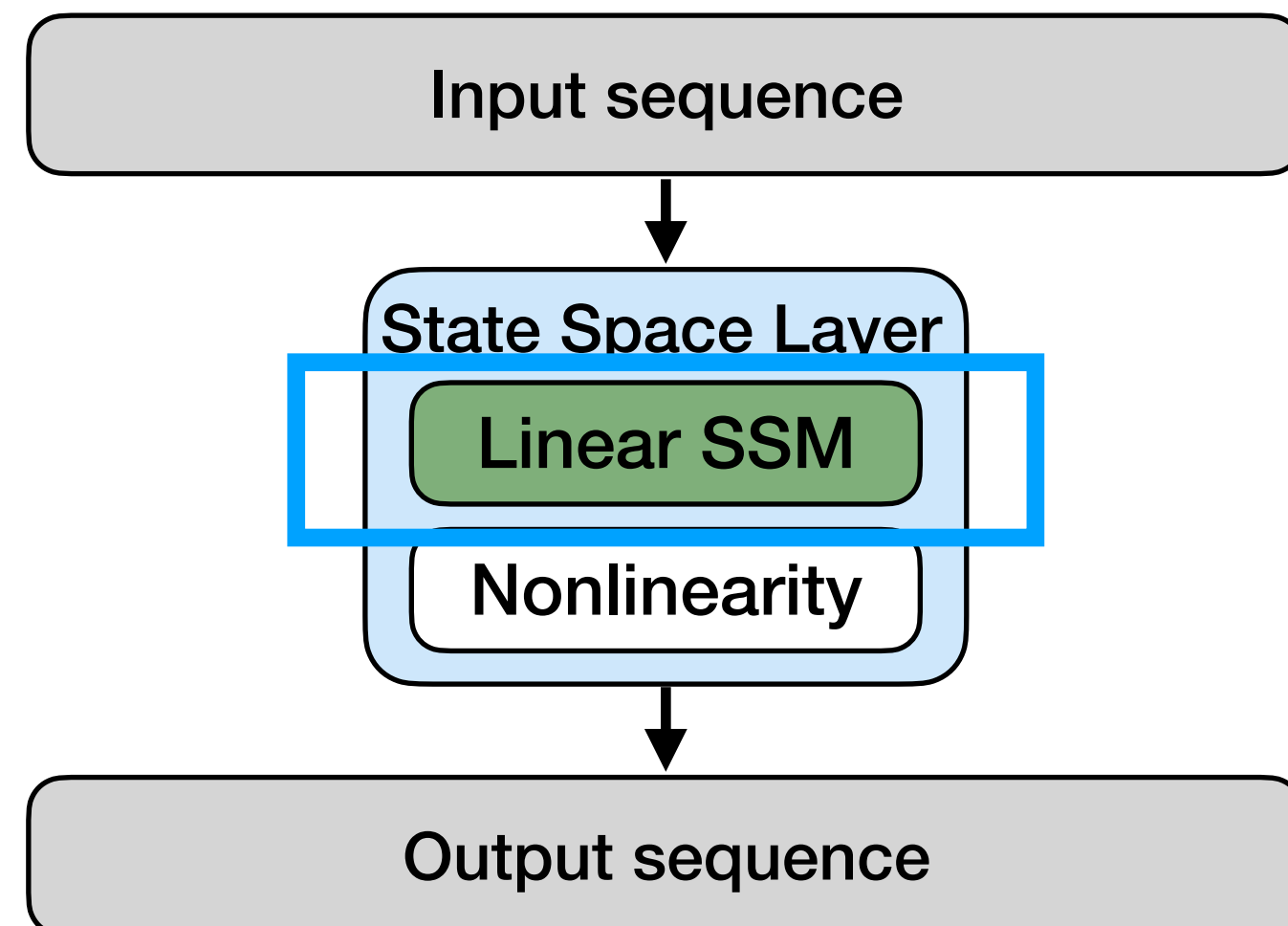
E.g. using Zero-order hold (ZOH):



$$\bar{\mathbf{A}} = e^{\mathbf{A}\Delta}, \quad \bar{\mathbf{B}} = \mathbf{A}^{-1}(\bar{\mathbf{A}} - \mathbf{I})\mathbf{B}, \quad \bar{\mathbf{C}} = \mathbf{C}, \quad \bar{\mathbf{D}} = \mathbf{D}$$

Key idea of Deep SSMs: Linear in time, nonlinear in depth

Continuous-time, linear state space model (SSM):



Input Signal: $\mathbf{u}(t) \in \mathbb{R}^U$

Hidden State: $\mathbf{x}(t) \in \mathbb{R}^N$

Output Signal: $\mathbf{y}(t) \in \mathbb{R}^M$

State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times U}$

Output Matrix: $\mathbf{C} \in \mathbb{R}^{M \times N}$

Feedthrough Matrix: $\mathbf{D} \in \mathbb{R}^{M \times U}$

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$$

$$\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t)$$

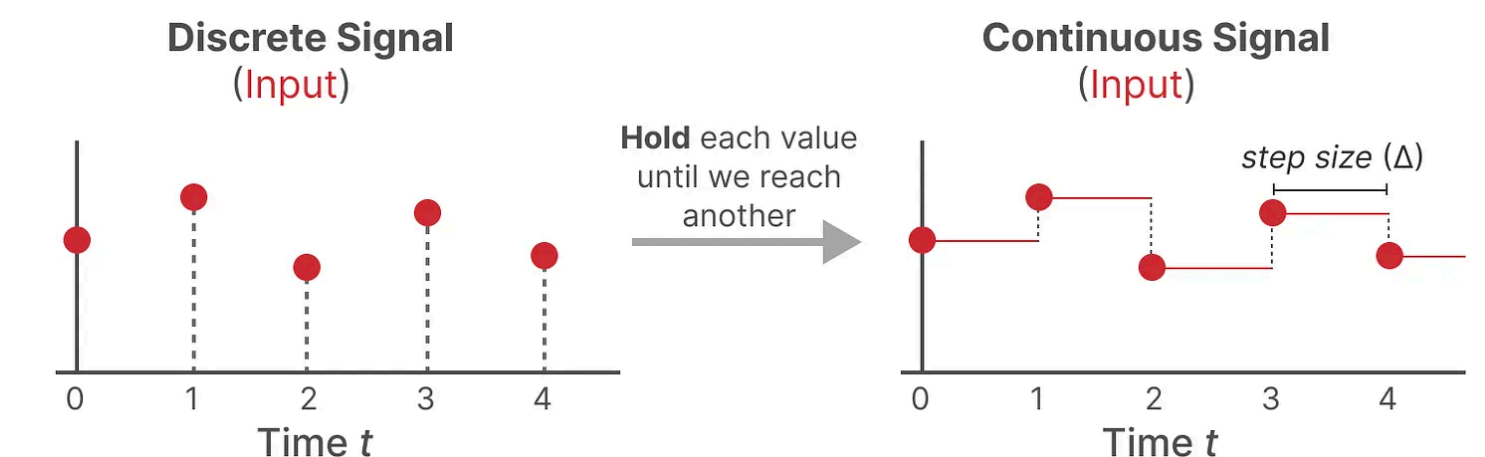
Discretized linear SSM:

$$\mathbf{x}_k = \bar{\mathbf{A}}\mathbf{x}_{k-1} + \bar{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \bar{\mathbf{C}}\mathbf{x}_k + \bar{\mathbf{D}}\mathbf{u}_k$$

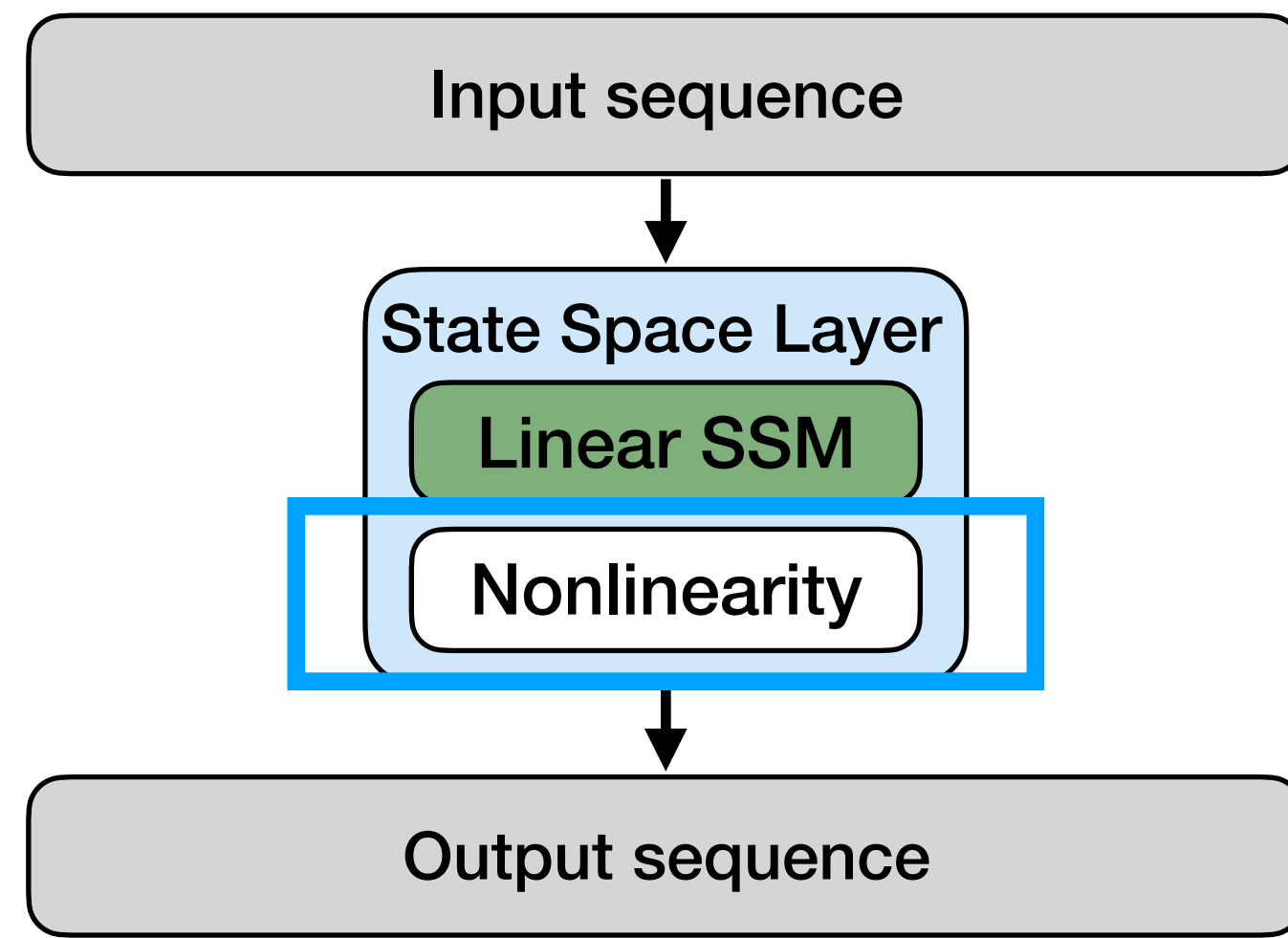
E.g. using Zero-order hold (ZOH):

Treat Δ as learnable parameter.



$$\bar{\mathbf{A}} = e^{\mathbf{A}\Delta}, \quad \bar{\mathbf{B}} = \mathbf{A}^{-1}(\bar{\mathbf{A}} - \mathbf{I})\mathbf{B}, \quad \bar{\mathbf{C}} = \mathbf{C}, \quad \bar{\mathbf{D}} = \mathbf{D}$$

Key idea of Deep SSMs: Linear in time, nonlinear in depth

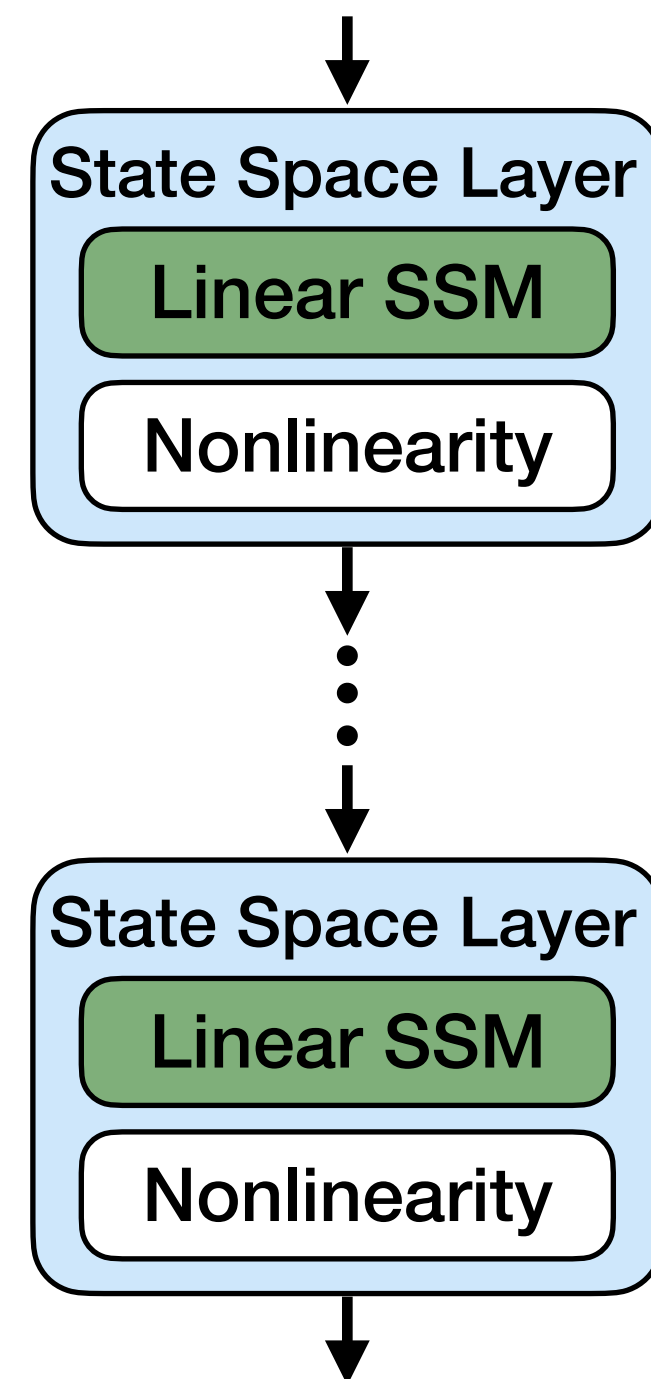


Nonlinear activation:

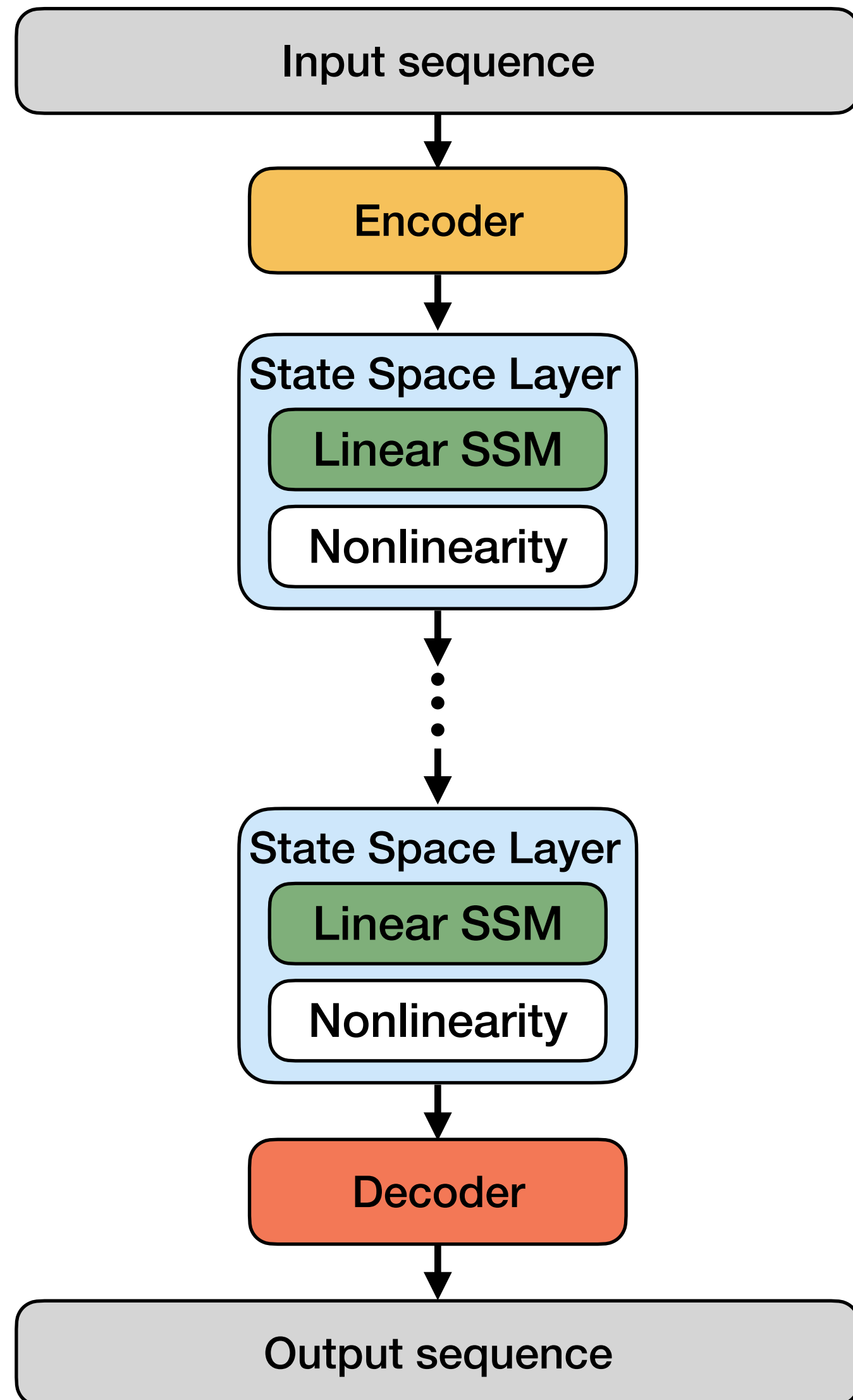
$$\mathbf{u}'_k = \mathbf{f}(\mathbf{y}_k)$$

E.g: gelu, GLU, layer norm, dropout, etc.

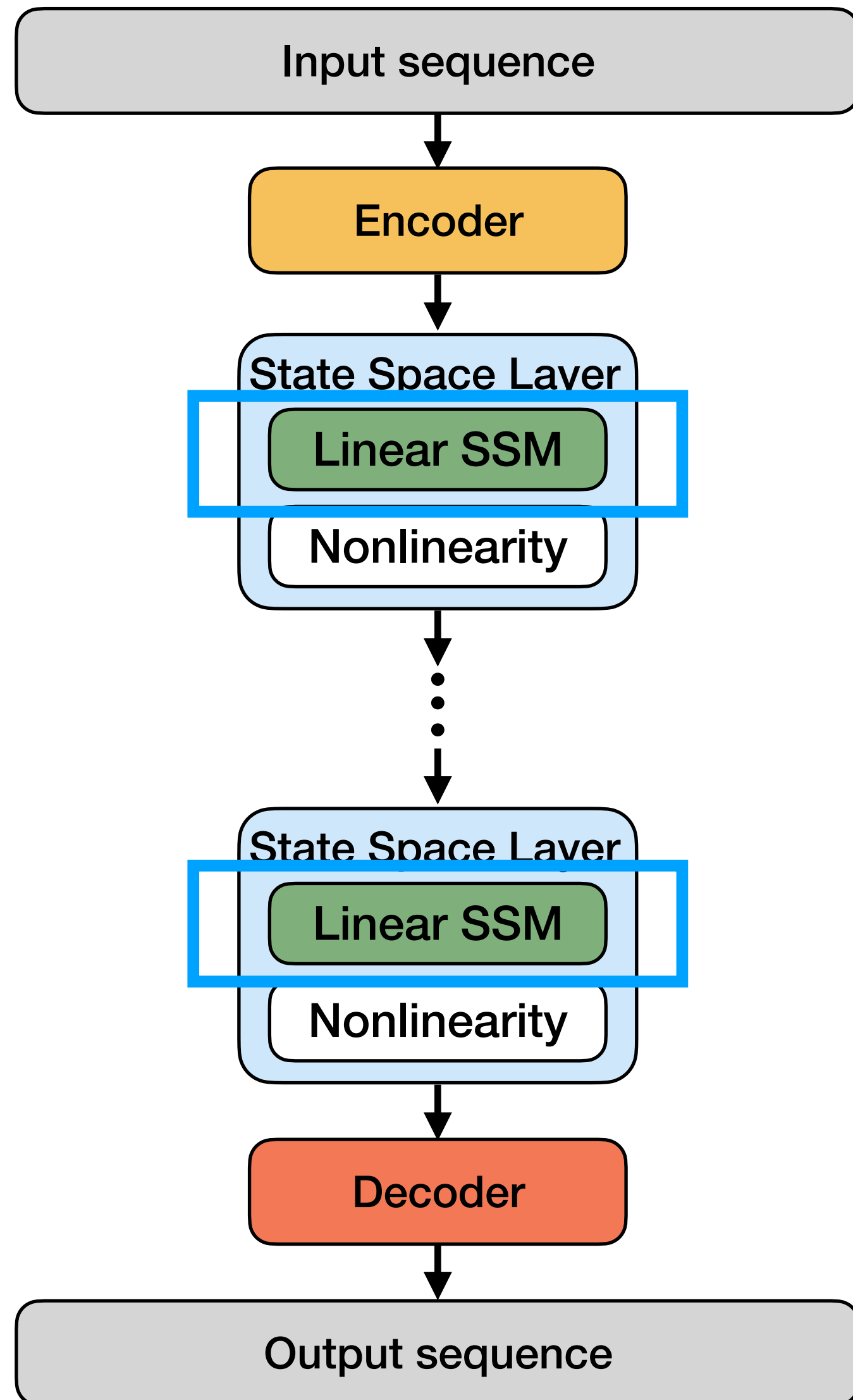
Key idea of Deep SSMs: Linear in time, nonlinear in depth



Key idea of Deep SSMs: Linear in time, nonlinear in depth

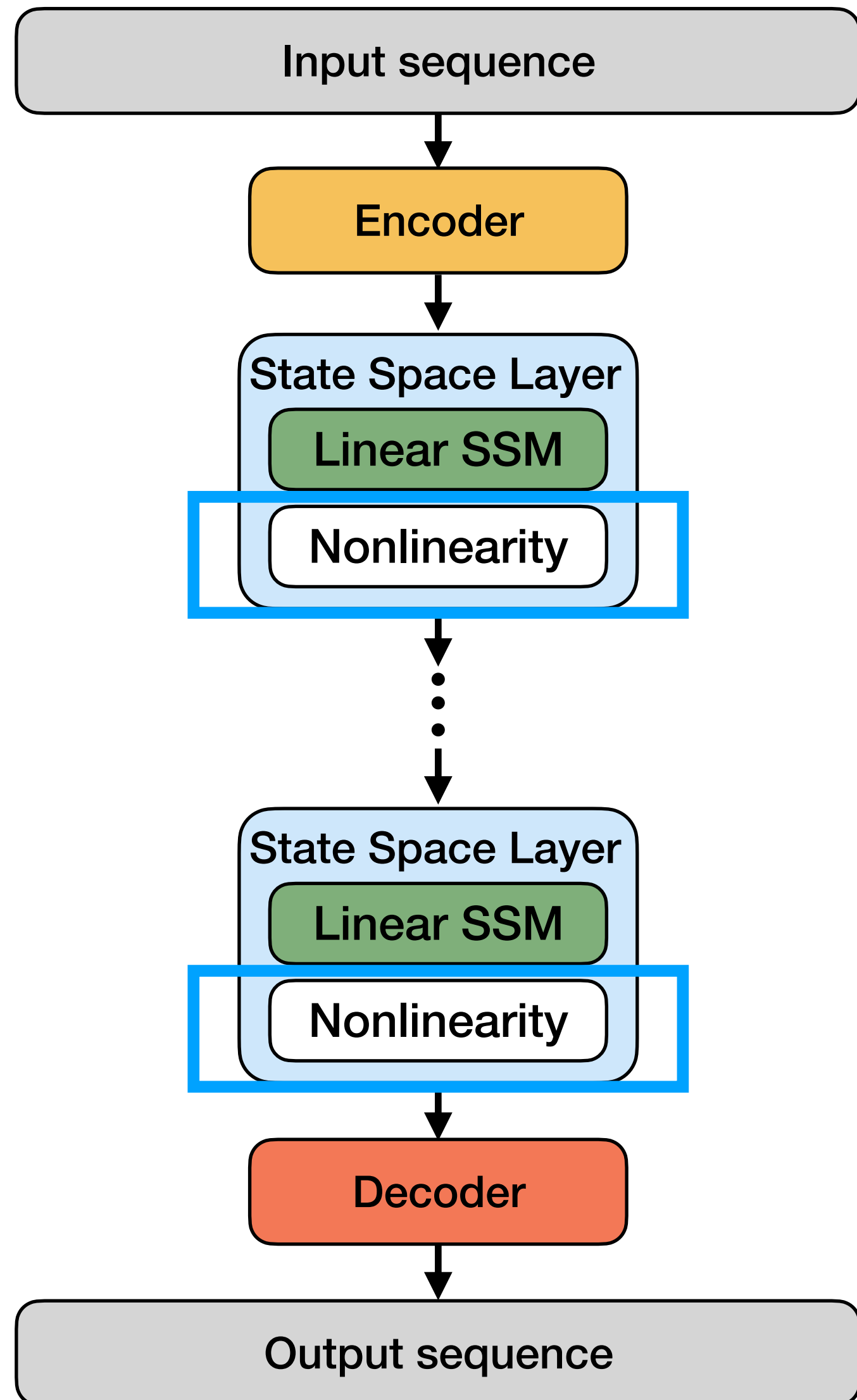


Key idea of Deep SSMs: Linear in time, nonlinear in depth



Linear in time: Efficient parallelization across the sequence

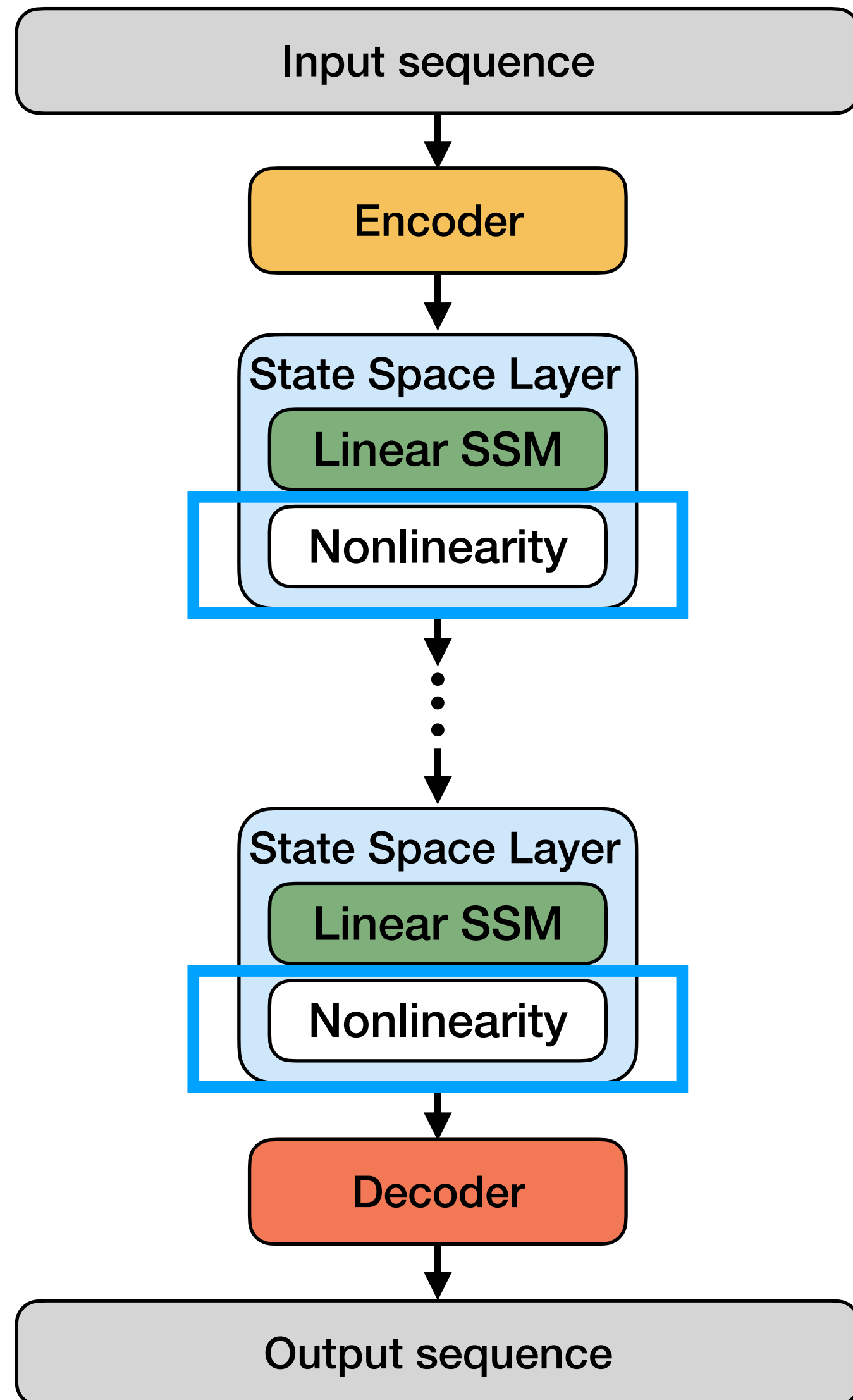
Key idea of Deep SSMs: Linear in time, nonlinear in depth



Linear in time: Efficient parallelization across the sequence

Nonlinear in depth: Stack of state space layers
can represent nonlinear systems

Key idea of Deep SSMs: Linear in time, nonlinear in depth



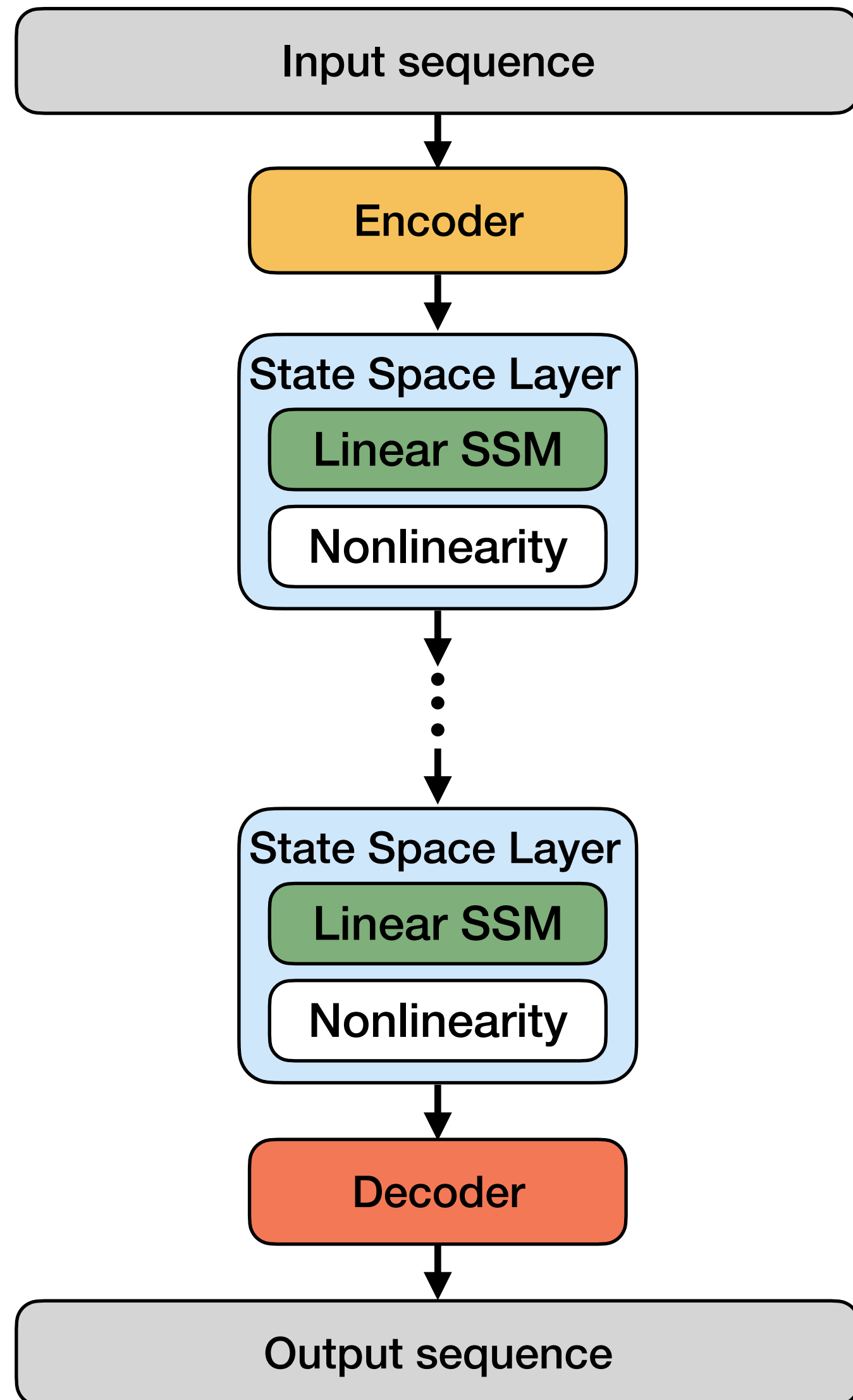
Linear in time: Efficient parallelization across the sequence

Nonlinear in depth: Stack of state space layers
can represent nonlinear systems

Expressivity Results:

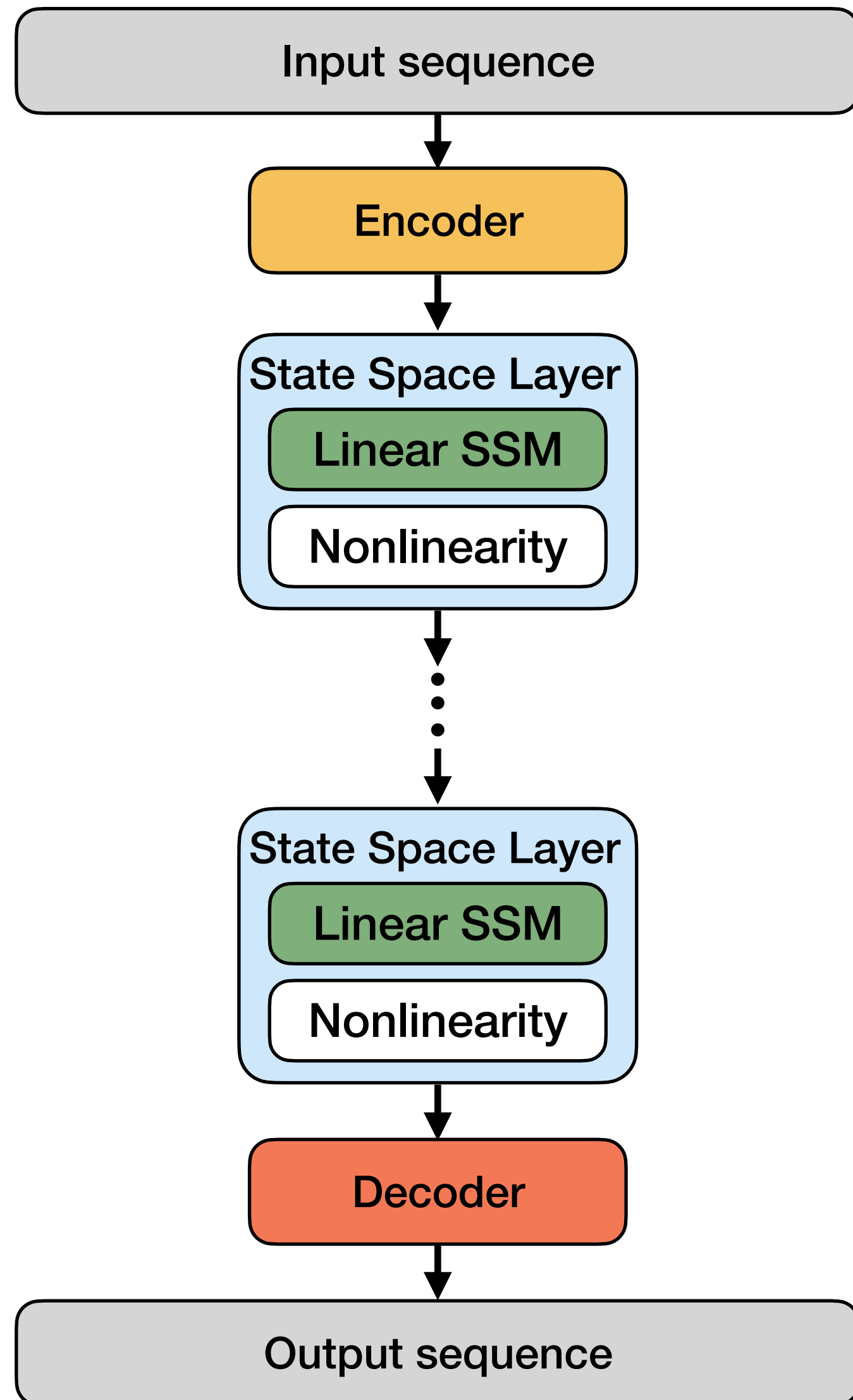
- Orvieto et al. 2023: <https://arxiv.org/abs/2307.11888>
- Wang et al. 2023: <https://arxiv.org/abs/2309.13414>

Key idea of Deep SSMs: Linear in time, nonlinear in depth



- Fast parallel processing
- Fast stateful autoregressive generation
- Can precisely parameterize to handle long-range dependencies (e.g. HiPPO framework, Gu et al. 2020)

Key idea of Deep SSMs: Linear in time, nonlinear in depth



Note, prior attempts at linear RNNs:

- QRNNs (Bradbury 2017)
- SRUs (Lei 2017)
- Linear surrogate RNNs (Martin 2018)

Likely reasons why more recent round of linear SSMs/RNNs have gained popularity:

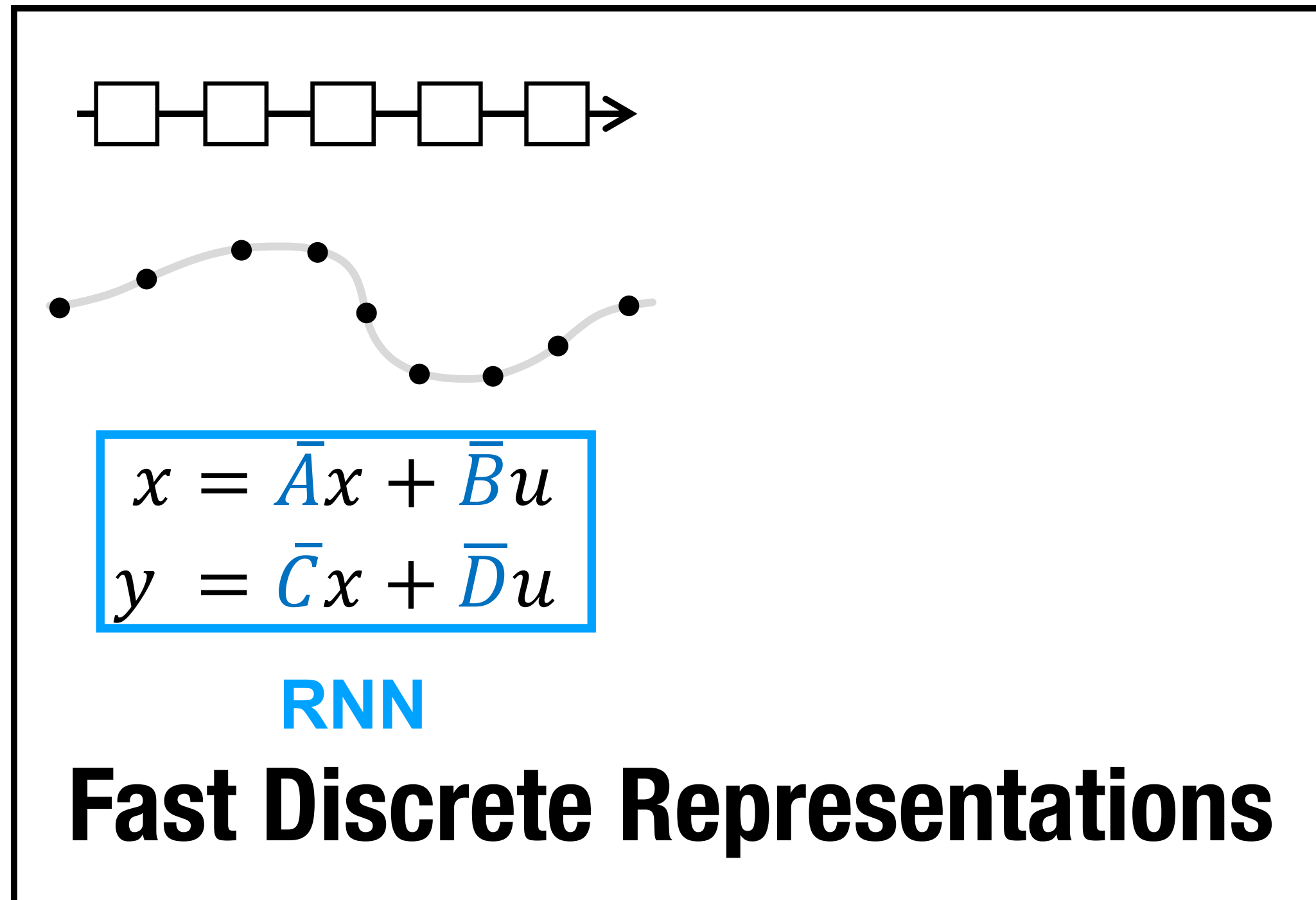
- Improved parameterizations
- Ideas from Transformers, e.g. backbones, layer normalizations, etc.
- Improved parallel algorithm implementations

Agenda

- Introduction, motivation, prior approaches
- Linear state space models (SSMs) overview
- **S4, convolutions, parameterization**
- S5, diagonalization, parallel scans
- S6/Mamba/Hawk/Griffin, data-dependent dynamics
- Challenges and opportunities

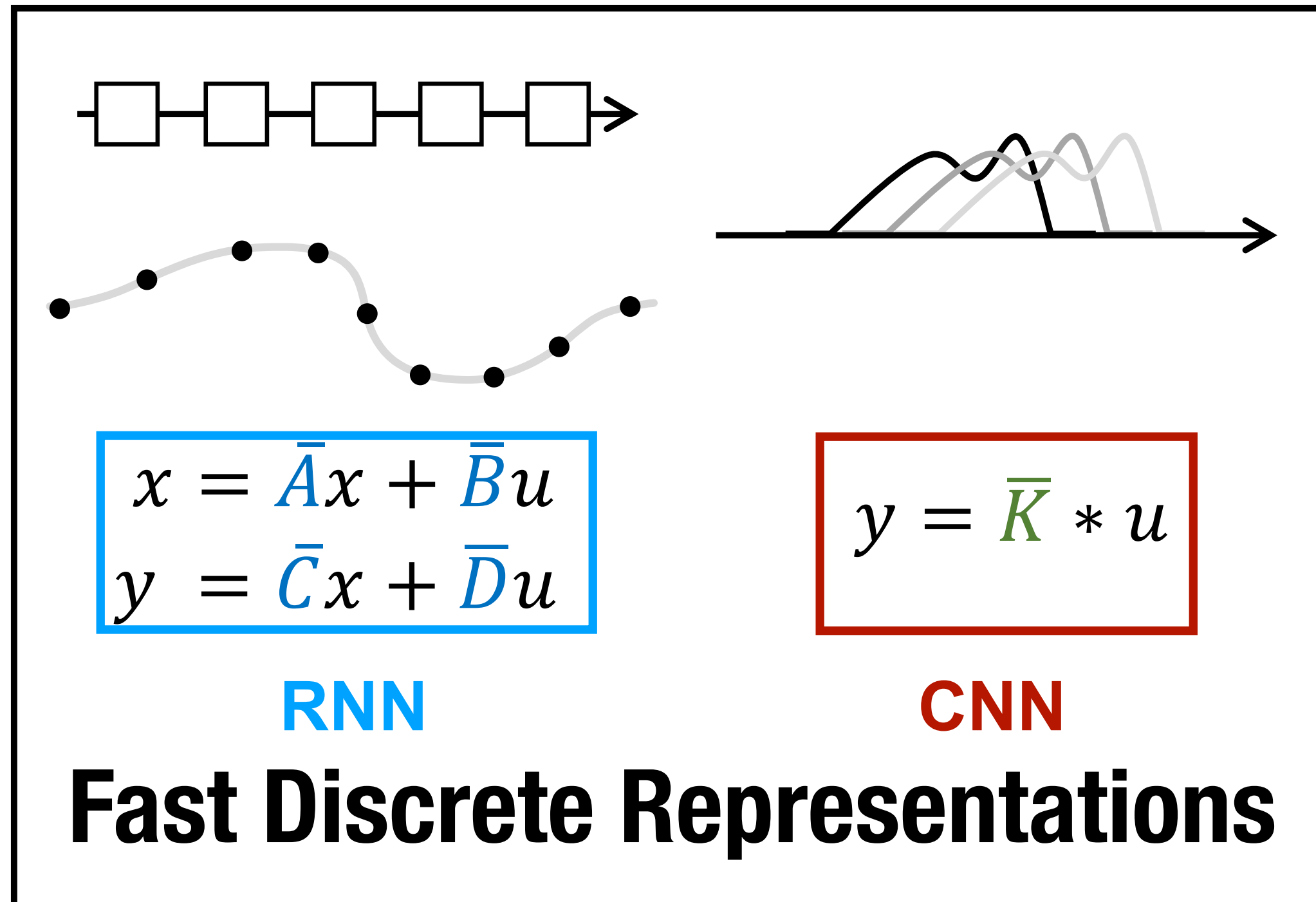
S4: Structured State Space Sequence Models

S4 can be an RNN



S4 can be run as either an **RNN** for fast autoregressive generation

S4 can be an RNN or a CNN



$$y_k = \bar{C}\bar{A}^k\bar{B}u_0 + \bar{C}\bar{A}^{k-1}\bar{B}u_1 + \cdots + \bar{C}\bar{A}\bar{B}u_{k-1} + \bar{C}\bar{B}u_k$$

$$y = \bar{K} * u.$$

Convolution

Kernel: $\bar{K} \in \mathbb{R}^L := (\bar{C}\bar{B}, \bar{C}\bar{A}\bar{B}, \dots, \bar{C}\bar{A}^{L-1}\bar{B})$

S4 can be run as either an **RNN** or a **CNN** for fast parallel processing

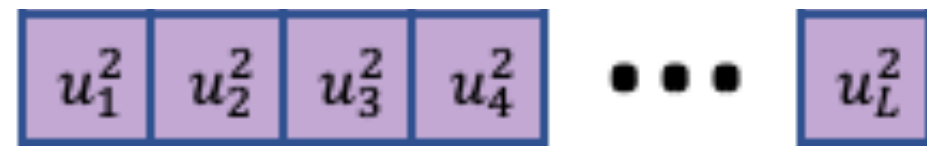
S4: Stack of single-input, single-output (SISO) SSMs

$$u_{1:L} \in \mathbb{R}^{L \times 3}$$

u_1^1	u_2^1	u_3^1	u_4^1	...	u_L^1
u_1^2	u_2^2	u_3^2	u_4^2		u_L^2
u_1^3	u_2^3	u_3^3	u_4^3		u_L^3

S4: Stack of single-input, single-output (SISO) SSMs

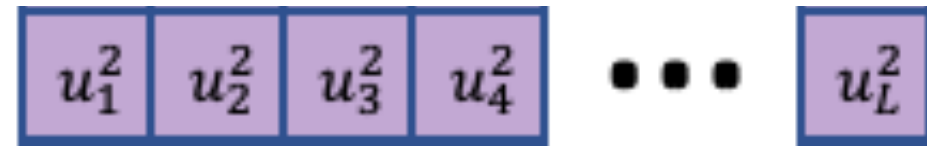
$$u_{1:L} \in \mathbb{R}^{L \times 3}$$



The diagram shows a sequence of input vectors. It starts with a row of four light purple boxes with dark blue borders, each containing the text u_1^2 , u_2^2 , u_3^2 , and u_4^2 respectively. To the right of these boxes are three black dots, followed by a single light purple box with a dark blue border containing the text u_L^2 .

S4: Stack of single-input, single-output (SISO) SSMs

$$u_{1:L} \in \mathbb{R}^{L \times 3}$$



A diagram showing a sequence of input vectors. It consists of four light blue boxes containing the expressions u_1^2 , u_2^2 , u_3^2 , and u_4^2 , followed by three black dots, and then a final light blue box containing u_L^2 .

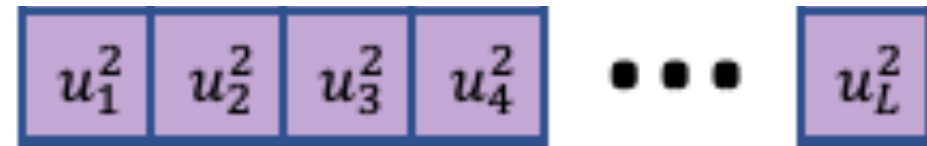
State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times 1}$

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

S4: Stack of single-input, single-output (SISO) SSMs

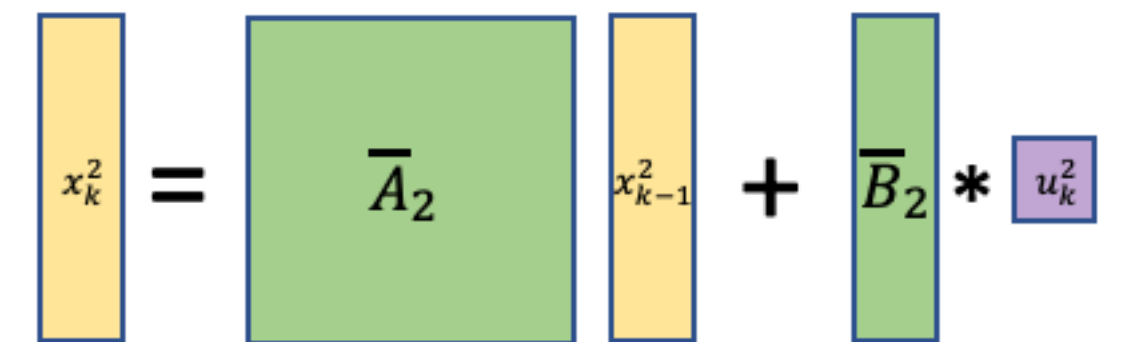
$$u_{1:L} \in \mathbb{R}^{L \times 3}$$



State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

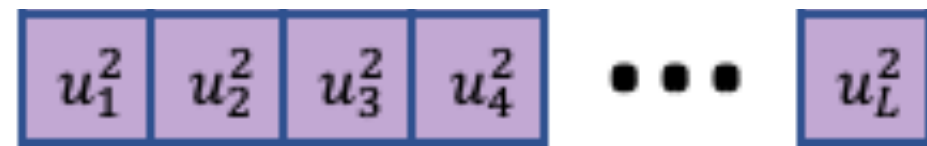
Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times 1}$

$$\mathbf{x}_k = \bar{\mathbf{A}}\mathbf{x}_{k-1} + \bar{\mathbf{B}}\mathbf{u}_k$$



S4: Stack of single-input, single-output (SISO) SSMs

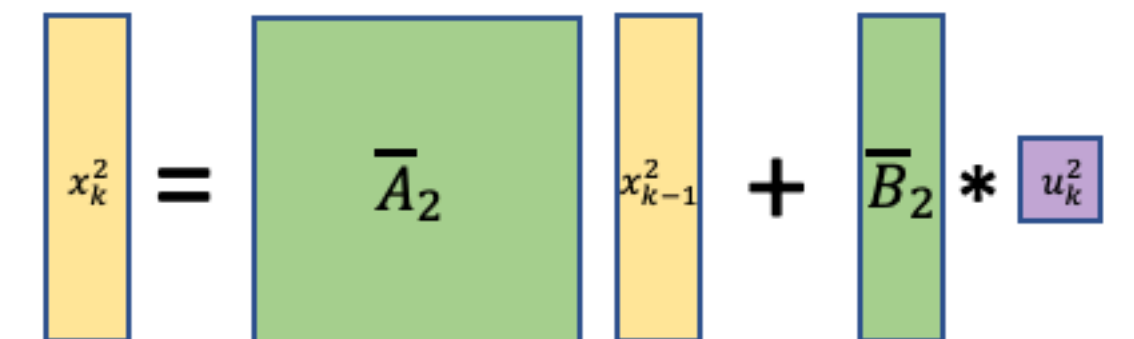
$$u_{1:L} \in \mathbb{R}^{L \times 3}$$



State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times 1}$

$$\mathbf{x}_k = \bar{\mathbf{A}}\mathbf{x}_{k-1} + \bar{\mathbf{B}}\mathbf{u}_k$$

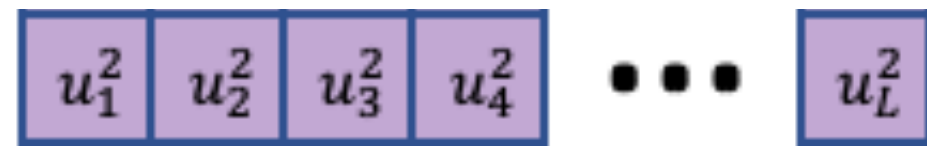


Output Matrix: $\mathbf{C} \in \mathbb{R}^{1 \times N}$

$$\mathbf{y}_k = \bar{\mathbf{C}}\mathbf{x}_k$$

S4: Stack of single-input, single-output (SISO) SSMs

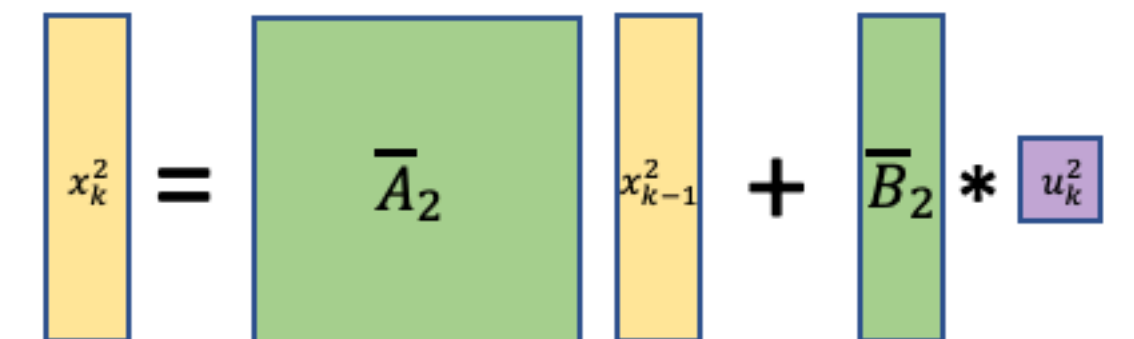
$$u_{1:L} \in \mathbb{R}^{L \times 3}$$



State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

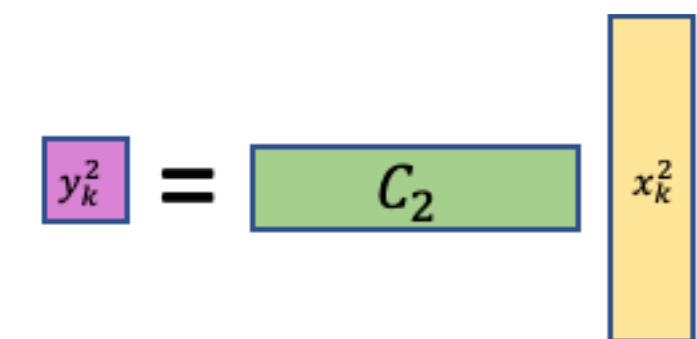
Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times 1}$

$$\mathbf{x}_k = \bar{\mathbf{A}}\mathbf{x}_{k-1} + \bar{\mathbf{B}}\mathbf{u}_k$$



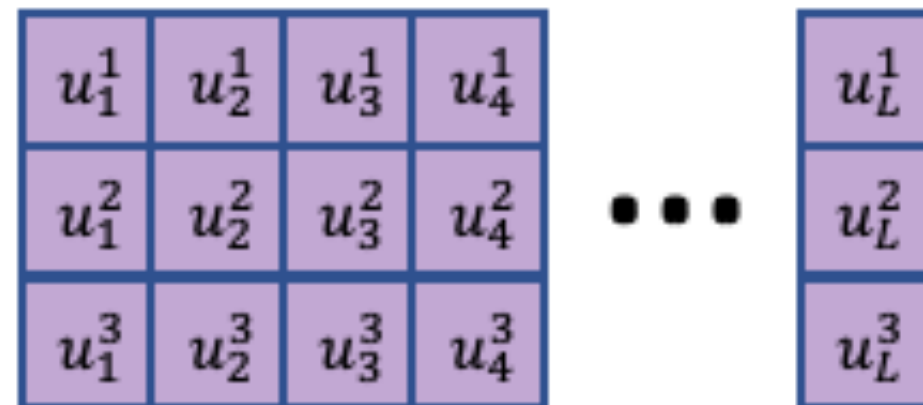
Output Matrix: $\mathbf{C} \in \mathbb{R}^{1 \times N}$

$$\mathbf{y}_k = \bar{\mathbf{C}}\mathbf{x}_k$$



S4: Stack of single-input, single-output (SISO) SSMs

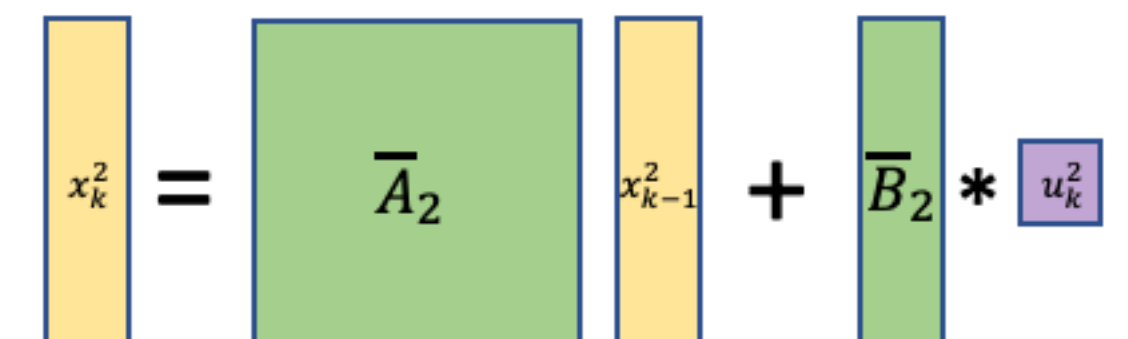
$$u_{1:L} \in \mathbb{R}^{L \times 3}$$



State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

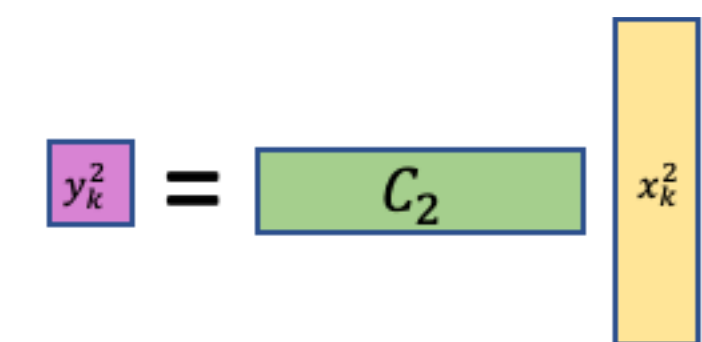
Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times 1}$

$$\mathbf{x}_k = \bar{\mathbf{A}}\mathbf{x}_{k-1} + \bar{\mathbf{B}}\mathbf{u}_k$$



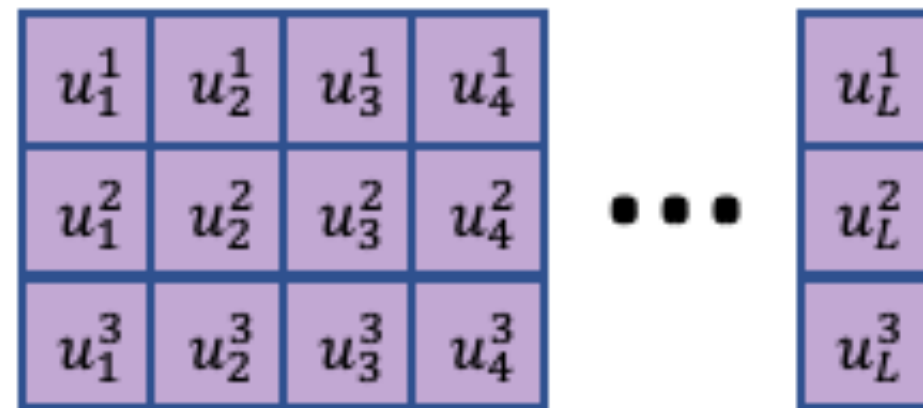
Output Matrix: $\mathbf{C} \in \mathbb{R}^{1 \times N}$

$$\mathbf{y}_k = \bar{\mathbf{C}}\mathbf{x}_k$$



S4: Stack of single-input, single-output (SISO) SSMs

$$u_{1:L} \in \mathbb{R}^{L \times 3}$$



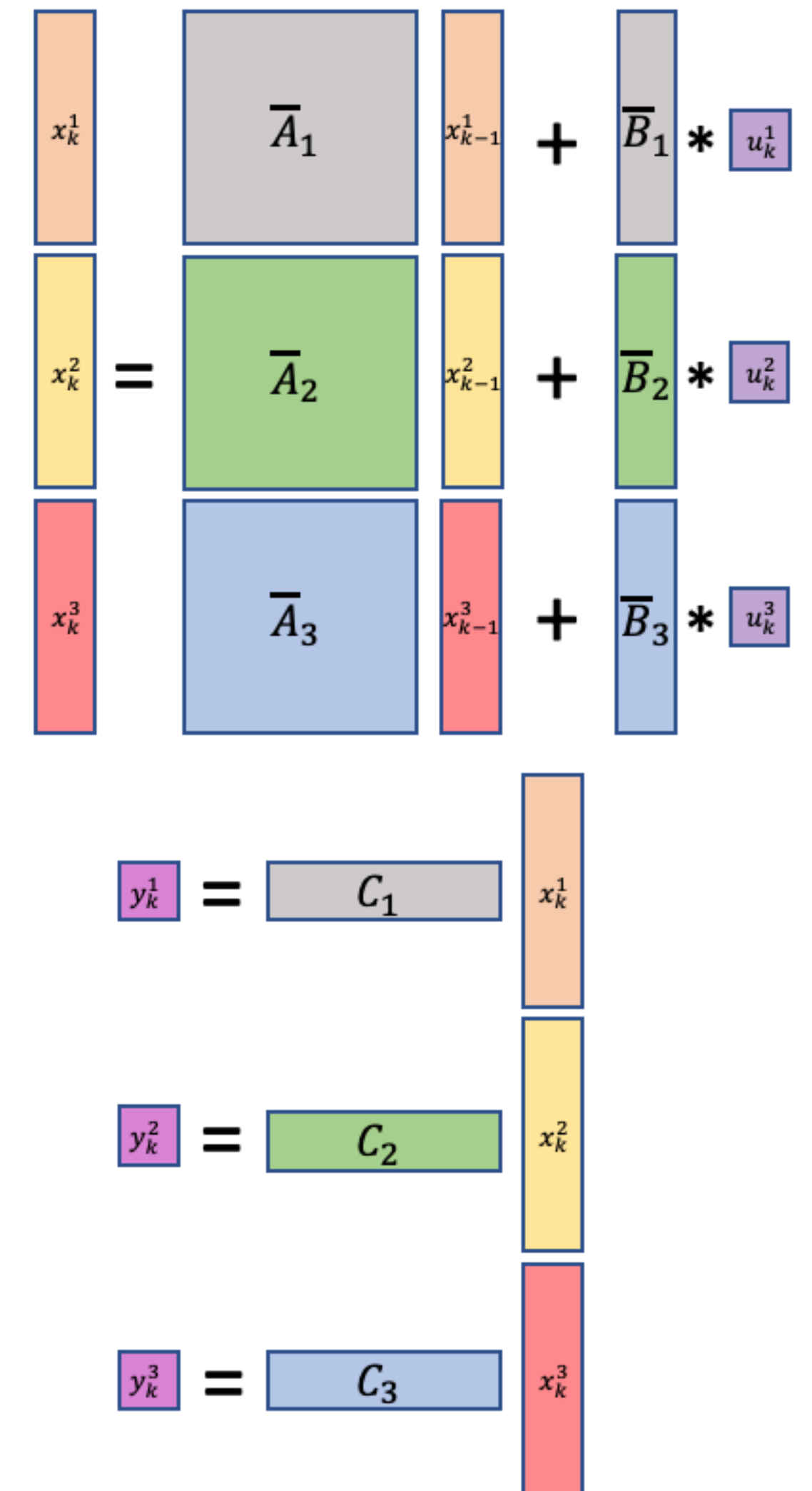
State Matrix: $\mathbf{A} \in \mathbb{R}^{N \times N}$

Input Matrix: $\mathbf{B} \in \mathbb{R}^{N \times 1}$

$$\mathbf{x}_k = \bar{\mathbf{A}}\mathbf{x}_{k-1} + \bar{\mathbf{B}}\mathbf{u}_k$$

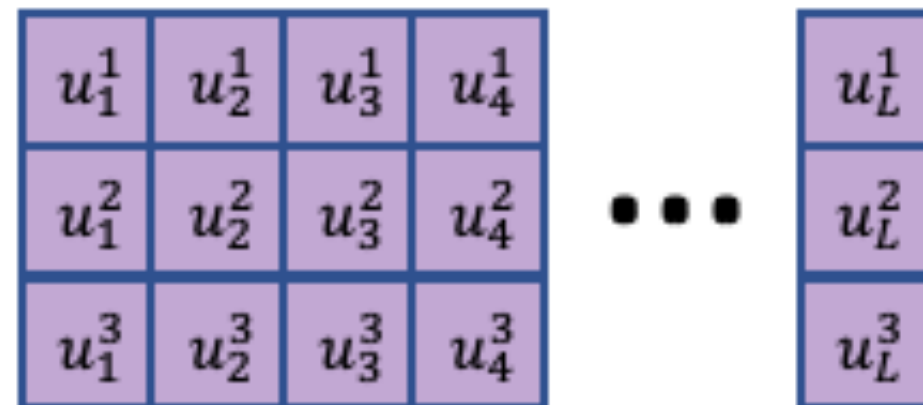
Output Matrix: $\mathbf{C} \in \mathbb{R}^{1 \times N}$

$$\mathbf{y}_k = \bar{\mathbf{C}}\mathbf{x}_k$$



S4: Stack of single-input, single-output (SISO) SSMs

$$u_{1:L} \in \mathbb{R}^{L \times 3}$$



State Matrix: $A \in \mathbb{R}^{N \times N}$

Input Matrix: $B \in \mathbb{R}^{N \times 1}$

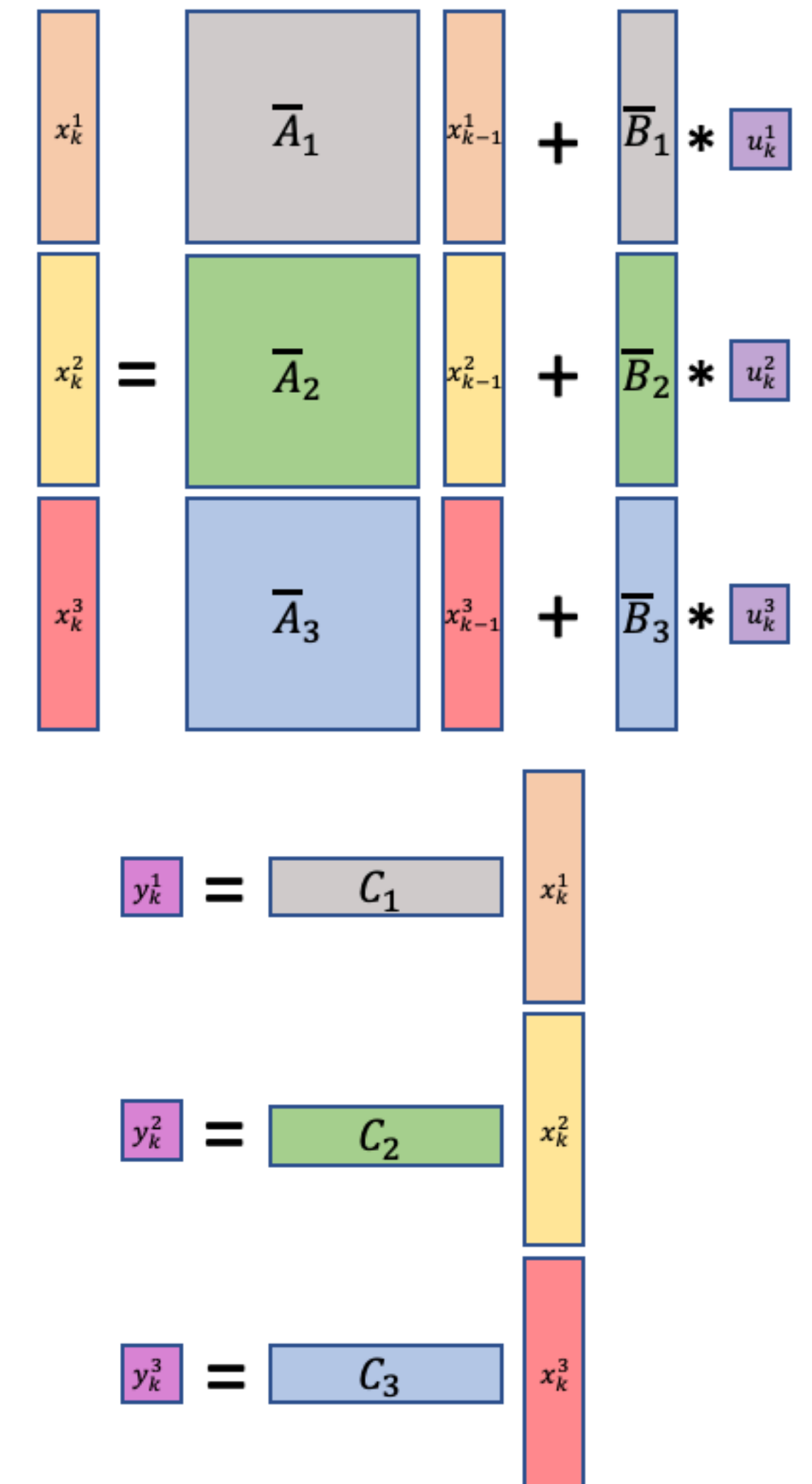
$$\mathbf{x}_k = \bar{\mathbf{A}}\mathbf{x}_{k-1} + \bar{\mathbf{B}}\mathbf{u}_k$$

Output Matrix: $C \in \mathbb{R}^{1 \times N}$

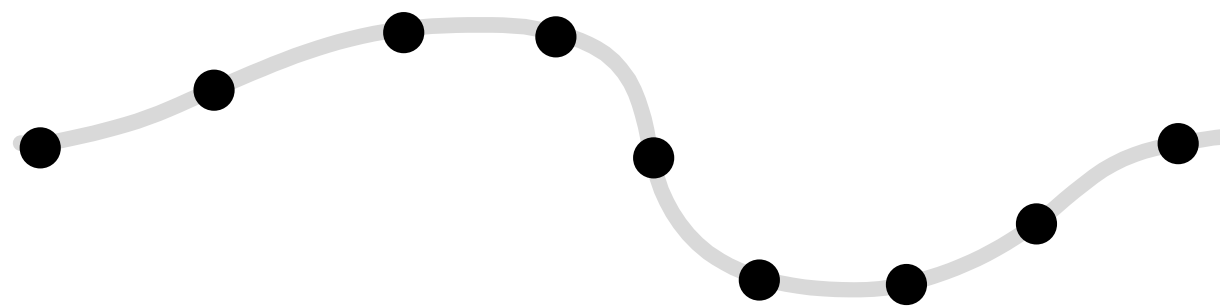
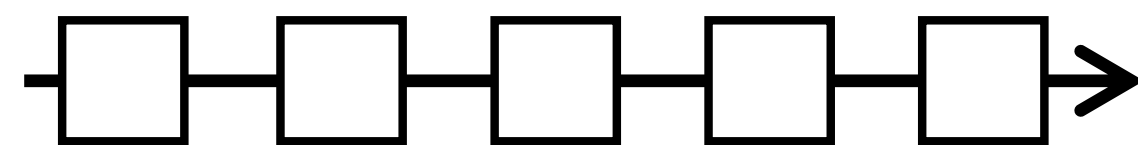
$$\mathbf{y}_k = \bar{\mathbf{C}}\mathbf{x}_k$$

Motivation for this structure:

- Computation (1D convolutions)
- Parameterization and Initialization (HiPPO, designed for SISO SSMs)
- Parameter efficient way to expand the state size



S4 Computation: RNN Mode



$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

Convolution equivalence holds for
any linear time-invariant (LTI) system

Unroll the recurrence:

$$x_0 = \overline{\mathbf{B}}u_0$$

$$y_0 = \overline{\mathbf{C}}\overline{\mathbf{B}}u_0$$

S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

Convolution equivalence holds for
any linear time-invariant (LTI) system

Unroll the recurrence:

$$x_0 = \overline{\mathbf{B}}u_0 \quad x_1 = \overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{B}}u_1$$

$$y_0 = \overline{\mathbf{C}}\overline{\mathbf{B}}u_0 \quad y_1 = \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_1$$

S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

Convolution equivalence holds for
any linear time-invariant (LTI) system

Unroll the recurrence:

$$\begin{aligned} x_0 &= \overline{\mathbf{B}}u_0 & x_1 &= \overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{B}}u_1 & x_2 &= \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_0 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2 \\ y_0 &= \overline{\mathbf{C}}\overline{\mathbf{B}}u_0 & y_1 &= \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_1 & y_2 &= \overline{\mathbf{C}}\overline{\mathbf{A}}^2\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_2 \end{aligned}$$

S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

Convolution equivalence holds for
any linear time-invariant (LTI) system

Unroll the recurrence:

$$\begin{aligned} x_0 &= \overline{\mathbf{B}}u_0 & x_1 &= \overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{B}}u_1 & x_2 &= \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_0 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2 \\ y_0 &= \overline{\mathbf{C}}\overline{\mathbf{B}}u_0 & y_1 &= \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_1 & y_2 &= \overline{\mathbf{C}}\overline{\mathbf{A}}^2\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_2 \end{aligned}$$

$$y_k = \overline{\mathbf{C}}\overline{\mathbf{A}}^k\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}^{k-1}\overline{\mathbf{B}}u_1 + \cdots + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_{k-1} + \overline{\mathbf{C}}\overline{\mathbf{B}}u_k$$

S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

Convolution equivalence holds for
any linear time-invariant (LTI) system

Unroll the recurrence:

$$\begin{aligned} x_0 &= \overline{\mathbf{B}}u_0 & x_1 &= \overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{B}}u_1 & x_2 &= \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_0 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2 \\ y_0 &= \overline{\mathbf{C}}\overline{\mathbf{B}}u_0 & y_1 &= \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_1 & y_2 &= \overline{\mathbf{C}}\overline{\mathbf{A}}^2\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_2 \end{aligned}$$

$$\begin{aligned} y_k &= \overline{\mathbf{C}}\overline{\mathbf{A}}^k\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}^{k-1}\overline{\mathbf{B}}u_1 + \cdots + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_{k-1} + \overline{\mathbf{C}}\overline{\mathbf{B}}u_k \\ y &= \overline{\mathbf{K}} * u. \end{aligned}$$

S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

Convolution equivalence holds for
any linear time-invariant (LTI) system

Unroll the recurrence:

$$x_0 = \overline{\mathbf{B}}u_0 \quad x_1 = \overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{B}}u_1 \quad x_2 = \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_0 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2$$

$$y_0 = \overline{\mathbf{C}}\overline{\mathbf{B}}u_0 \quad y_1 = \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_1 \quad y_2 = \overline{\mathbf{C}}\overline{\mathbf{A}}^2\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_2$$

$$y_k = \overline{\mathbf{C}}\overline{\mathbf{A}}^k\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}^{k-1}\overline{\mathbf{B}}u_1 + \cdots + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_{k-1} + \overline{\mathbf{C}}\overline{\mathbf{B}}u_k$$

$$y = \overline{\mathbf{K}} * u.$$

Convolution kernel: $\overline{\mathbf{K}} \in \mathbb{R}^L := (\overline{\mathbf{C}}\overline{\mathbf{B}}, \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}, \dots, \overline{\mathbf{C}}\overline{\mathbf{A}}^{L-1}\overline{\mathbf{B}})$

S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

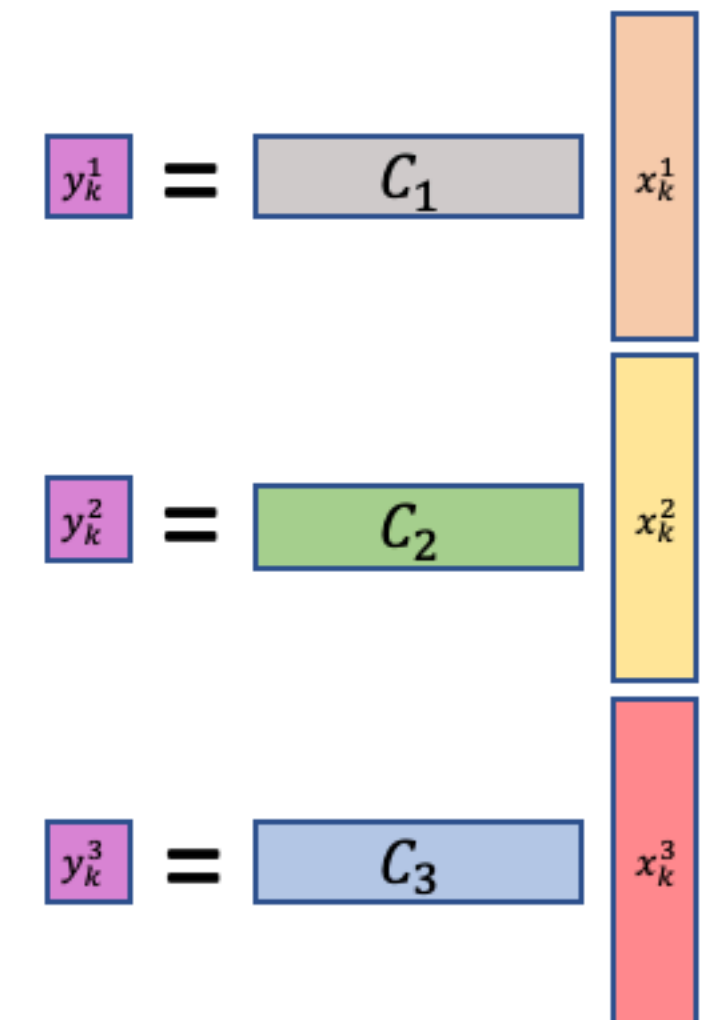
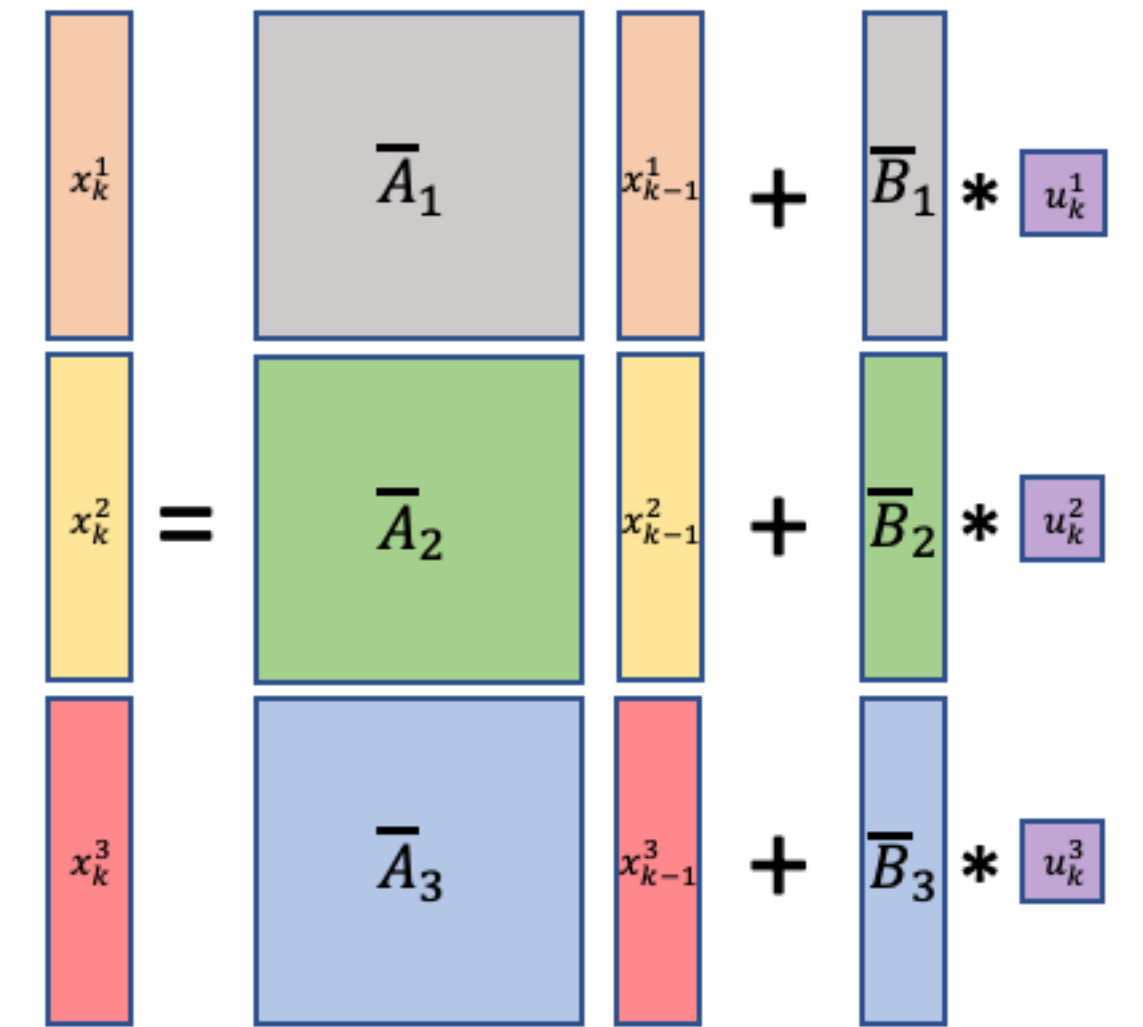
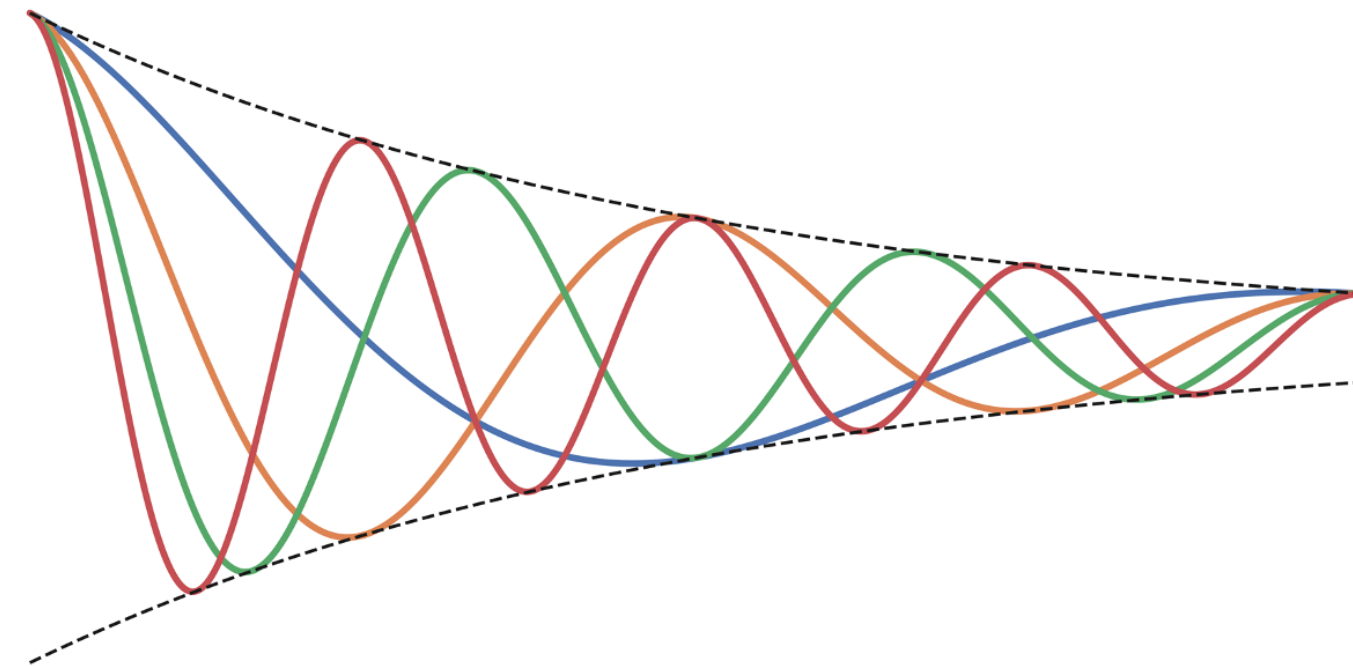
$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

$$y_k = \overline{\mathbf{C}}\overline{\mathbf{A}}^k\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}^{k-1}\overline{\mathbf{B}}u_1 + \cdots + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_{k-1} + \overline{\mathbf{C}}\overline{\mathbf{B}}u_k$$

$$y = \overline{\mathbf{K}} * u.$$

Convolution kernel: $\overline{\mathbf{K}} \in \mathbb{R}^L := (\overline{\mathbf{C}}\overline{\mathbf{B}}, \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}, \dots, \overline{\mathbf{C}}\overline{\mathbf{A}}^{L-1}\overline{\mathbf{B}})$

Stack of SISO SSMs gives a range of 1D convolution kernels that
that can capture different timescales:



S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

$$y_k = \overline{\mathbf{C}}\overline{\mathbf{A}}^k\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}^{k-1}\overline{\mathbf{B}}u_1 + \cdots + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_{k-1} + \overline{\mathbf{C}}\overline{\mathbf{B}}u_k$$

$$y = \overline{\mathbf{K}} * u.$$

Convolution kernel: $\overline{\mathbf{K}} \in \mathbb{R}^L := (\overline{\mathbf{C}}\overline{\mathbf{B}}, \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}, \dots, \overline{\mathbf{C}}\overline{\mathbf{A}}^{L-1}\overline{\mathbf{B}})$

Convolution Theorem:

$$\mathcal{F}[f * g] = \mathcal{F}[f]\mathcal{F}[g].$$

S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

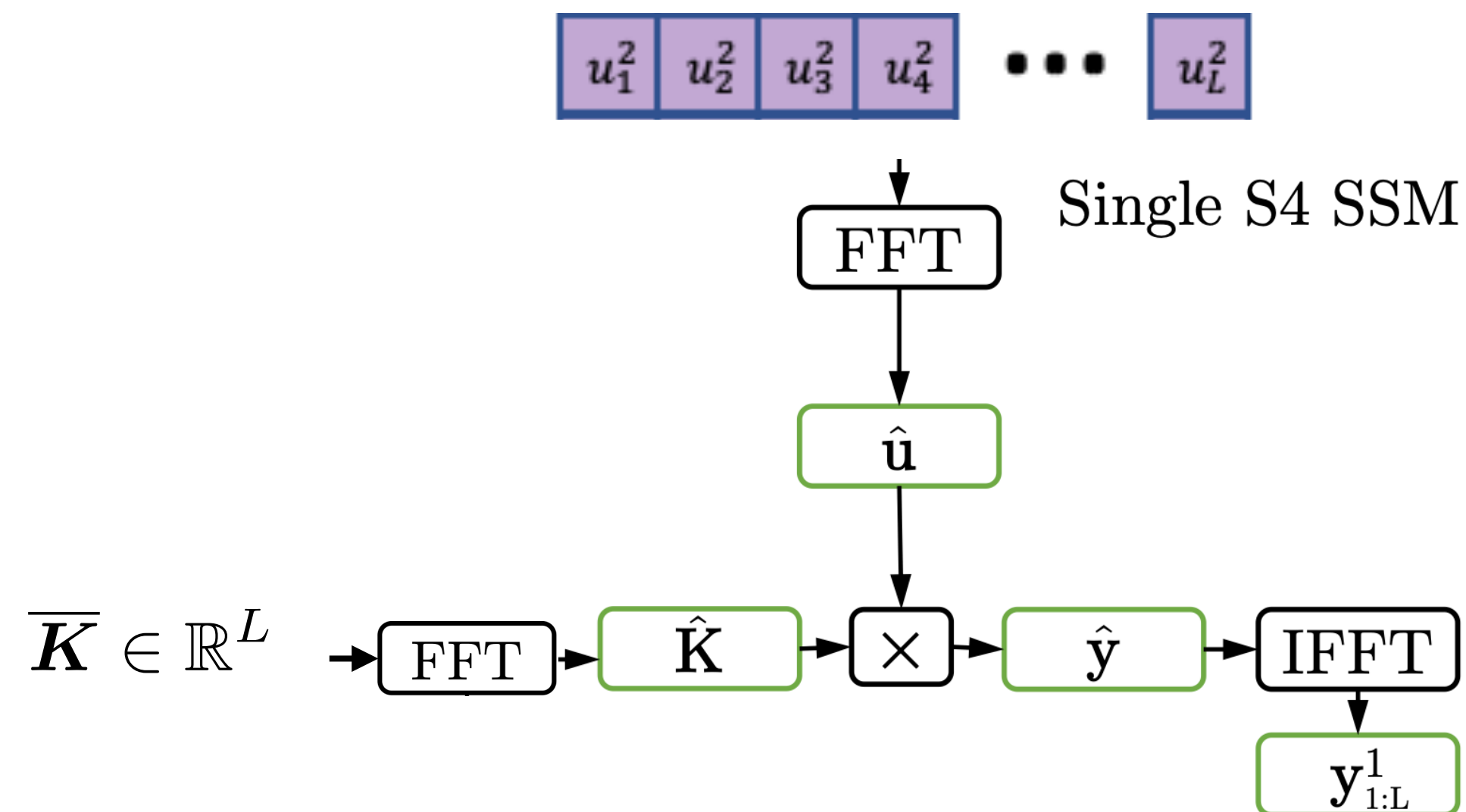
$$y_k = \overline{\mathbf{C}}\overline{\mathbf{A}}^k\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}^{k-1}\overline{\mathbf{B}}u_1 + \cdots + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_{k-1} + \overline{\mathbf{C}}\overline{\mathbf{B}}u_k$$

$$y = \overline{\mathbf{K}} * u.$$

Convolution kernel: $\overline{\mathbf{K}} \in \mathbb{R}^L := (\overline{\mathbf{C}}\overline{\mathbf{B}}, \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}, \dots, \overline{\mathbf{C}}\overline{\mathbf{A}}^{L-1}\overline{\mathbf{B}})$

Convolution Theorem:

$$\mathcal{F}[f * g] = \mathcal{F}[f]\mathcal{F}[g].$$



S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

$$\mathbf{y}_k = \overline{\mathbf{C}}\mathbf{x}_k + \overline{\mathbf{D}}\mathbf{u}_k$$

$$y_k = \overline{\mathbf{C}}\overline{\mathbf{A}}^k\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}^{k-1}\overline{\mathbf{B}}u_1 + \cdots + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_{k-1} + \overline{\mathbf{C}}\overline{\mathbf{B}}u_k$$

$$y = \overline{\mathbf{K}} * u.$$

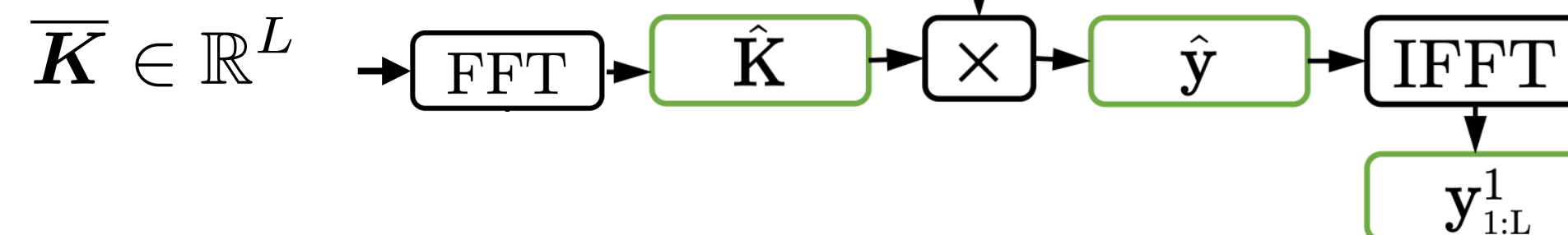
Convolution kernel: $\overline{\mathbf{K}} \in \mathbb{R}^L := (\overline{\mathbf{C}}\overline{\mathbf{B}}, \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}, \dots, \overline{\mathbf{C}}\overline{\mathbf{A}}^{L-1}\overline{\mathbf{B}})$

Convolution Theorem:

$$\mathcal{F}[f * g] = \mathcal{F}[f]\mathcal{F}[g].$$

Given kernel, the convolution can be computed with $O(L \log L)$ cost and $O(L)$ space

Importantly, can be parallelized across the sequence



S4 Computation: Convolution Mode

Consider a single S4 SSM:

$$y_k = \overline{CA}^k \overline{B} u_0 + \overline{CA}^{k-1} \overline{B} u_1 + \dots + \overline{CAB} u_{k-1} + \overline{CB} u_k$$

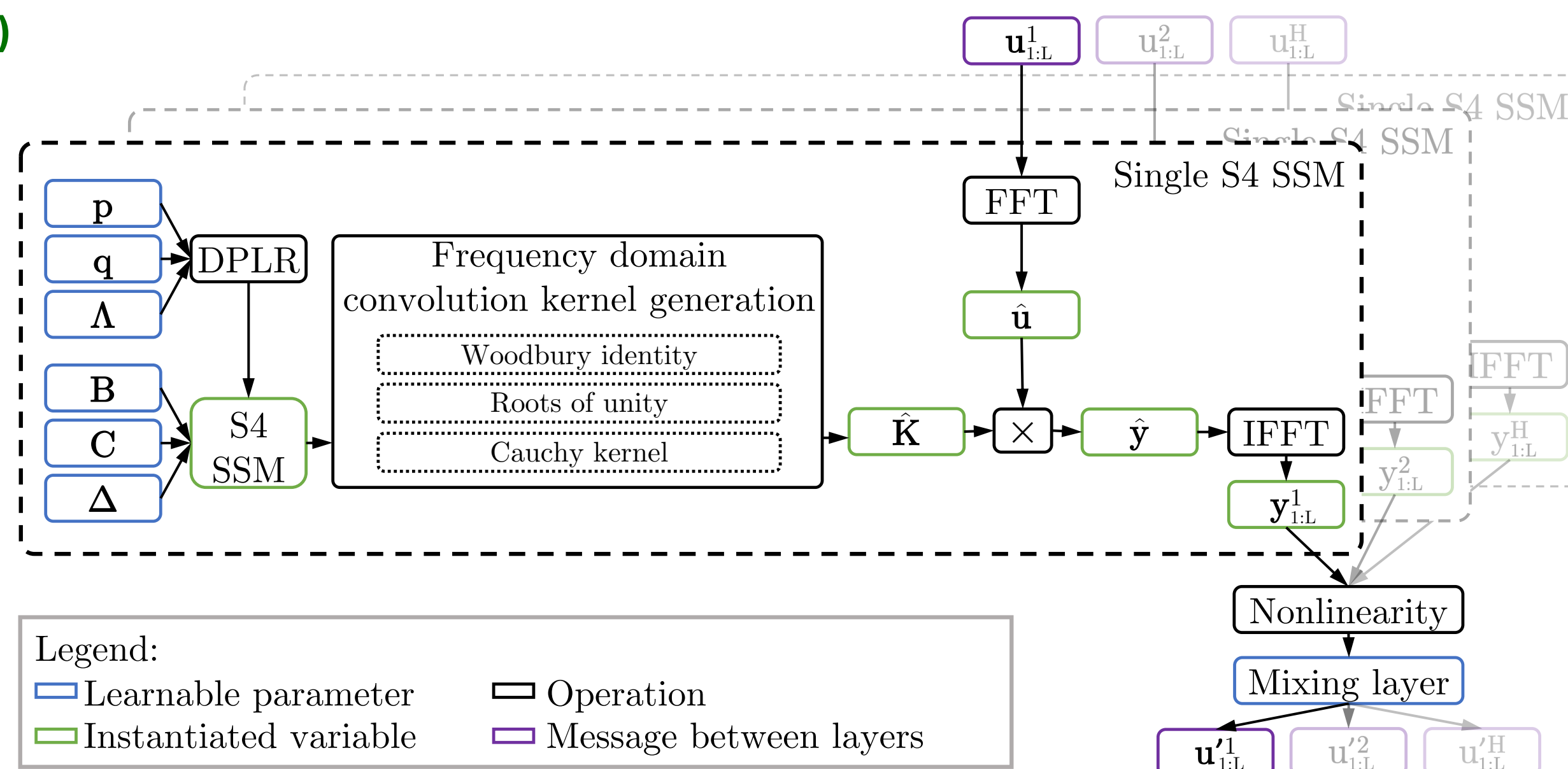
$$y = \overline{K} * u.$$

Convolution kernel: $\overline{K} \in \mathbb{R}^L := (\overline{CB}, \overline{CAB}, \dots, \overline{CA}^{L-1} \overline{B})$

Total cost of S4 layer: $\mathcal{O}(\boxed{H^2 L} + \boxed{HL \log L})$ (If we can compute kernel efficiently.... but this requires successive powers of A...)

Nonlinear FFN Convolutions (H channels)

- Naively, computing the kernel requires $\mathcal{O}(N^2 L)$ operations
- If dynamics matrix is diagonal, Vandermonde matrices can be used, $\mathcal{O}(NL)$ time and space naively, but can in theory be cheaper
- S4 used a diagonal plus low rank (DPLR) dynamics matrix, so required a sophisticated algorithm which resulted in the use of Cauchy kernels

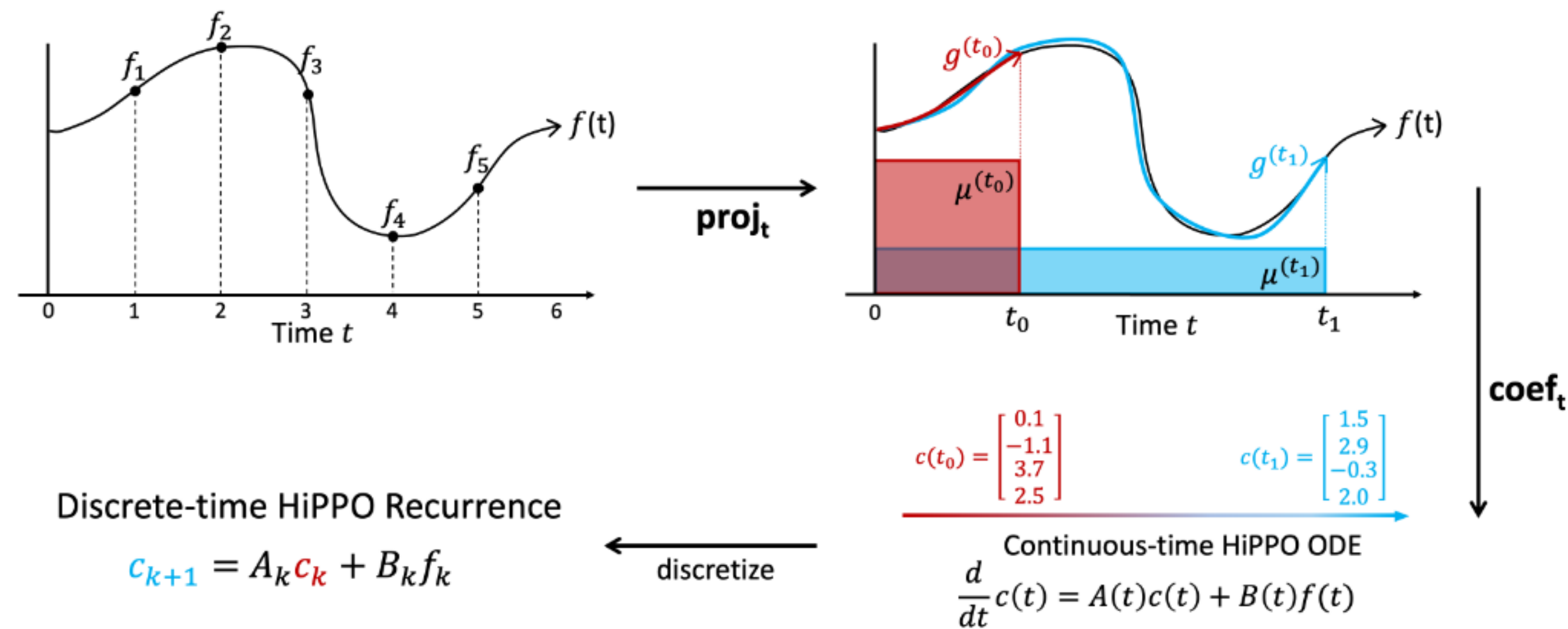


S4 Parameterization and Initialization

Main idea: Using HiPPO Theory (Gu et al. 2020), can represent history of a scalar signal using a SISO linear SSM with special state matrix, HiPPO matrices.

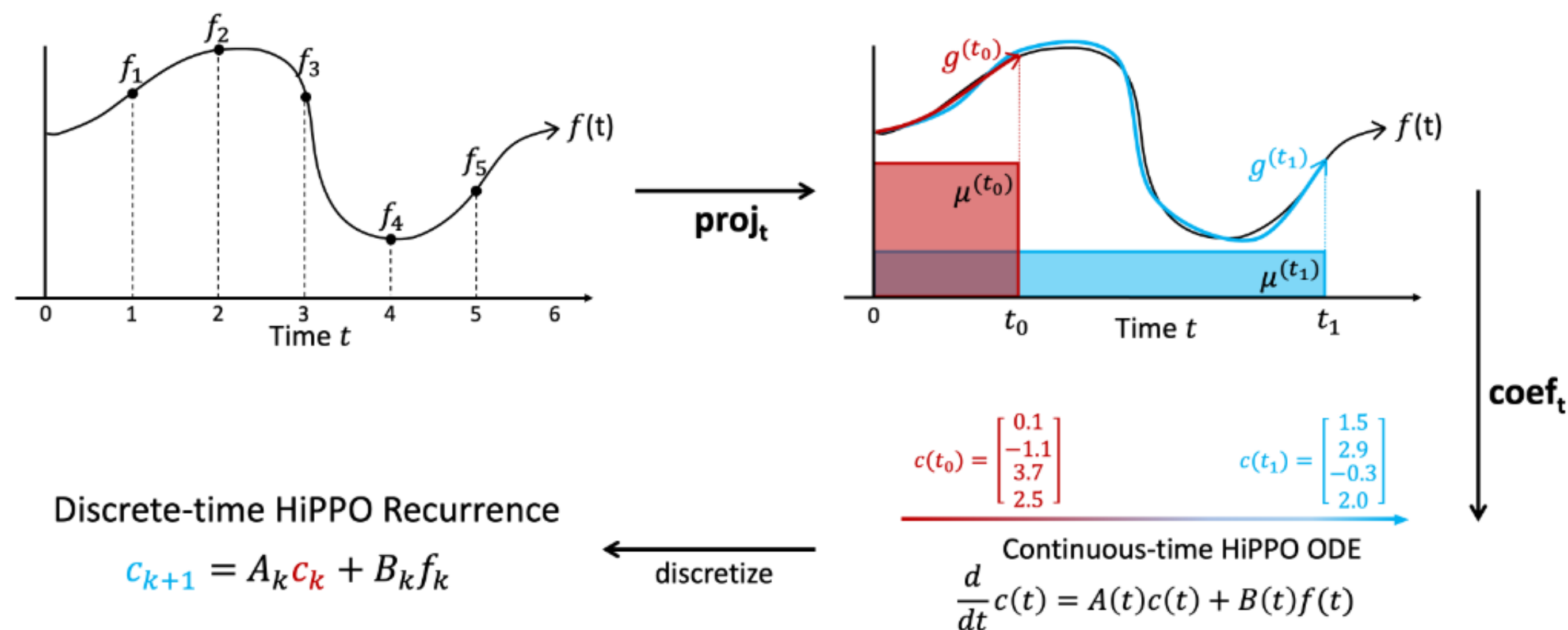
S4 Parameterization and Initialization

Main idea: Using HiPPO Theory (Gu et al. 2020), can represent history of a scalar signal using a SISO linear SSM with special state matrix, HiPPO matrices.



S4 Parameterization and Initialization

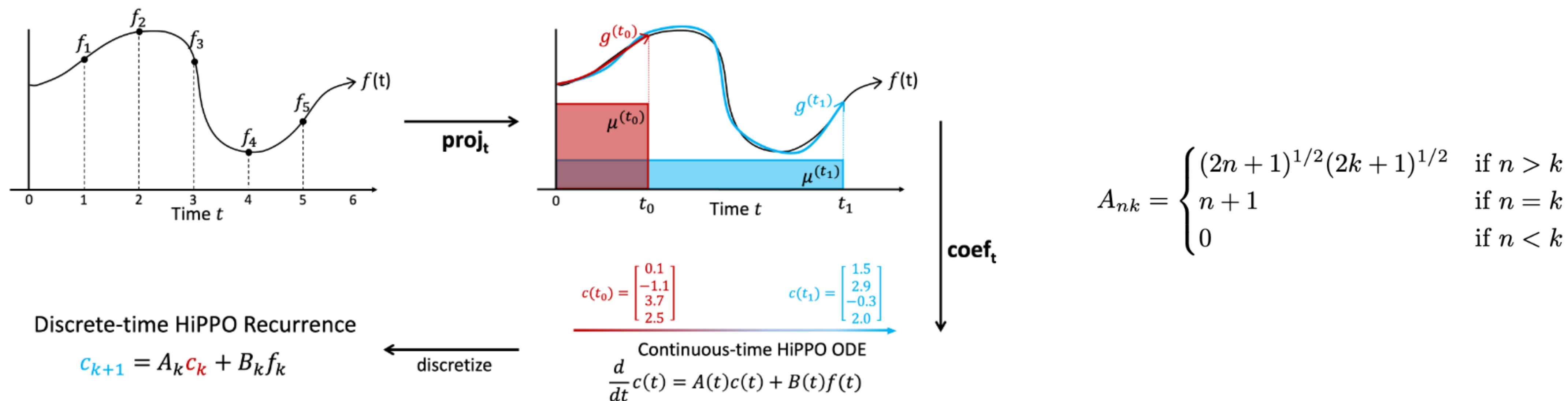
Main idea: Using HiPPO Theory (Gu et al. 2020), can represent history of a scalar signal using a SISO linear SSM with special state matrix, HiPPO matrices.



Intuitively, we can think of the memory representation $c(t) \in \mathbb{R}^N$ as being the *coefficient vector* of the optimal polynomial approximation to the history of $f(t)$.

S4 Parameterization and Initialization

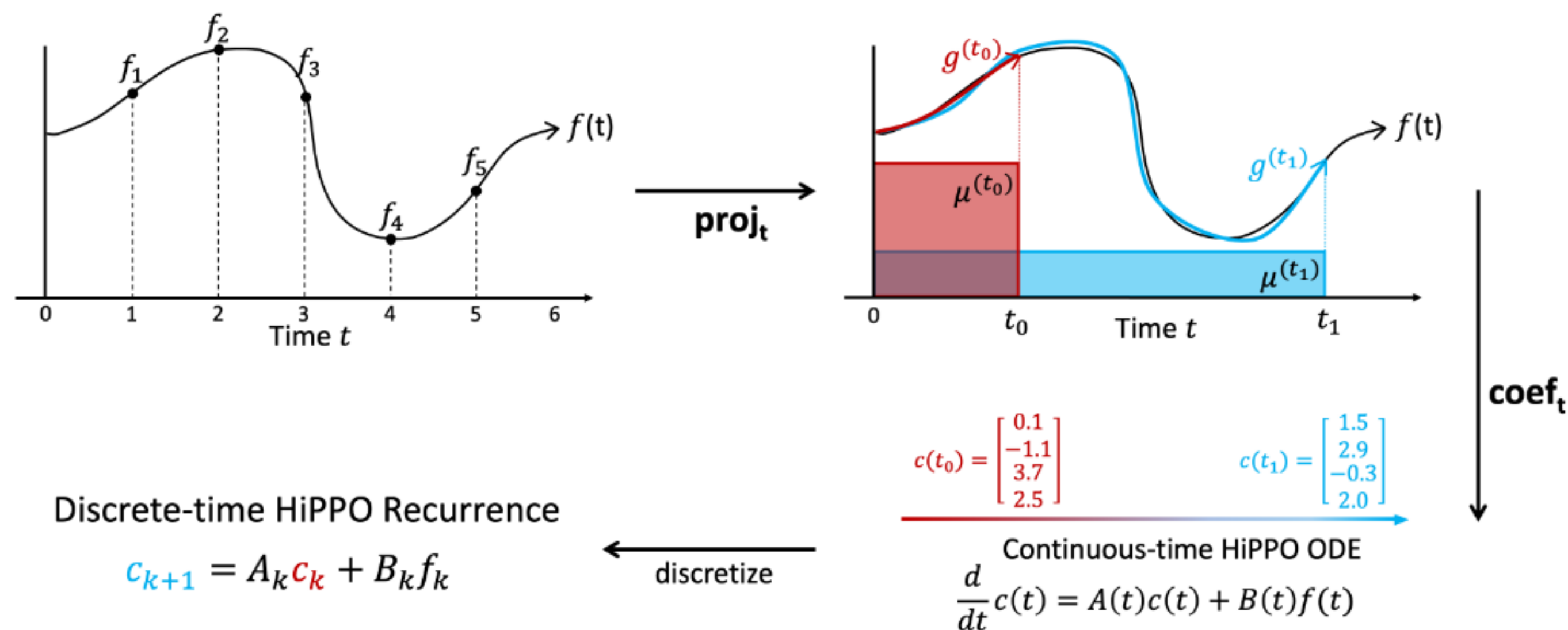
Main idea: Using HiPPO Theory (Gu et al. 2020), can represent history of a scalar signal using a SISO linear SSM with special state matrix, HiPPO matrices.



Intuitively, we can think of the memory representation $c(t) \in \mathbb{R}^N$ as being the *coefficient vector* of the optimal polynomial approximation to the history of $f(t)$.

S4 Parameterization and Initialization

Main idea: Using HiPPO Theory (Gu et al. 2020), can represent history of a scalar signal using a SISO linear SSM with special state matrix, HiPPO matrices.



$$A_{nk} = \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

S4 work found using these matrices were really important for LRA

Intuitively, we can think of the memory representation $c(t) \in \mathbb{R}^N$ as being the *coefficient vector* of the optimal polynomial approximation to the history of $f(t)$.

S4 Parameterization and Initialization

Main idea: Using HiPPO Theory (Gu et al. 2020), can represent history of a scalar signal using a SISO linear SSM with special state matrix, HiPPO matrices.

$$A_{nk} = \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

S4 work found using these matrices were really important for LRA

$$\mathbf{A}_{\text{LegS}} = \mathbf{A}_{\text{LegS}}^{\text{Normal}} - \mathbf{P}_{\text{Legs}} \mathbf{P}_{\text{Legs}}^{\top}$$

S4 Parameterization and Initialization

Main idea: Using HiPPO Theory (Gu et al. 2020), can represent history of a scalar signal using a SISO linear SSM with special state matrix, HiPPO matrices.

$$A_{nk} = \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

S4 work found using these matrices were really important for LRA

$$\mathbf{A}_{\text{LegS}} = \mathbf{A}_{\text{LegS}}^{\text{Normal}} - \mathbf{P}_{\text{Legs}} \mathbf{P}_{\text{Legs}}^{\top}$$

These matrices cannot be diagonalized numerically, but can be conjugated into diagonal plus low-rank form:

S4 Parameterization and Initialization

Main idea: Using HiPPO Theory (Gu et al. 2020), can represent history of a scalar signal using a SISO linear SSM with special state matrix, HiPPO matrices.

$$A_{nk} = \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

S4 work found using these matrices were really important for LRA

$$\mathbf{A}_{\text{LegS}} = \mathbf{A}_{\text{LegS}}^{\text{Normal}} - \mathbf{P}_{\text{Legs}} \mathbf{P}_{\text{Legs}}^{\top}$$

These matrices cannot be diagonalized numerically, but can be conjugated into diagonal plus low-rank form:

From S4 paper:

Theorem 1. *All HiPPO matrices from [16] have a NPLR representation*

$$\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^* - \mathbf{P} \mathbf{Q}^{\top} = \mathbf{V} (\mathbf{\Lambda} - (\mathbf{V}^* \mathbf{P}) (\mathbf{V}^* \mathbf{Q})^*) \mathbf{V}^* \quad (6)$$

for unitary $\mathbf{V} \in \mathbb{C}^{N \times N}$, diagonal $\mathbf{\Lambda}$, and low-rank factorization $\mathbf{P}, \mathbf{Q} \in \mathbb{R}^{N \times r}$. These matrices HiPPO- LegS, LegT, LagT all satisfy $r = 1$ or $r = 2$. In particular, equation (2) is NPLR with $r = 1$.

S4 Parameterization and Initialization

Main idea: Using HiPPO Theory (Gu et al. 2020), can represent history of a scalar signal using a SISO linear SSM with special state matrix, HiPPO matrices.

$$A_{nk} = \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

S4 work found using these matrices were really important for LRA

$$\mathbf{A}_{\text{LegS}} = \mathbf{A}_{\text{LegS}}^{\text{Normal}} - \mathbf{P}_{\text{LegS}} \mathbf{P}_{\text{LegS}}^{\top}$$

These matrices cannot be diagonalized numerically, but can be conjugated into diagonal plus low-rank form:

From S4 paper:

Theorem 1. All HiPPO matrices from [16] have a NPLR representation

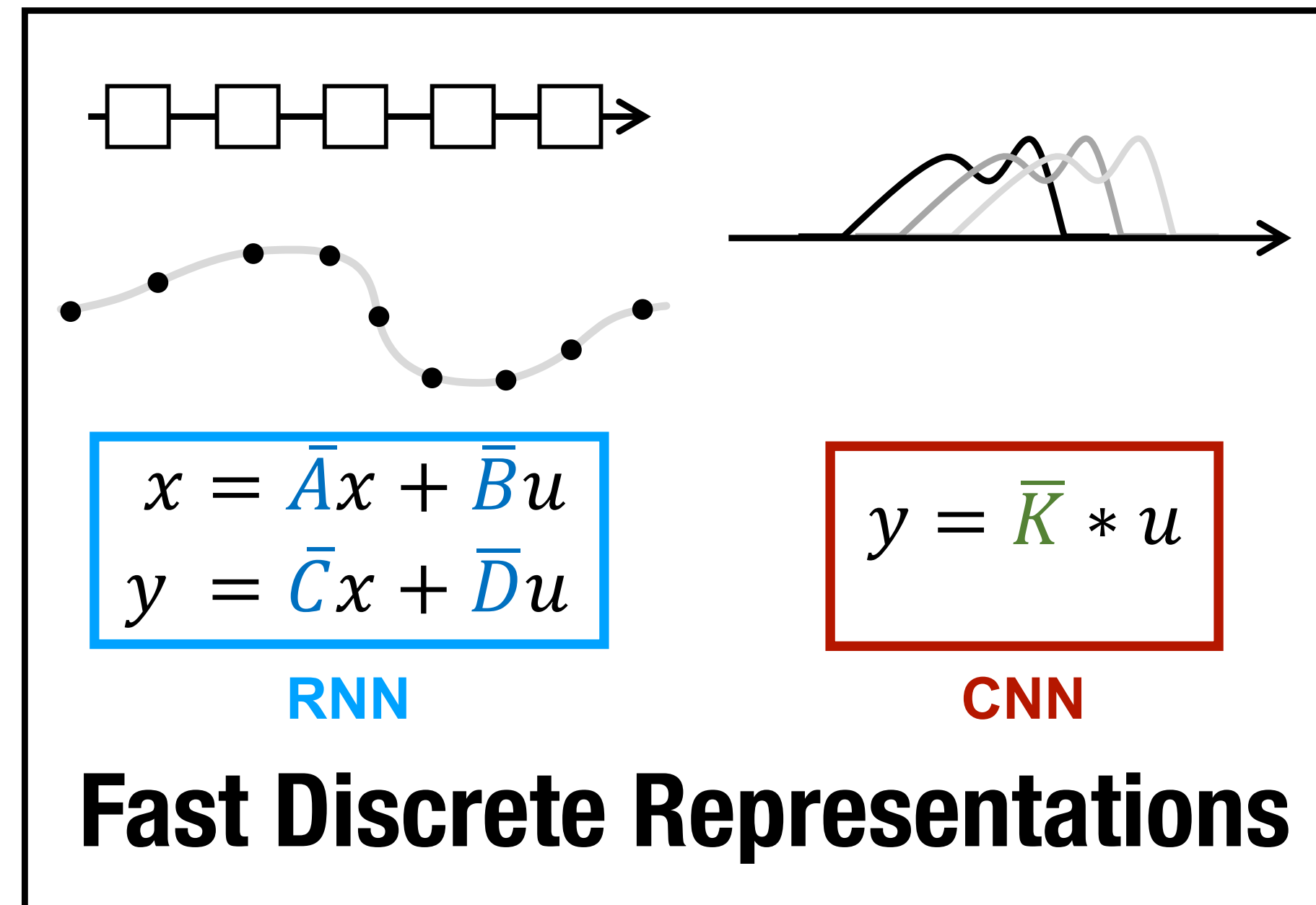
$$\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^* - \mathbf{P} \mathbf{Q}^{\top} = \mathbf{V} (\mathbf{\Lambda} - (\mathbf{V}^* \mathbf{P}) (\mathbf{V}^* \mathbf{Q})^*) \mathbf{V}^* \quad (6)$$

for unitary $\mathbf{V} \in \mathbb{C}^{N \times N}$, diagonal $\mathbf{\Lambda}$, and low-rank factorization $\mathbf{P}, \mathbf{Q} \in \mathbb{R}^{N \times r}$. These matrices HiPPO- LegS, LegT, LagT all satisfy $r = 1$ or $r = 2$. In particular, equation (2) is NPLR with $r = 1$.

$$\mathbf{A} = \begin{array}{|c|} \hline \text{diagonal matrix} \\ \hline \end{array} + \begin{array}{|c|} \hline \text{low-rank matrix} \\ \hline \end{array}$$

We will discuss what these matrices are doing more later.

Background: S4 needs to be an RNN and CNN



S4 needs to be both an **RNN** and a **CNN**

This is elegant! But there are limitations:

- **CNN mode requires time-invariant system**
- **CNN mode cannot easily access states**
- **Complicated implementation**

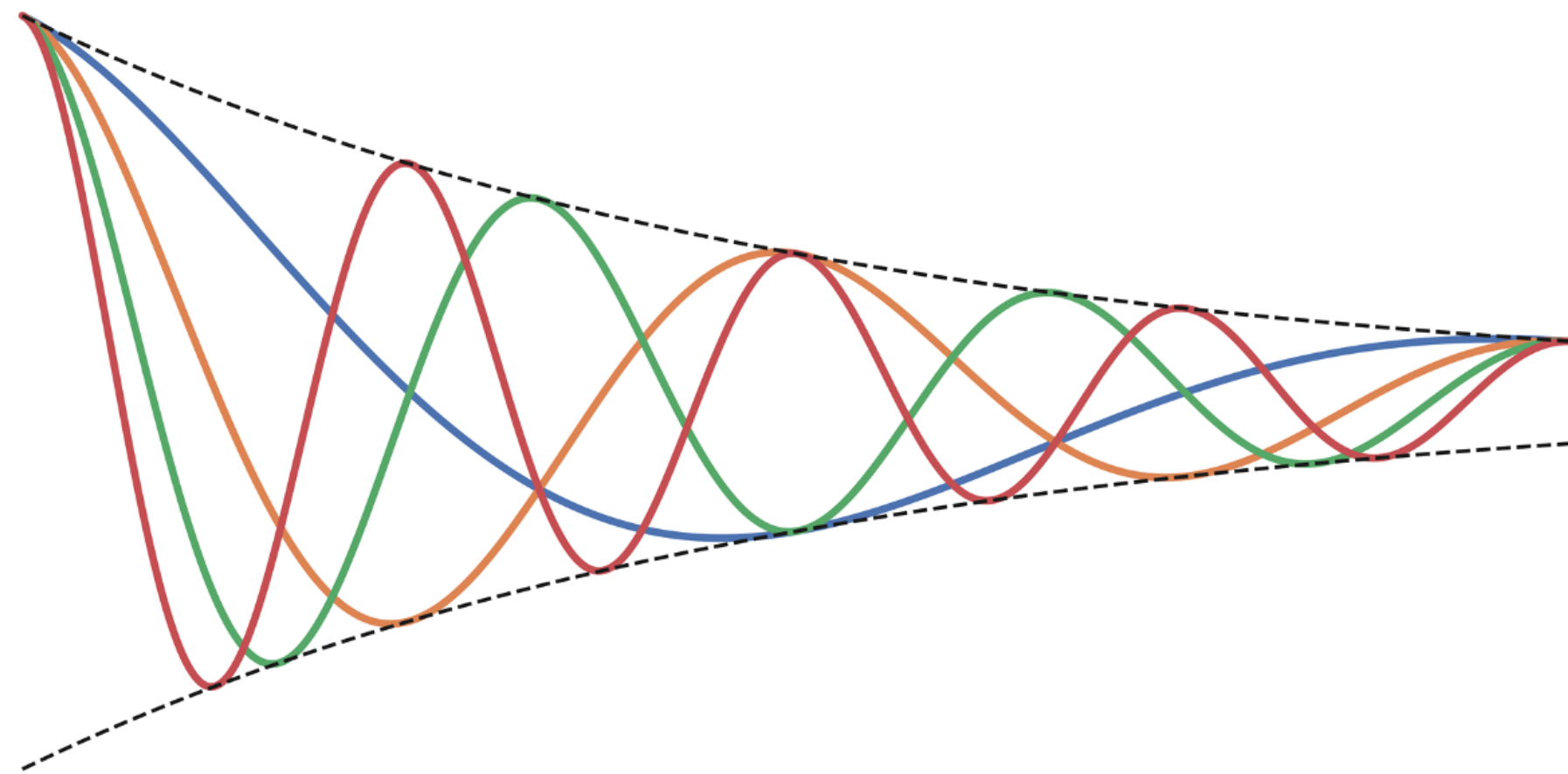
Agenda

- Introduction, motivation, prior approaches
- Linear state space models (SSMs) overview
- S4, convolutions, parameterization
- **S5, diagonalization, parallel scans**
- S6/Mamba/Hawk/Griffin, data-dependent dynamics
- Challenges and opportunities

**Can we get the same
parallelizability, efficiency and performance,
as S4 while addressing these limitations?**

From S4 to S5: Fully Recurrent

Convolution

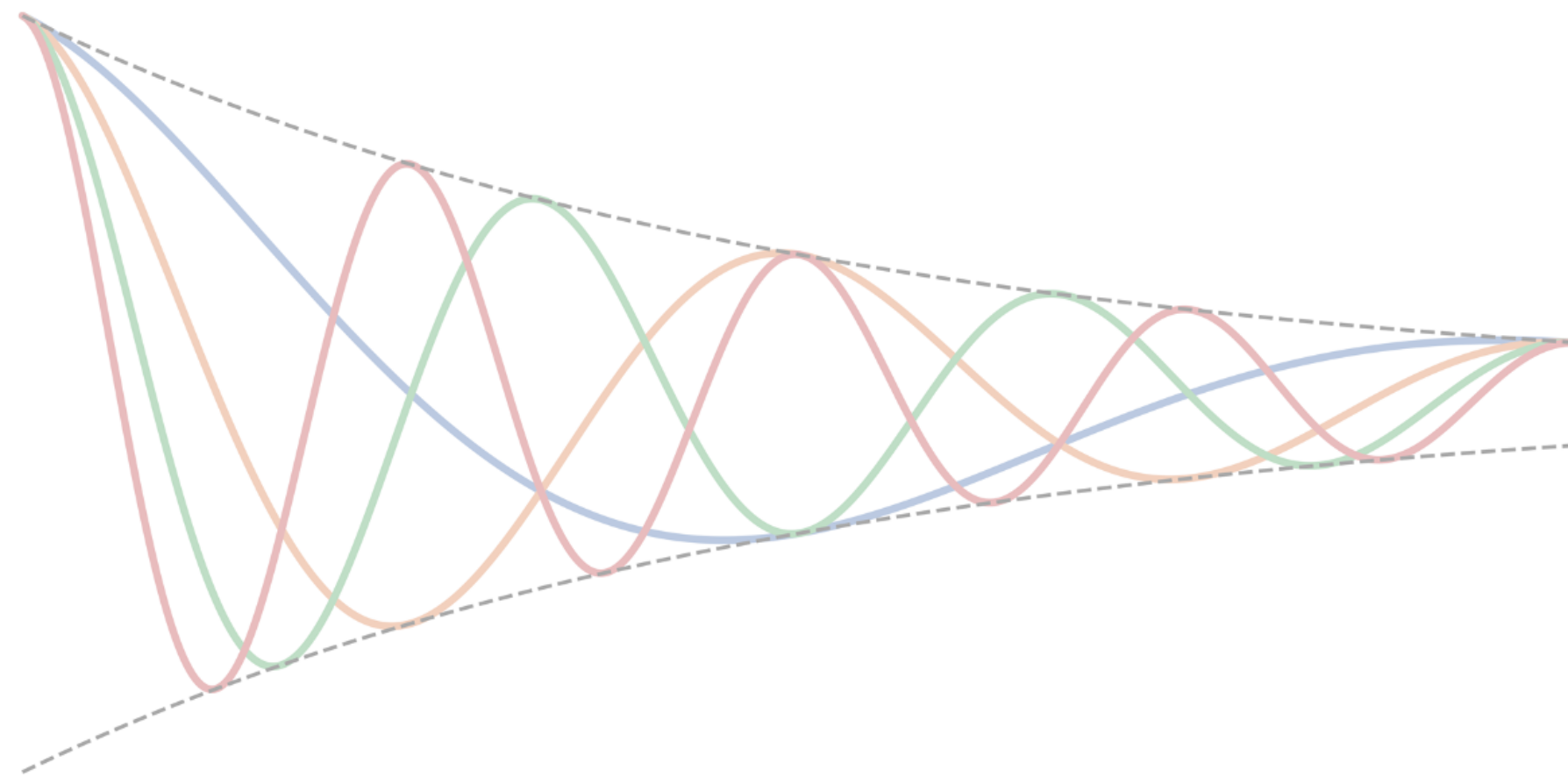


Convolution limitations:

- Requires time-invariant system
- Cannot easily access states

From S4 to S5: Fully Recurrent

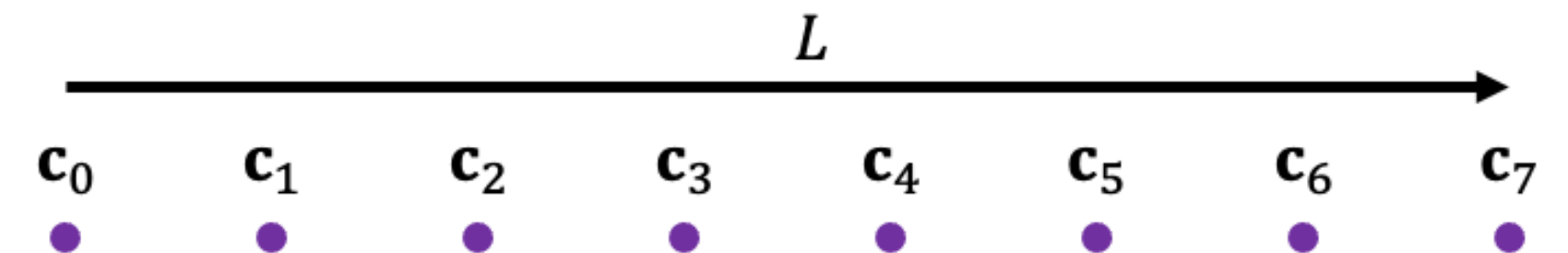
Convolution



Convolution limitations:

- Requires time-invariant system
- Cannot easily access states

Parallel scan (prefix-sum)



$\mathcal{O}(\log L)$

Scan allows:

- Time-varying systems
- Access to states (parallel or autoregressive)

From S4 to S5:

From S4 to S5: SIS0 to MIMO

Independent
single-input, single-output (SISO)
sequence maps

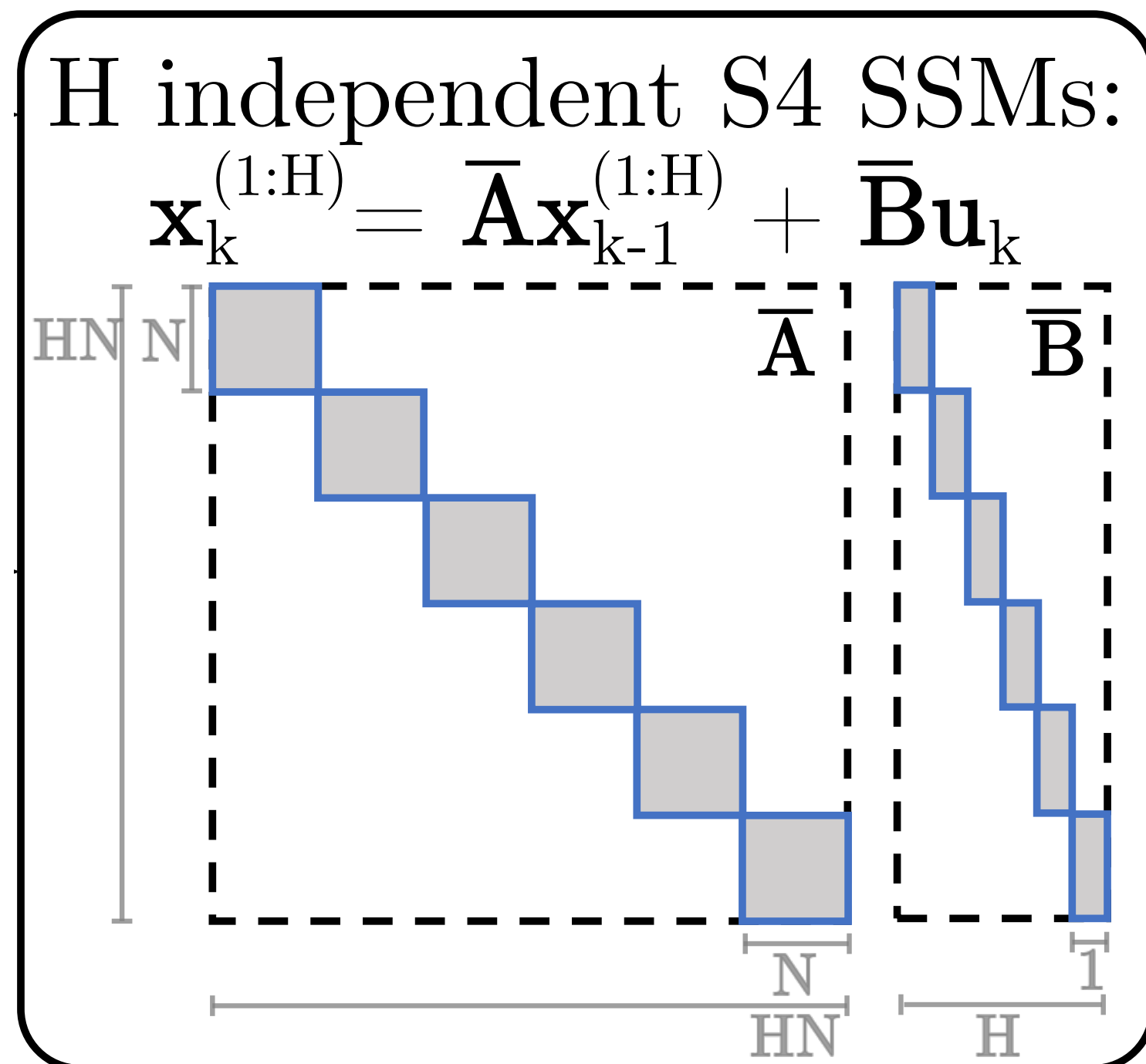
Diagram illustrating independent SISO sequence maps for three channels (1, 2, 3):

$$\begin{aligned} \begin{bmatrix} x_k^1 \\ x_k^2 \\ x_k^3 \end{bmatrix} &= \begin{bmatrix} \bar{A}_1 \\ \bar{A}_2 \\ \bar{A}_3 \end{bmatrix} \begin{bmatrix} x_{k-1}^1 \\ x_{k-1}^2 \\ x_{k-1}^3 \end{bmatrix} + \begin{bmatrix} \bar{B}_1 \\ \bar{B}_2 \\ \bar{B}_3 \end{bmatrix} * \begin{bmatrix} u_k^1 \\ u_k^2 \\ u_k^3 \end{bmatrix} \end{aligned}$$

•
•
•

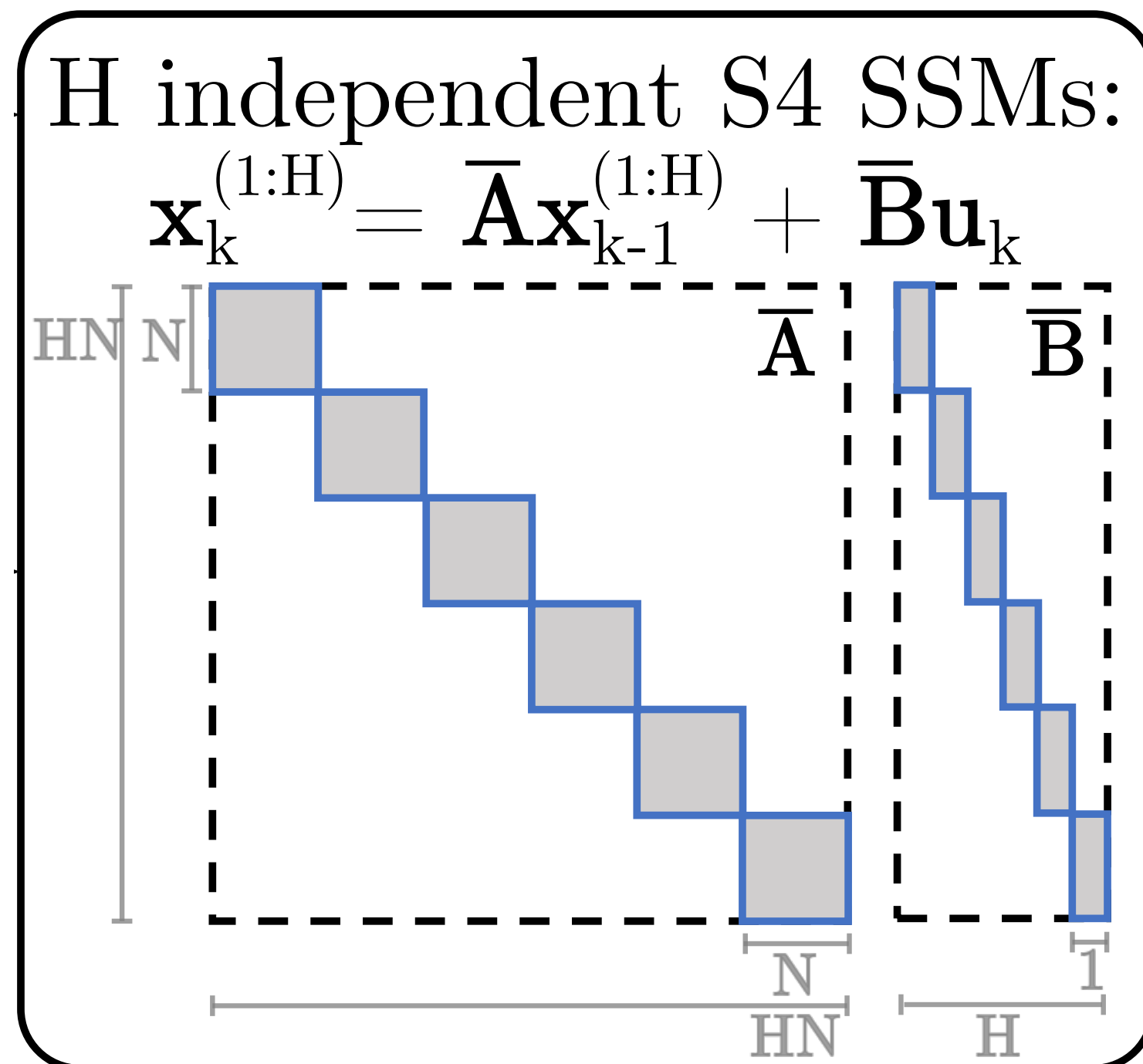
From S4 to S5: SIS0 to MIMO

**Independent
single-input, single-output (SISO)
sequence maps**



From S4 to S5: SIS0 to MIMO

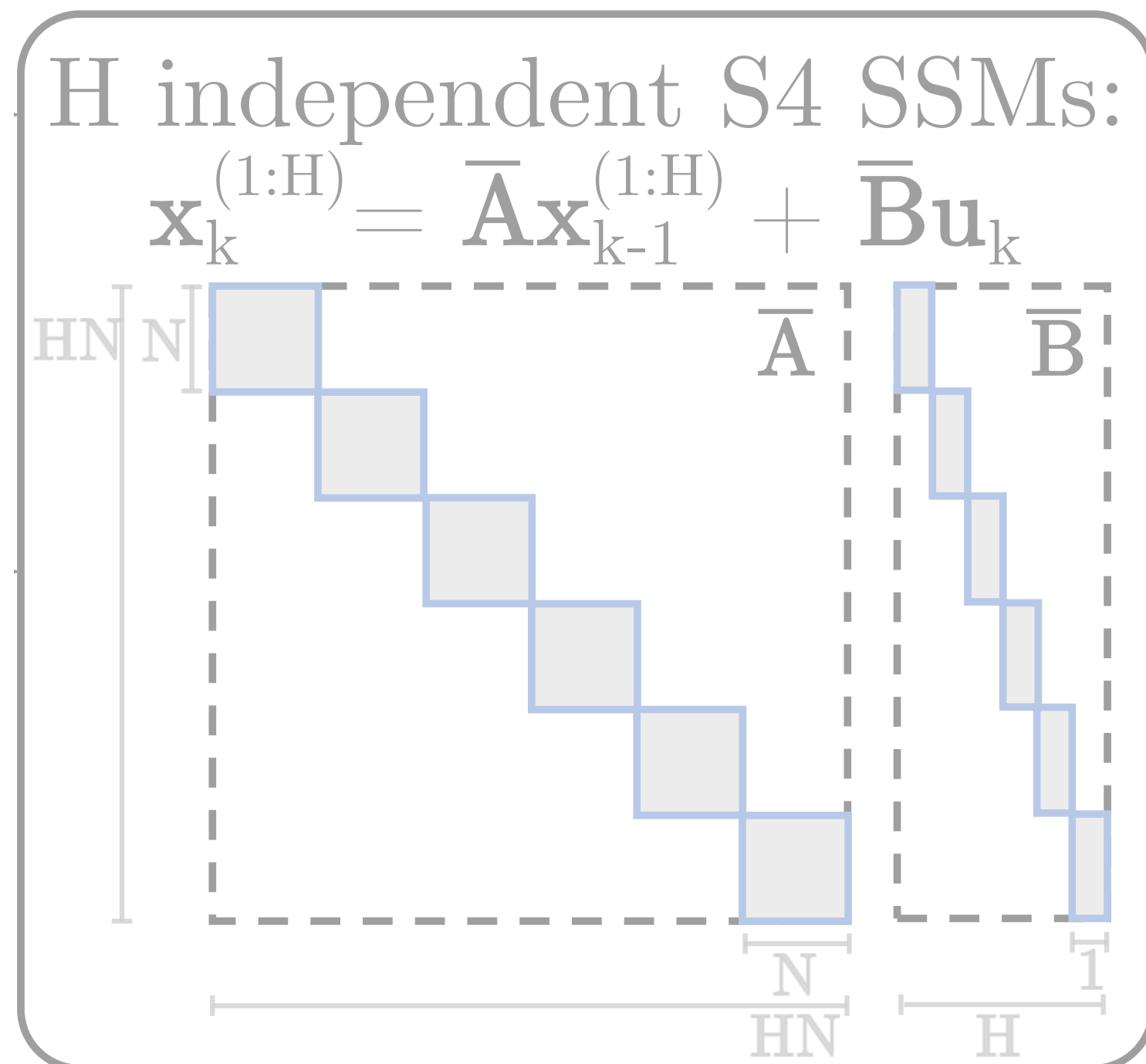
**Independent
single-input, single-output (SISO)
sequence maps**



**Large effective state size
prevents the use of (basic) parallel scans.**

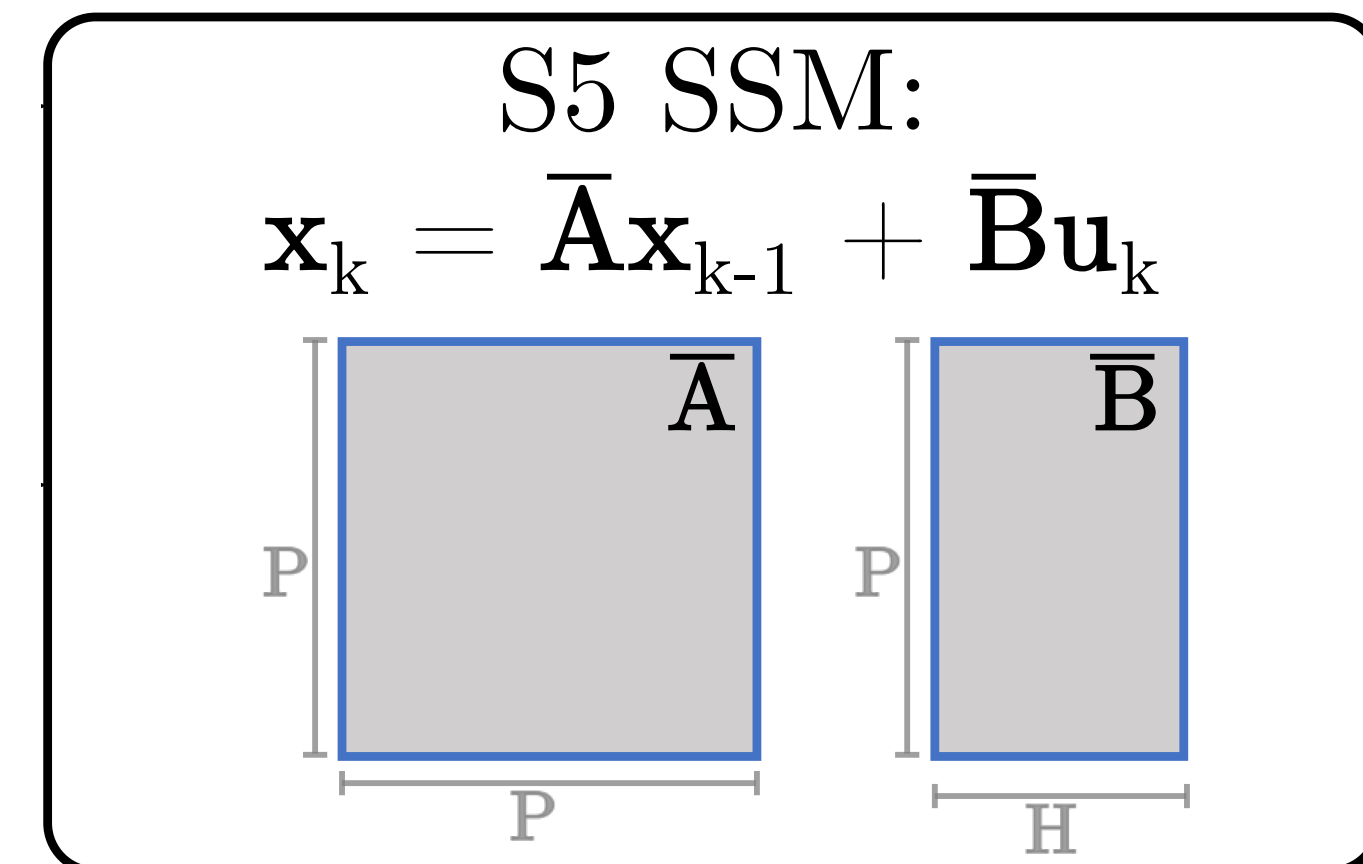
From S4 to S5: SIS0 to MIMO

Independent
single-input, single-output (SISO)
sequence maps



Large effective state size
prevents the use of (basic) parallel scans.

One
multi-input, multi-output (MIMO)
sequence map



From S4 to S5: SISO to MIMO

Independent
single-input, single-output (SISO)

One
multi-input, multi-output (MIMO)

Proposition 2. (Informal) *The output of a MIMO S5 SSM is a different projection of the same underlying dynamics of an S4 system.*

Proof. See Appendix D.2.

Implication: S5 can leverage parameterization schemes from S4

Large effective state size
prevents the use of parallel scans.

From S4 to S5: Diagonalized dynamics

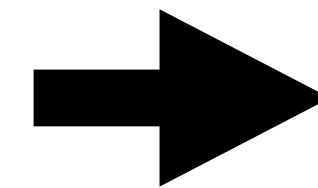
Diagonal plus low-rank
state matrix

$$A = \begin{bmatrix} \square & & & & \\ & \square & & & \\ & & \square & & \\ & & & \square & \\ & & & & \square \end{bmatrix} + \begin{bmatrix} \square \\ \square \\ \square \\ \square \\ \square \end{bmatrix} \begin{bmatrix} \square & \square & \square & \square & \square \end{bmatrix}$$

From S4 to S5: Diagonalized dynamics

Diagonal plus low-rank
state matrix

$$A = \begin{bmatrix} \square & & & & \\ & \square & & & \\ & & \square & & \\ & & & \square & \\ & & & & \square \end{bmatrix} + \begin{bmatrix} \square \\ \square \\ \square \\ \square \\ \square \end{bmatrix} \begin{bmatrix} \square & \square & \square & \square & \square \end{bmatrix}$$



Diagonal
state matrix

$$A = \begin{bmatrix} \square & & & & \\ & \square & & & \\ & & \square & & \\ & & & \square & \\ & & & & \square \end{bmatrix}$$

Similar to findings from DSS (Gupta et al. 2022)
and S4D (Gu et al. 2022)

Importance of parameterization (e.g. HiPPO)?

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

Importance of parameterization (e.g. HiPPO)?

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

diagonal, complex

ZOH Discretization:

$$\mathbf{x}_k = e^{\mathbf{A}\Delta}\mathbf{x}_{k-1} + \mathbf{A}^{-1}(\overline{\mathbf{A}} - \mathbf{I})\mathbf{B}\mathbf{u}_k$$

Importance of parameterization (e.g. HiPPO)?

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

diagonal, complex

ZOH Discretization:

$$\mathbf{x}_k = e^{\mathbf{A}\Delta}\mathbf{x}_{k-1} + \mathbf{A}^{-1}(\overline{\mathbf{A}} - \mathbf{I})\mathbf{B}\mathbf{u}_k$$

Note: real-valued diagonal matrices would be restricted in expressivity in terms of which dynamics can be represented.

But almost all square matrices are diagonalizable over the complex plane: <https://chiasme.wordpress.com/2013/09/03/almost-all-matrices-are-diagonalizable/>

Importance of parameterization (e.g. HiPPO)?

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

diagonal, complex

ZOH Discretization:

$$\mathbf{x}_k = e^{\mathbf{A}\Delta}\mathbf{x}_{k-1} + \mathbf{A}^{-1}(\overline{\mathbf{A}} - \mathbf{I})\mathbf{B}\mathbf{u}_k$$

Stability criteria:

- To avoid exploding, discrete eigenvalues should be within the complex unit circle

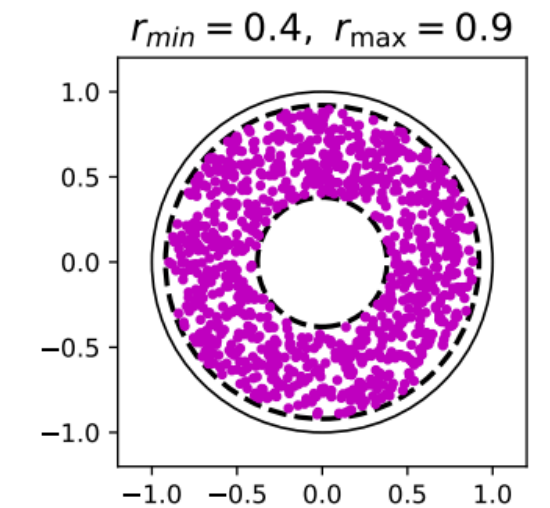


Figure 3 | Eigenvalues of a diagonal matrix A with entries sampled using Lemma 3.2. For $r_{\min} = 0$, $r_{\max} = 1$, the distribution coincides with Glorot init. in the limit.

Importance of parameterization (e.g. HiPPO)?

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

diagonal, complex

ZOH Discretization:

$$\mathbf{x}_k = e^{\mathbf{A}\Delta}\mathbf{x}_{k-1} + \mathbf{A}^{-1}(\overline{\mathbf{A}} - \mathbf{I})\mathbf{B}\mathbf{u}_k$$

Stability criteria:

- To avoid exploding, discrete eigenvalues should be within the complex unit circle
- To avoid vanishing too quickly, eigenvalues should be close to complex unit circle (same for vanishing gradients due to linearity).

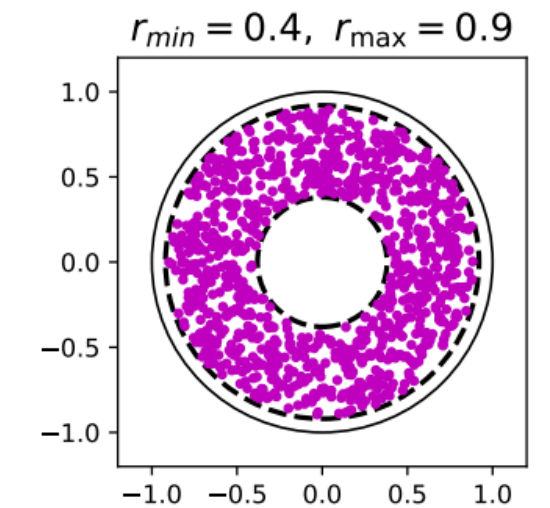


Figure 3 | Eigenvalues of a diagonal matrix A with entries sampled using Lemma 3.2. For $r_{\min} = 0$, $r_{\max} = 1$, the distribution coincides with Glorot init. in the limit.

Importance of parameterization (e.g. HiPPO)?

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

diagonal, complex

ZOH Discretization:

$$\mathbf{x}_k = e^{\mathbf{A}\Delta}\mathbf{x}_{k-1} + \mathbf{A}^{-1}(\overline{\mathbf{A}} - \mathbf{I})\mathbf{B}\mathbf{u}_k$$

Stability criteria:

- To avoid exploding, discrete eigenvalues should be within the complex unit circle
- To avoid vanishing too quickly, eigenvalues should be close to complex unit circle (same for vanishing gradients due to linearity).

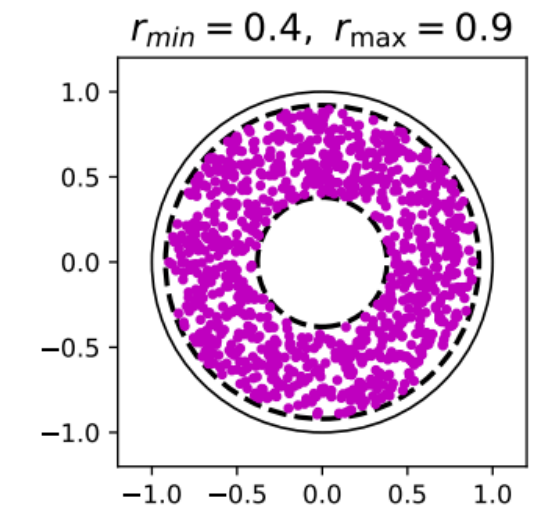
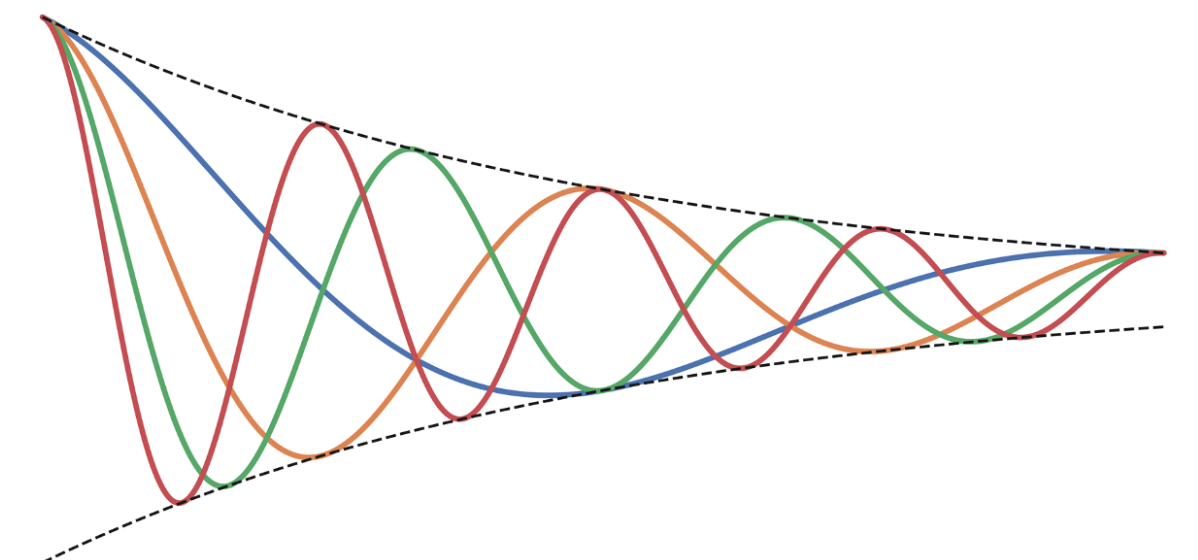


Figure 3 | Eigenvalues of a diagonal matrix A with entries sampled using Lemma 3.2. For $r_{\min} = 0$, $r_{\max} = 1$, the distribution coincides with Glorot init. in the limit.

HiPPO initialization gives these nice properties with stable, slowly decaying eigenvalues



Importance of parameterization (e.g. HiPPO)?

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

diagonal, complex

ZOH Discretization:

$$\mathbf{x}_k = e^{\mathbf{A}\Delta}\mathbf{x}_{k-1} + \mathbf{A}^{-1}(\overline{\mathbf{A}} - \mathbf{I})\mathbf{B}\mathbf{u}_k$$

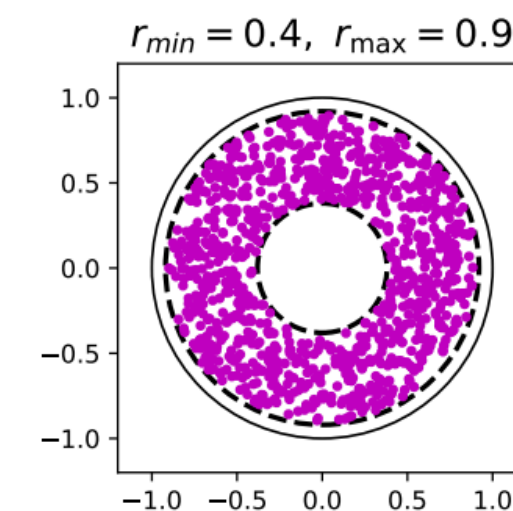
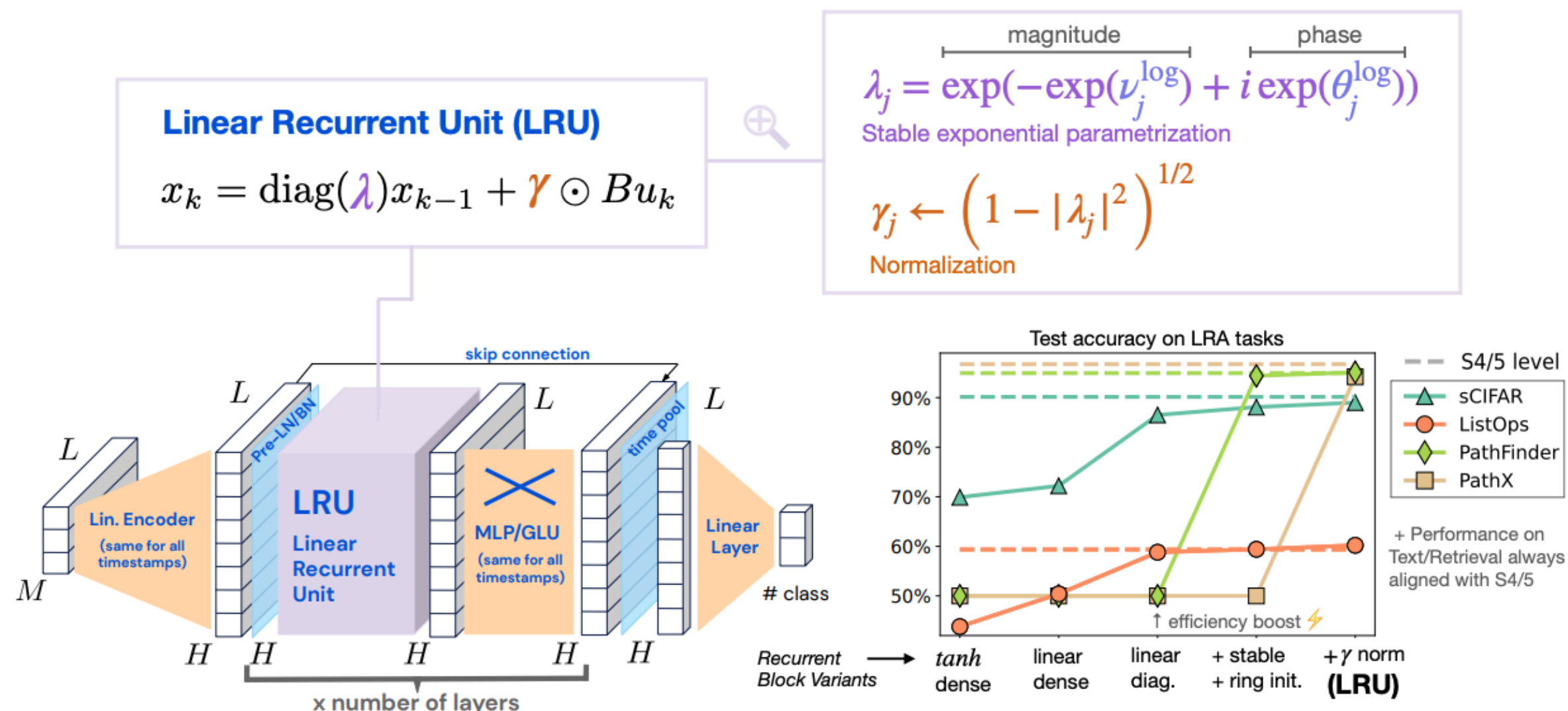


Figure 3 | Eigenvalues of a diagonal matrix \mathbf{A} with entries sampled using Lemma 3.2. For $r_{\min} = 0$, $r_{\max} = 1$, the distribution coincides with Glorot init. in the limit.

LRU paper (Orvieto 2023), shows S4/S5 style linear RNNs can be parameterized without explicit discretization or HiPPO, but still achieve similar performance on benchmarks



Importance of parameterization (e.g. HiPPO)?

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

diagonal, complex

ZOH Discretization:

$$\mathbf{x}_k = e^{\overline{\mathbf{A}}\Delta} \mathbf{x}_{k-1} + \overline{\mathbf{A}}^{-1}(\overline{\mathbf{A}} - \mathbf{I})\overline{\mathbf{B}}\mathbf{u}_k$$

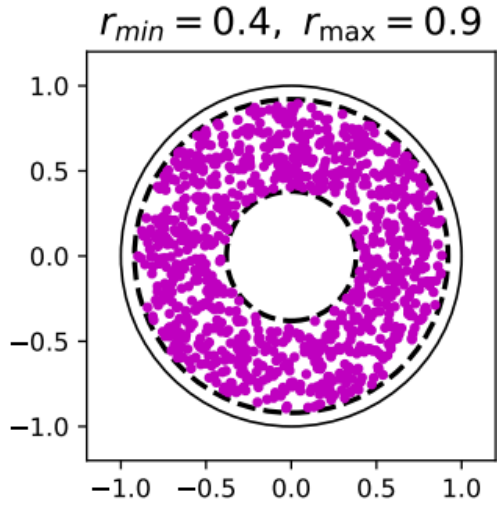
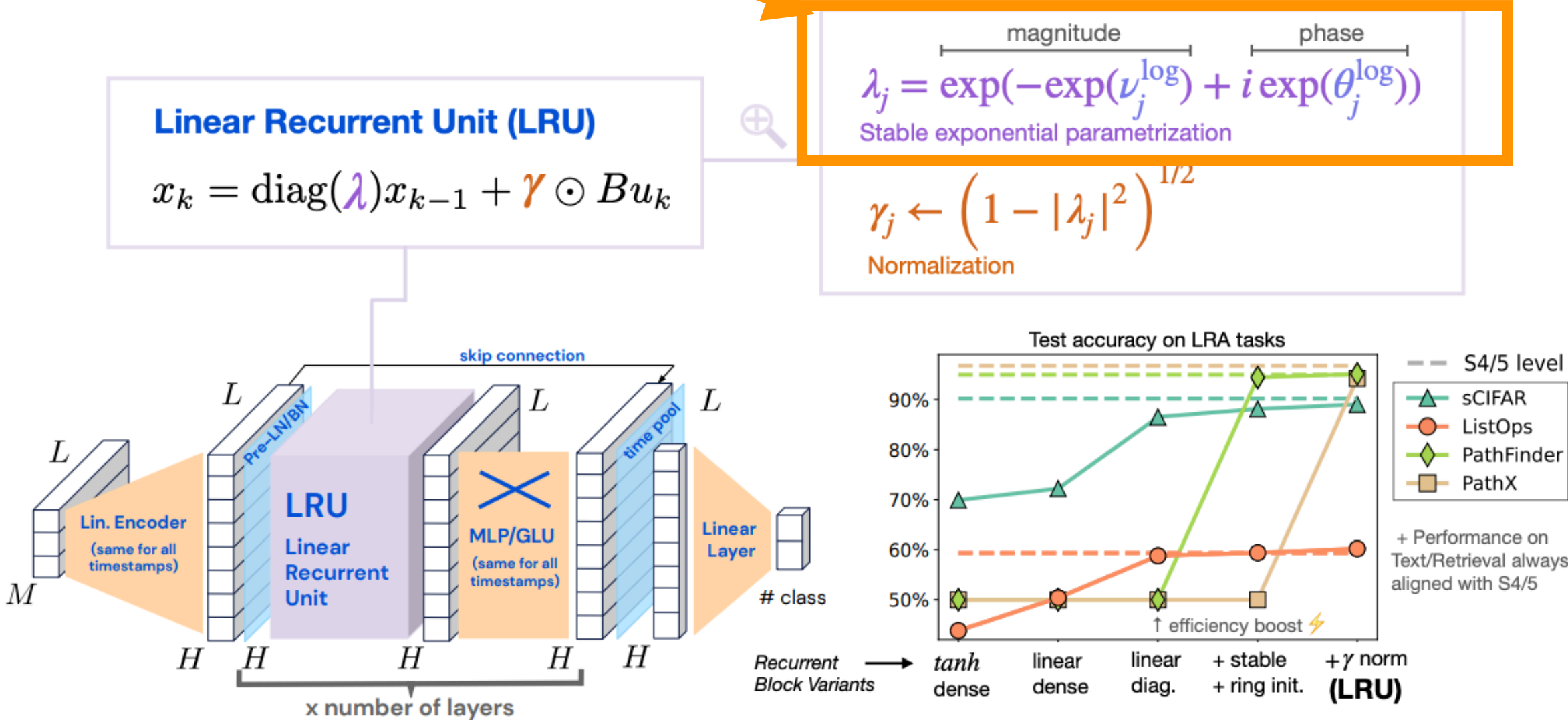


Figure 3 | Eigenvalues of a diagonal matrix A with entries sampled using Lemma 3.2. For $r_{\min} = 0$, $r_{\max} = 1$, the distribution coincides with Glorot init. in the limit.

LRU paper (Orvieto 2023), shows S4/S5 style linear RNNs can be parameterized without explicit discretization or HiPPO, but still achieve similar performance on benchmarks

Takeaways:

- Complex parameterization important for some problems



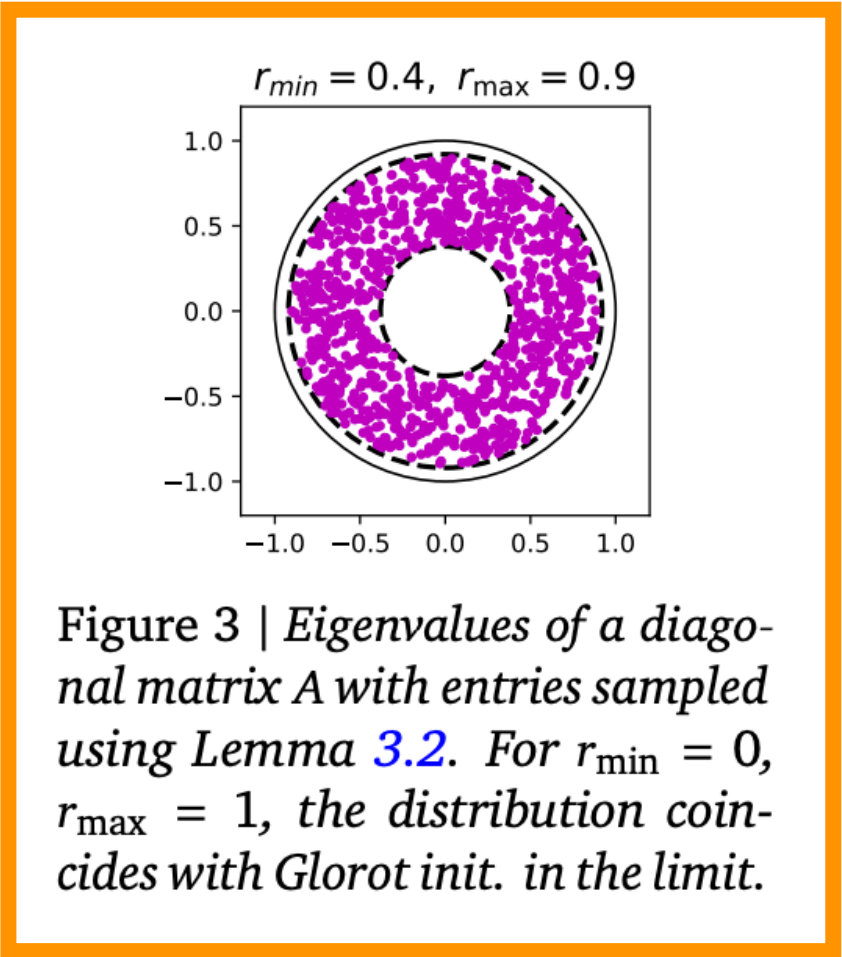
Importance of parameterization (e.g. HiPPO)?

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

diagonal, complex

ZOH Discretization:

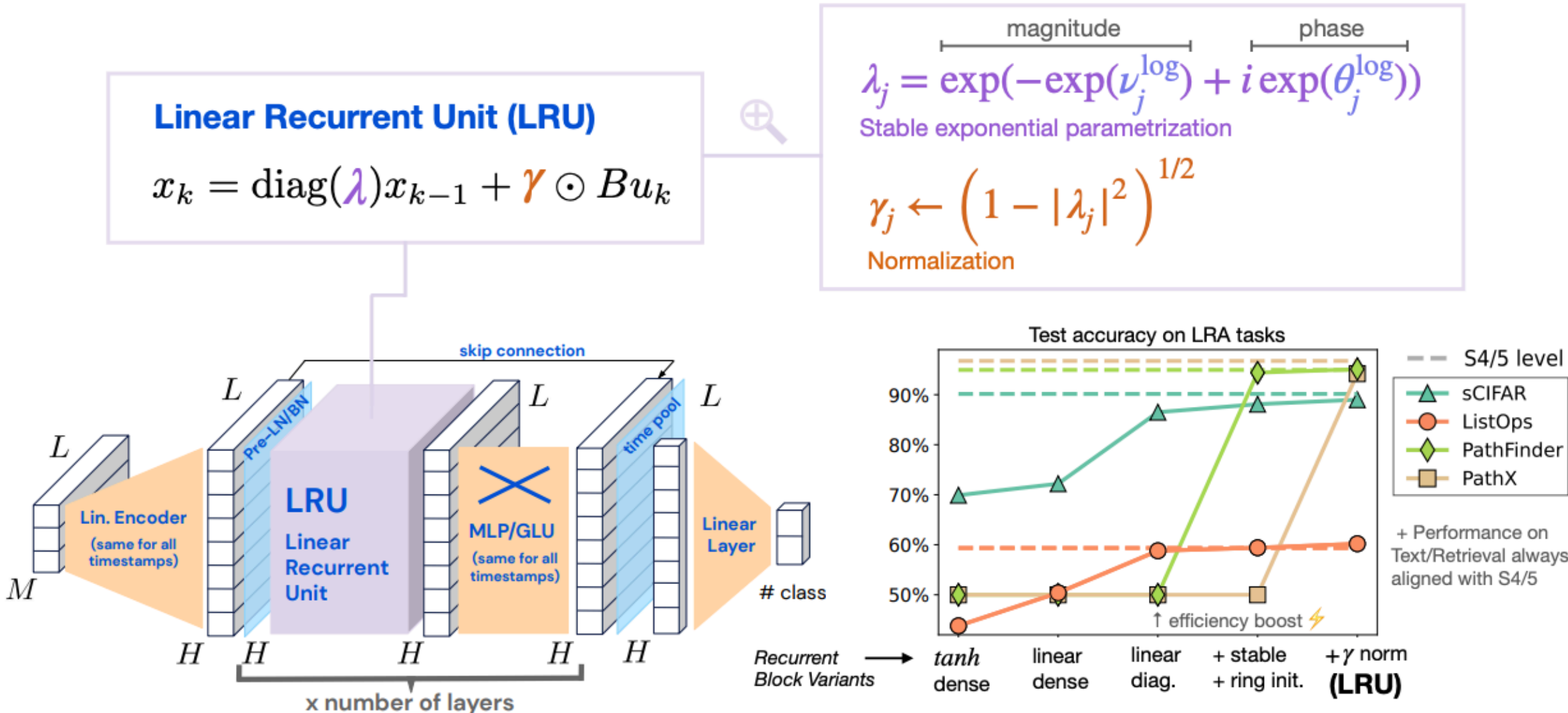
$$\mathbf{x}_k = e^{\mathbf{A}\Delta}\mathbf{x}_{k-1} + \mathbf{A}^{-1}(\overline{\mathbf{A}} - \mathbf{I})\overline{\mathbf{B}}\mathbf{u}_k$$



LRU paper (Orvieto 2023), shows S4/S5 style linear RNNs can be parameterized without explicit discretization or HiPPO, but still achieve similar performance on benchmarks

Takeaways:

- Complex parameterization important for some problems
- HiPPO+discretization gives a smart eigenvalue initialization



Importance of parameterization (e.g. HiPPO)?

$$\mathbf{x}_k = \overline{\mathbf{A}}\mathbf{x}_{k-1} + \overline{\mathbf{B}}\mathbf{u}_k$$

diagonal, complex

ZOH Discretization:

$$\mathbf{x}_k = e^{\mathbf{A}\Delta}\mathbf{x}_{k-1} + \mathbf{A}^{-1}(\overline{\mathbf{A}} - \mathbf{I})\overline{\mathbf{B}}\mathbf{u}_k$$

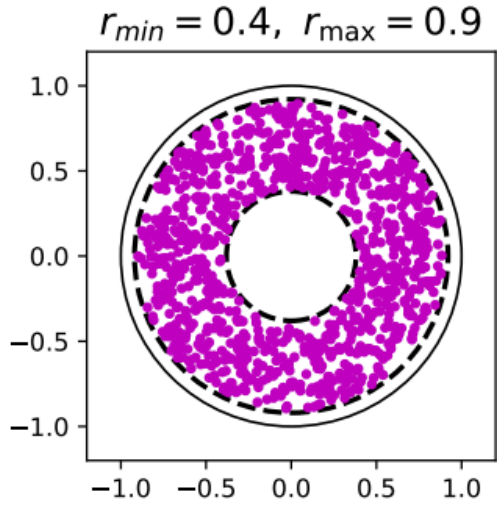
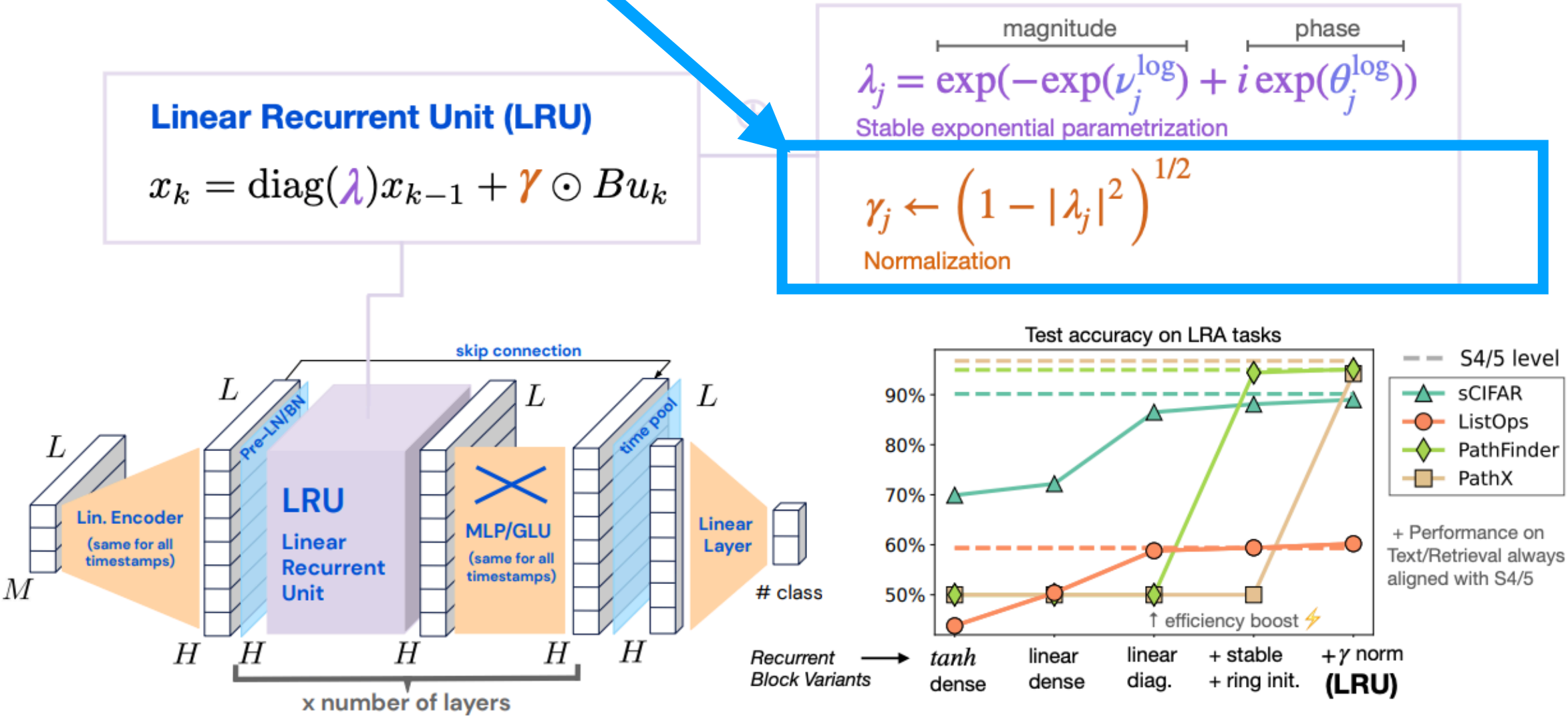


Figure 3 | Eigenvalues of a diagonal matrix \mathbf{A} with entries sampled using Lemma 3.2. For $r_{\min} = 0$, $r_{\max} = 1$, the distribution coincides with Glorot init. in the limit.

LRU paper (Orvieto 2023), shows S4/S5 style linear RNNs can be parameterized without explicit discretization or HiPPO, but still achieve similar performance on benchmarks

Takeaways:

- Complex parameterization important for some problems
- HiPPO+discretization gives a smart eigenvalue initialization
- Discretization provides normalizing effect on effective inputs



From S4 to S5: Diagonalized dynamics

Diagonal plus low-rank
state matrix

Diagonal
state matrix

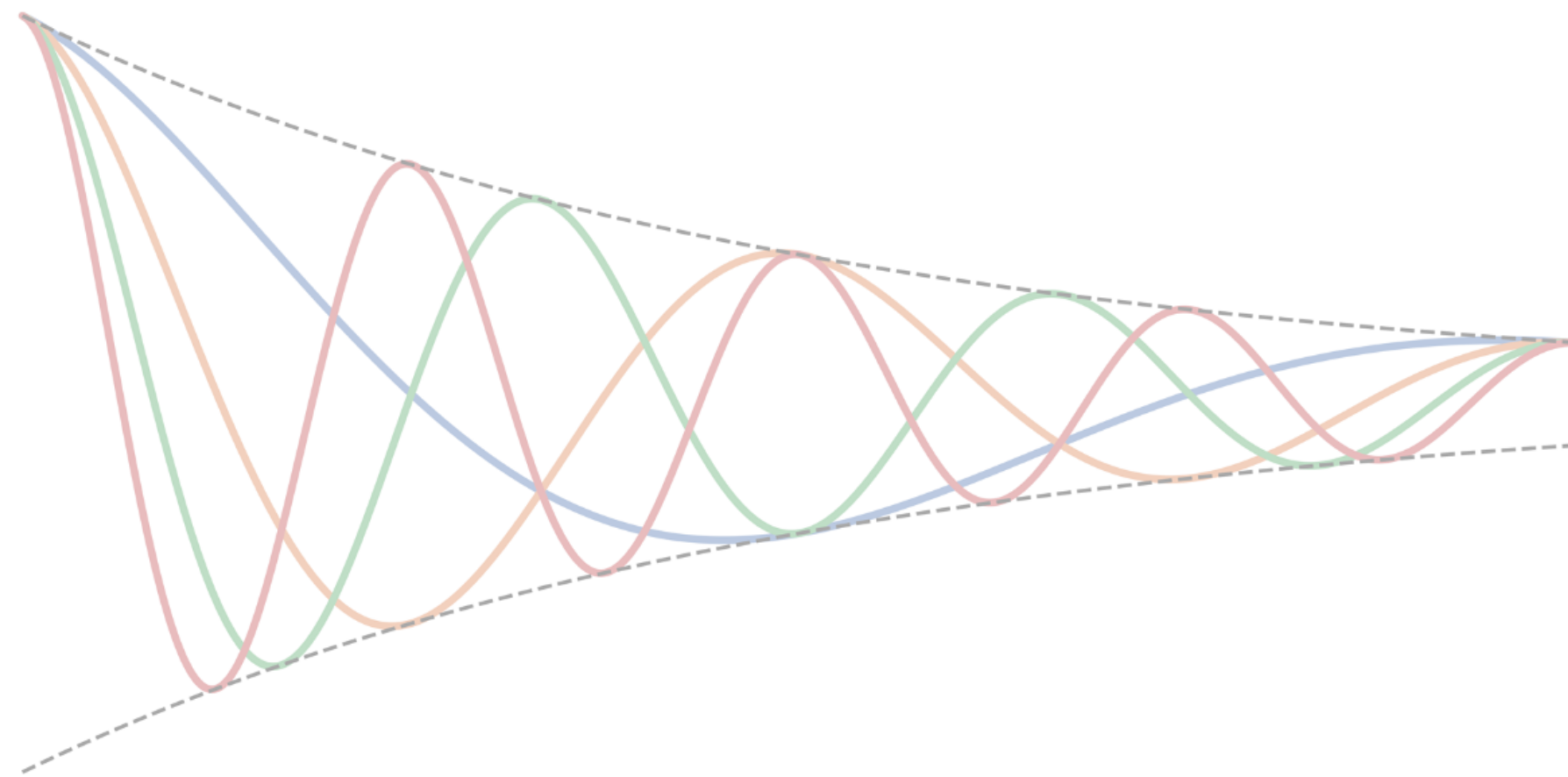
**MIMO + diagonalized dynamics:
enables the use of efficient parallel scans.**



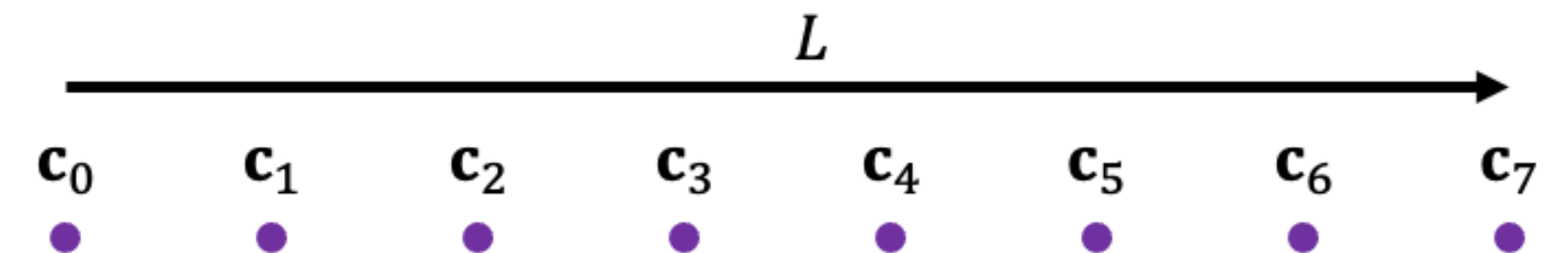
Similar findings to DSS (Gupta et al. 2022)
and S4D (Gu et al. 2022)

From S4 to S5: Fully Recurrent

Convolution



Parallel scan



$\mathcal{O}(\log L)$

Convolution limitations:

- Requires time-invariant system
- Cannot access states

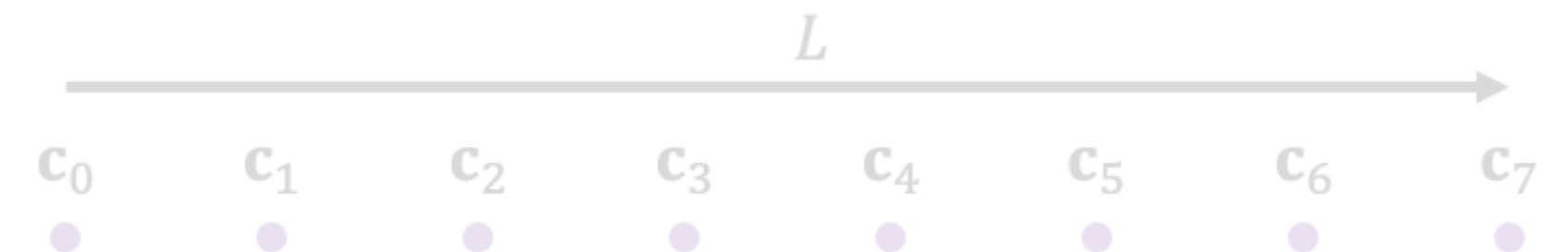
Scan allows:

- Time-varying systems
- Access to states (parallel or autoregressive)

From S4 to S5: Fully Recurrent

Convolution

Parallel scan



Proposition 1. *(Informal) An S5 layer is as efficient as an S4 layer.*

Proof. See Appendix C.1.

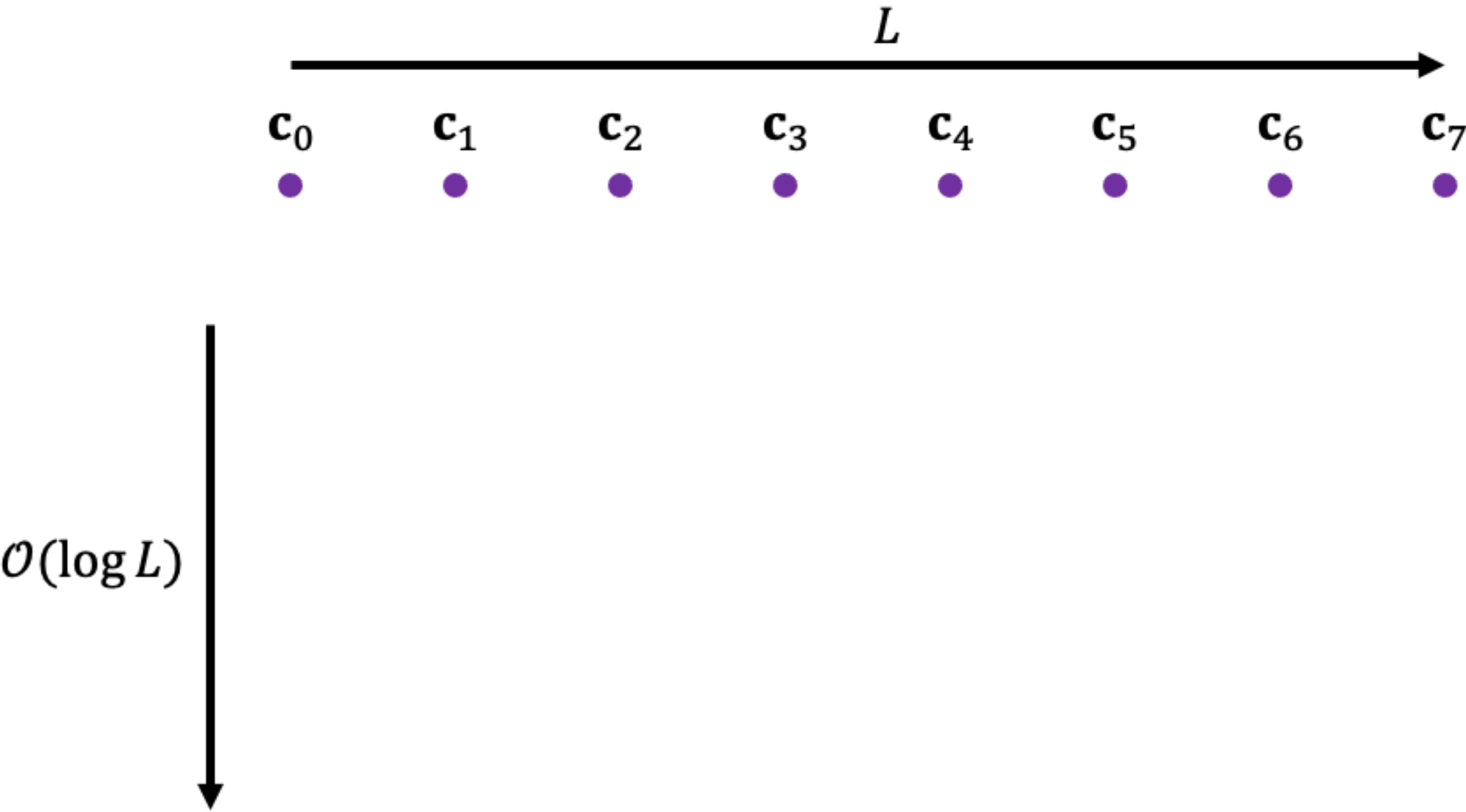
Convolution limitations:

- Requires time-invariant system
- Cannot access states

Scan allows:

- Time-varying systems
- Access to states (parallel or autoregressive)

From S4 to S5: Parallel Scans



From S4 to S5: Parallel Scans

Consider the scalar sequence: $[a, b, c, d]$ and the addition operator $+$.

Performing a scan (all prefix sum) on this sequence using $+$ returns the cumulative sum:

$[a, a+b, a+b+c, a+b+c+d]$

From S4 to S5: Parallel Scans

Consider the scalar sequence: $[a, b, c, d]$ and the addition operator $+$.

Performing a scan (all prefix sum) on this sequence using $+$ returns the cumulative sum:

$[a, a+b, a+b+c, a+b+c+d]$

Sequential Scan

(3 sequential steps required)

a

From S4 to S5: Parallel Scans

Consider the scalar sequence: $[a, b, c, d]$ and the addition operator $+$.

Performing a scan (all prefix sum) on this sequence using $+$ returns the cumulative sum:

$[a, a+b, a+b+c, a+b+c+d]$

Sequential Scan

(3 sequential steps required)

a

$a+b$

From S4 to S5: Parallel Scans

Consider the scalar sequence: $[a, b, c, d]$ and the addition operator $+$.

Performing a scan (all prefix sum) on this sequence using $+$ returns the cumulative sum:

$[a, a+b, a+b+c, a+b+c+d]$

Sequential Scan

(3 sequential steps required)

a

$a+b$

$a+b+c$

From S4 to S5: Parallel Scans

Consider the scalar sequence: $[a, b, c, d]$ and the addition operator $+$.

Performing a scan (all prefix sum) on this sequence using $+$ returns the cumulative sum:

$[a, a+b, a+b+c, a+b+c+d]$

Sequential Scan

(3 sequential steps required)

a

$a+b$

$a+b+c$

$a+b+c+d$

From S4 to S5: Parallel Scans

Consider the scalar sequence: [a,b,c,d] and the addition operator +.
Performing a scan (all prefix sum) on this sequence using + returns the cumulative sum:

[a, a+b, a+b+c, a+b+c+d]

Sequential Scan
(3 sequential steps required)

a
a+b
a+b+c
a+b+c+d

Parallel Scan
(2 sequential steps required)



From S4 to S5: Parallel Scans

Consider the scalar sequence: $[a, b, c, d]$ and the addition operator $+$.

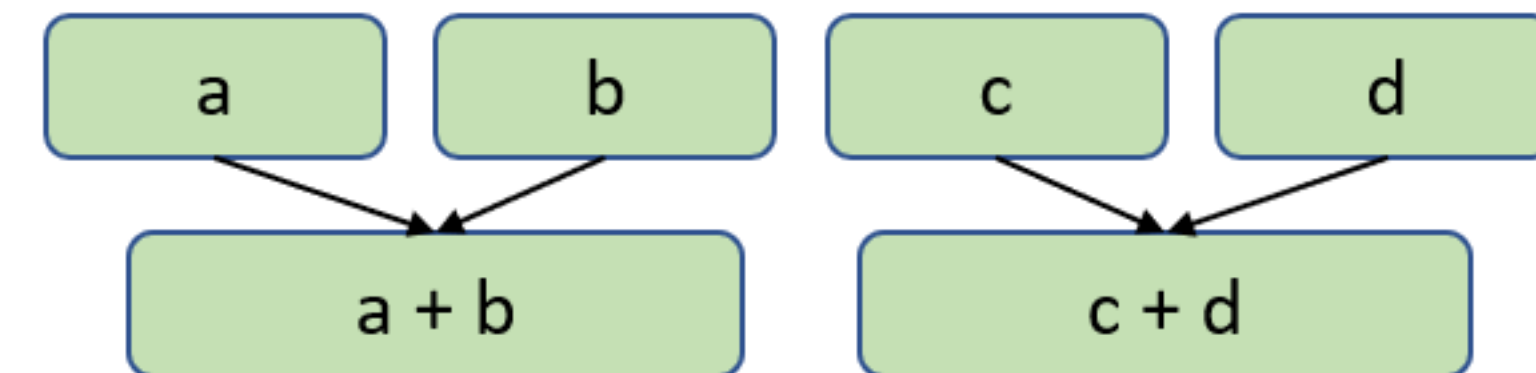
Performing a scan (all prefix sum) on this sequence using $+$ returns the cumulative sum:

$[a, a+b, a+b+c, a+b+c+d]$

Sequential Scan
(3 sequential steps required)

a
 $a+b$
 $a+b+c$
 $a+b+c+d$

Parallel Scan
(2 sequential steps required)



From S4 to S5: Parallel Scans

Consider the scalar sequence: $[a, b, c, d]$ and the addition operator $+$.

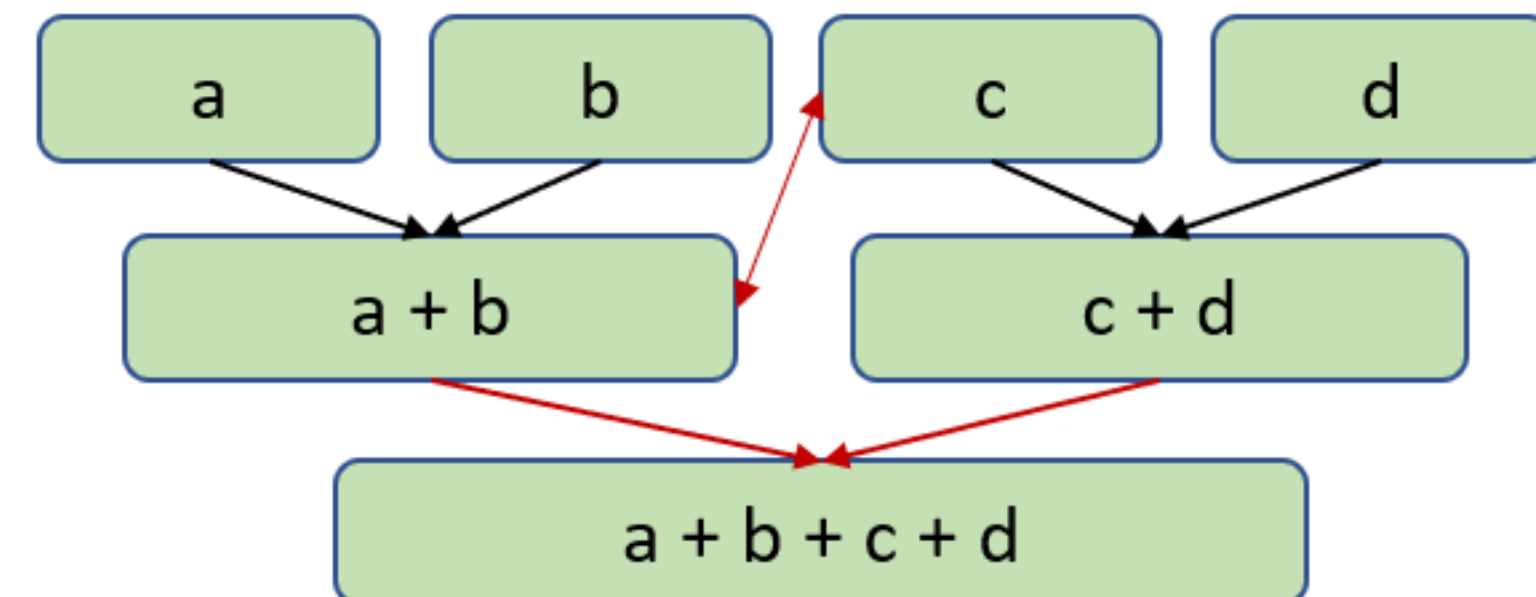
Performing a scan (all prefix sum) on this sequence using $+$ returns the cumulative sum:

$[a, a+b, a+b+c, a+b+c+d]$

Sequential Scan
(3 sequential steps required)

a
a+b
a+b+c
a+b+c+d

Parallel Scan
(2 sequential steps required)



From S4 to S5: Parallel Scans

Consider the scalar sequence: $[a, b, c, d]$ and the addition operator $+$.

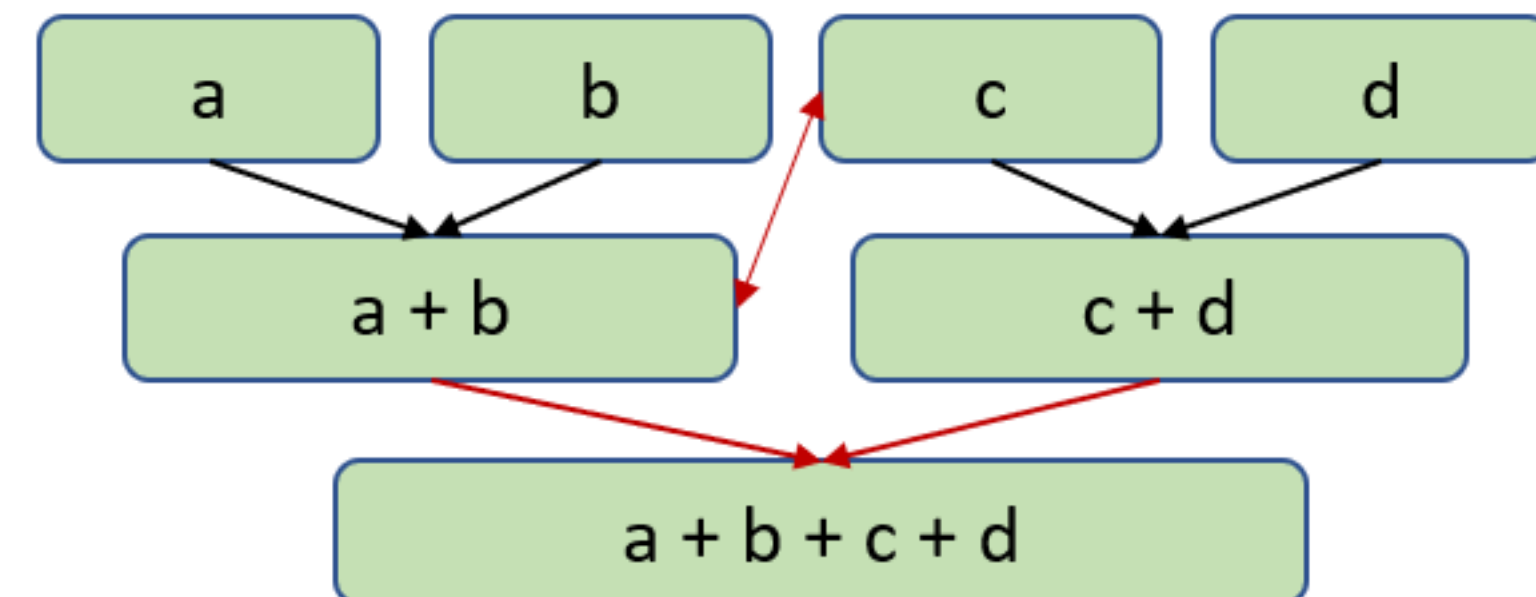
Performing a scan (all prefix sum) on this sequence using $+$ returns the cumulative sum:

$[a, a+b, a+b+c, a+b+c+d]$

Sequential Scan
(3 sequential steps required)

a
a+b
a+b+c
a+b+c+d

Parallel Scan
(2 sequential steps required)



**Given sufficient processors, number of sequential steps
scales logarithmically in sequence length**

From S4 to S5: Parallel Scans

$$x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k$$

Sequential Scan
(3 sequential steps required)

$$x_1 = \overline{\mathbf{B}}u_1$$

Parallel Scan
(2 sequential steps required)

From S4 to S5: Parallel Scans

$$x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k$$

Sequential Scan
(3 sequential steps required)

$$x_1 = \overline{\mathbf{B}}u_1$$

$$x_2 = \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2$$

Parallel Scan
(2 sequential steps required)

From S4 to S5: Parallel Scans

$$x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k$$

Sequential Scan
(3 sequential steps required)

$$x_1 = \overline{\mathbf{B}}u_1$$

$$x_2 = \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2$$

$$x_3 = \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_2 + \overline{\mathbf{B}}u_3$$

Parallel Scan
(2 sequential steps required)

From S4 to S5: Parallel Scans

$$x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k$$

Sequential Scan
(3 sequential steps required)

$$x_1 = \overline{\mathbf{B}}u_1$$

$$x_2 = \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2$$

$$x_3 = \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_2 + \overline{\mathbf{B}}u_3$$

$$x_4 = \overline{\mathbf{A}}^3\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_2 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_3 + \overline{\mathbf{B}}u_4$$

Parallel Scan
(2 sequential steps required)

From S4 to S5: Parallel Scans

$$x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k$$

Binary associative operator: $q_i \bullet q_j := (q_{j,a} \odot q_{i,a}, q_{j,a} \otimes q_{i,b} + q_{j,b})$

Sequential Scan
(3 sequential steps required)

$$x_1 = \overline{\mathbf{B}}u_1$$

$$x_2 = \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2$$

$$x_3 = \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_2 + \overline{\mathbf{B}}u_3$$

$$x_4 = \overline{\mathbf{A}}^3\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_2 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_3 + \overline{\mathbf{B}}u_4$$

Parallel Scan
(2 sequential steps required)

(A, Bu₁)

(A, Bu₂)

(A, Bu₃)

(A, Bu₄)

From S4 to S5: Parallel Scans

$$x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k$$

Binary associative operator: $q_i \bullet q_j := (q_{j,a} \odot q_{i,a}, q_{j,a} \otimes q_{i,b} + q_{j,b})$

Sequential Scan
(3 sequential steps required)

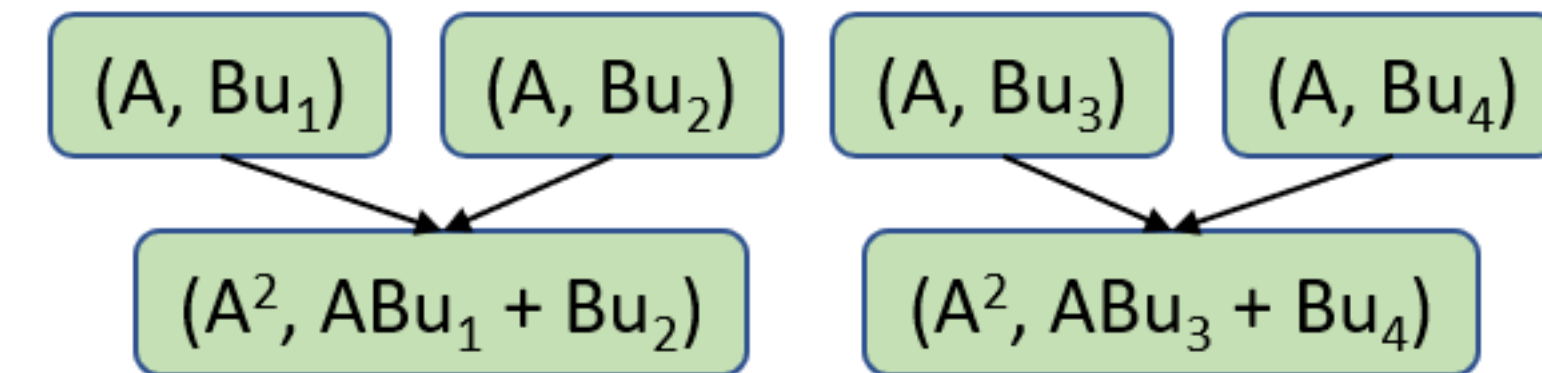
$$x_1 = \overline{\mathbf{B}}u_1$$

$$x_2 = \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2$$

$$x_3 = \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_2 + \overline{\mathbf{B}}u_3$$

$$x_4 = \overline{\mathbf{A}}^3\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_2 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_3 + \overline{\mathbf{B}}u_4$$

Parallel Scan
(2 sequential steps required)



From S4 to S5: Parallel Scans

$$x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k$$

Binary associative operator: $q_i \bullet q_j := (q_{j,a} \odot q_{i,a}, q_{j,a} \otimes q_{i,b} + q_{j,b})$

Sequential Scan
(3 sequential steps required)

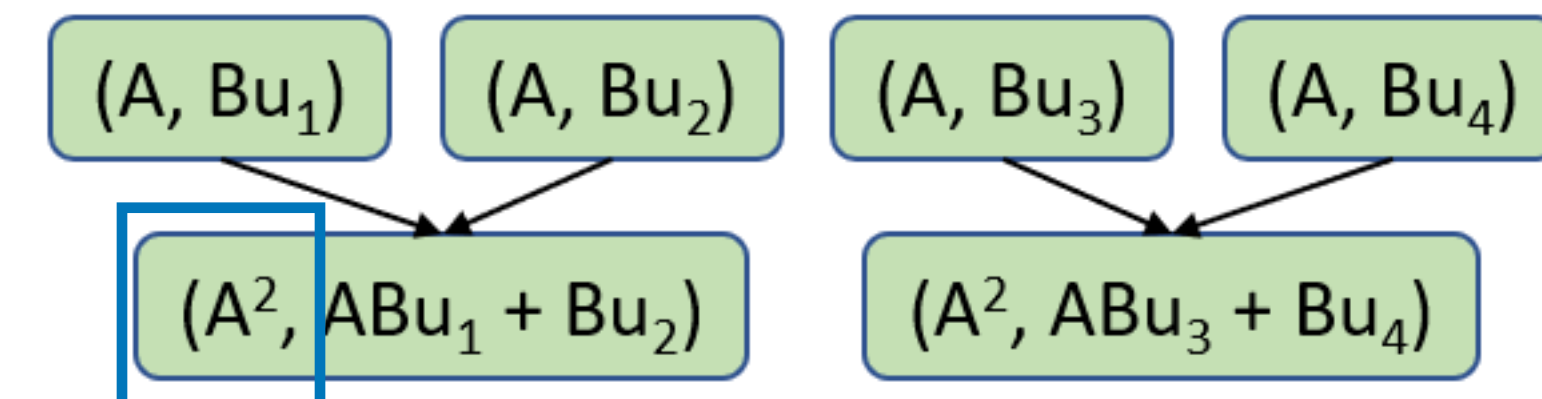
$$x_1 = \overline{\mathbf{B}}u_1$$

$$x_2 = \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2$$

$$x_3 = \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_2 + \overline{\mathbf{B}}u_3$$

$$x_4 = \overline{\mathbf{A}}^3\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_2 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_3 + \overline{\mathbf{B}}u_4$$

Parallel Scan
(2 sequential steps required)



**Note: matrix-matrix multiplication,
this is why diagonalization is important to avoid cubic cost!**

From S4 to S5: Parallel Scans

$$x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k$$

Binary associative operator: $q_i \bullet q_j := (q_{j,a} \odot q_{i,a}, q_{j,a} \otimes q_{i,b} + q_{j,b})$

Sequential Scan
(3 sequential steps required)

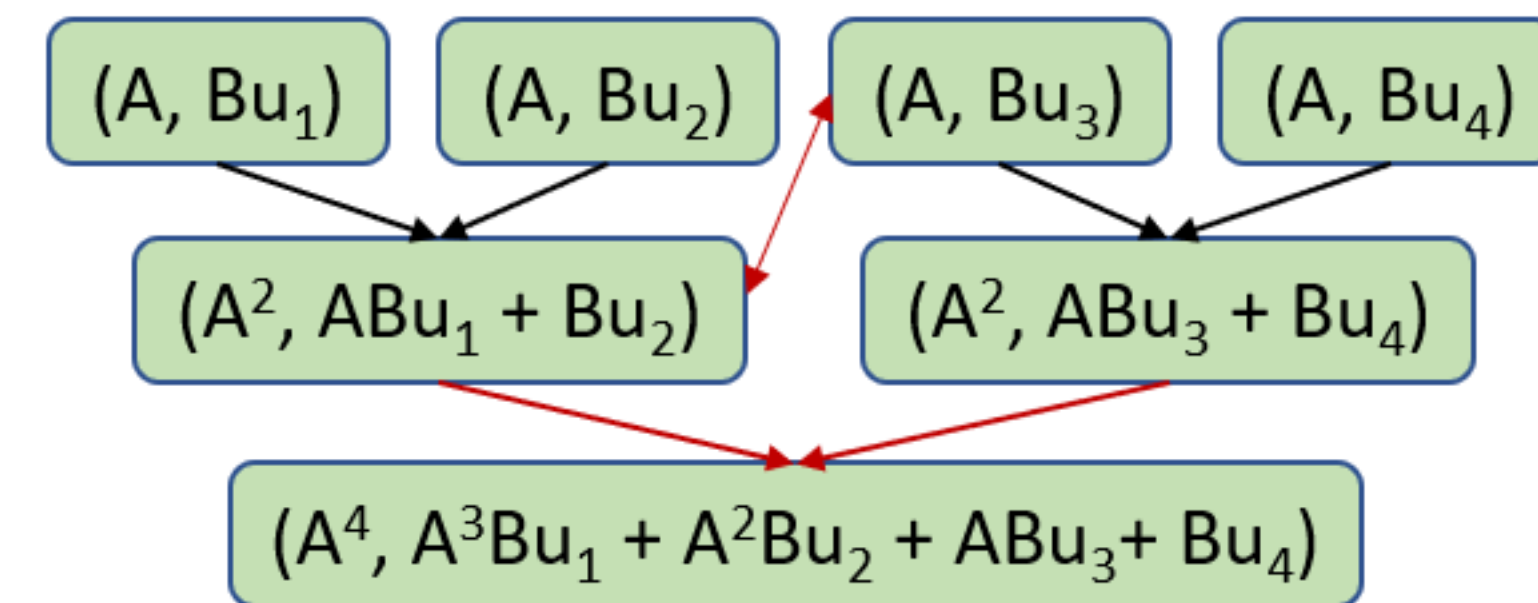
$$x_1 = \overline{\mathbf{B}}u_1$$

$$x_2 = \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2$$

$$x_3 = \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_2 + \overline{\mathbf{B}}u_3$$

$$x_4 = \overline{\mathbf{A}}^3\overline{\mathbf{B}}u_1 + \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_2 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_3 + \overline{\mathbf{B}}u_4$$

Parallel Scan
(2 sequential steps required)



Given sufficient processors, number of sequential steps
scales logarithmically in sequence length

S5 retains S4’s high performance

Long Range Arena

Model (Input length)	ListOps (2,048)	Text (4,096)	Retrieval (4,000)	Image (1,024)	Pathfinder (1,024)	Path-X (16,384)	Avg.
Transformer	36.37	64.27	57.46	42.44	71.40	✗	53.66
S4D-LegS	<u>60.47</u>	86.18	89.46	<u>88.19</u>	93.06	91.95	84.89
S4-LegS	59.60	<u>86.82</u>	<u>90.90</u>	88.65	<u>94.20</u>	<u>96.35</u>	<u>86.09</u>
S5	62.15	89.31	91.40	88.00	95.33	98.58	87.46

S5 retains S4’s high performance

Speech Commands 35-way Raw Speech Classification

Model (Input length)	Parameters	16kHz (16,000)	8kHz (8,000)
InceptionNet	481K	61.24	05.18
ResNet-1	216K	77.86	08.74
XResNet-50	904K	83.01	07.72
ConvNet	26.2M	95.51	07.26
S4-LegS	307K	<u>96.08</u>	<u>91.32</u>
S4D-LegS	306K	95.83	91.08
S5	280K	96.52	94.53

Neural Latents Benchmark

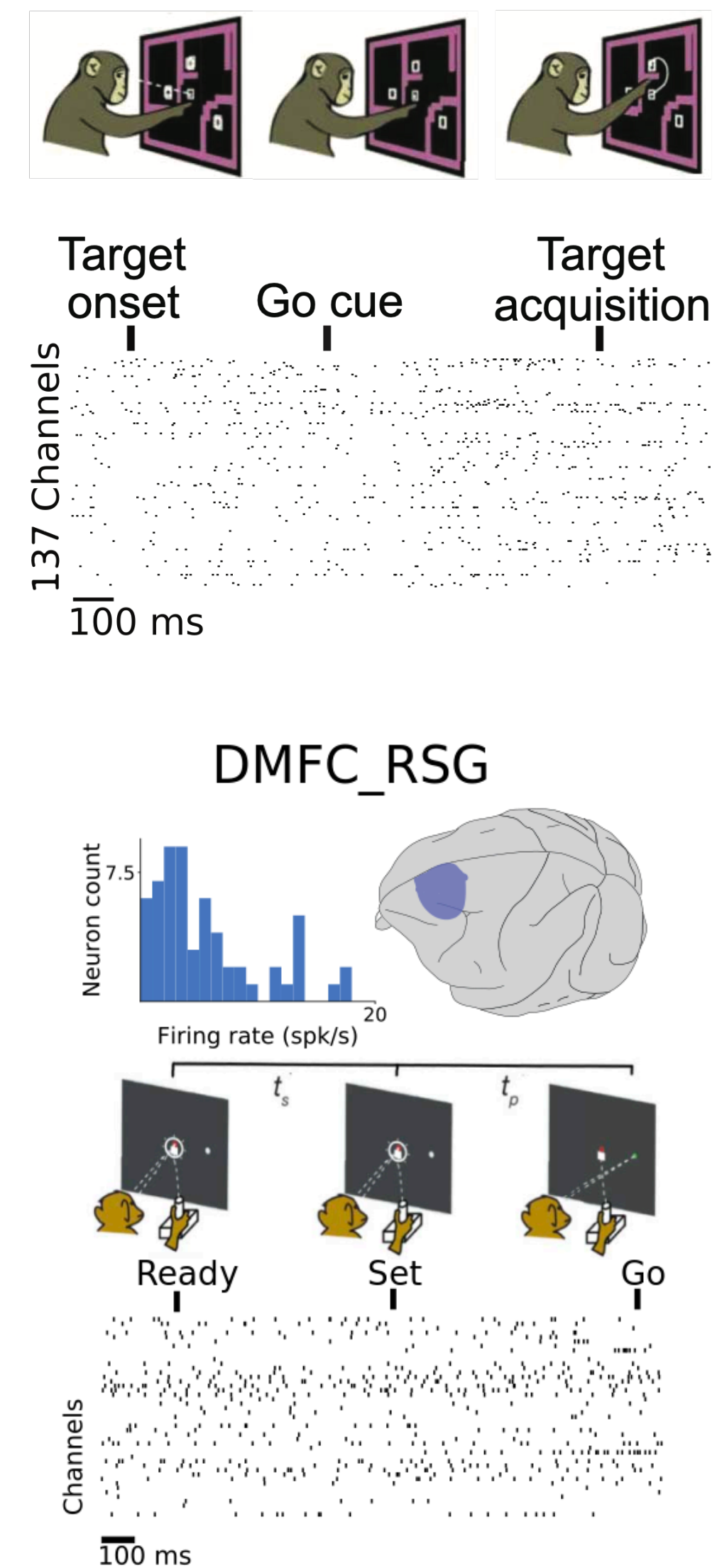
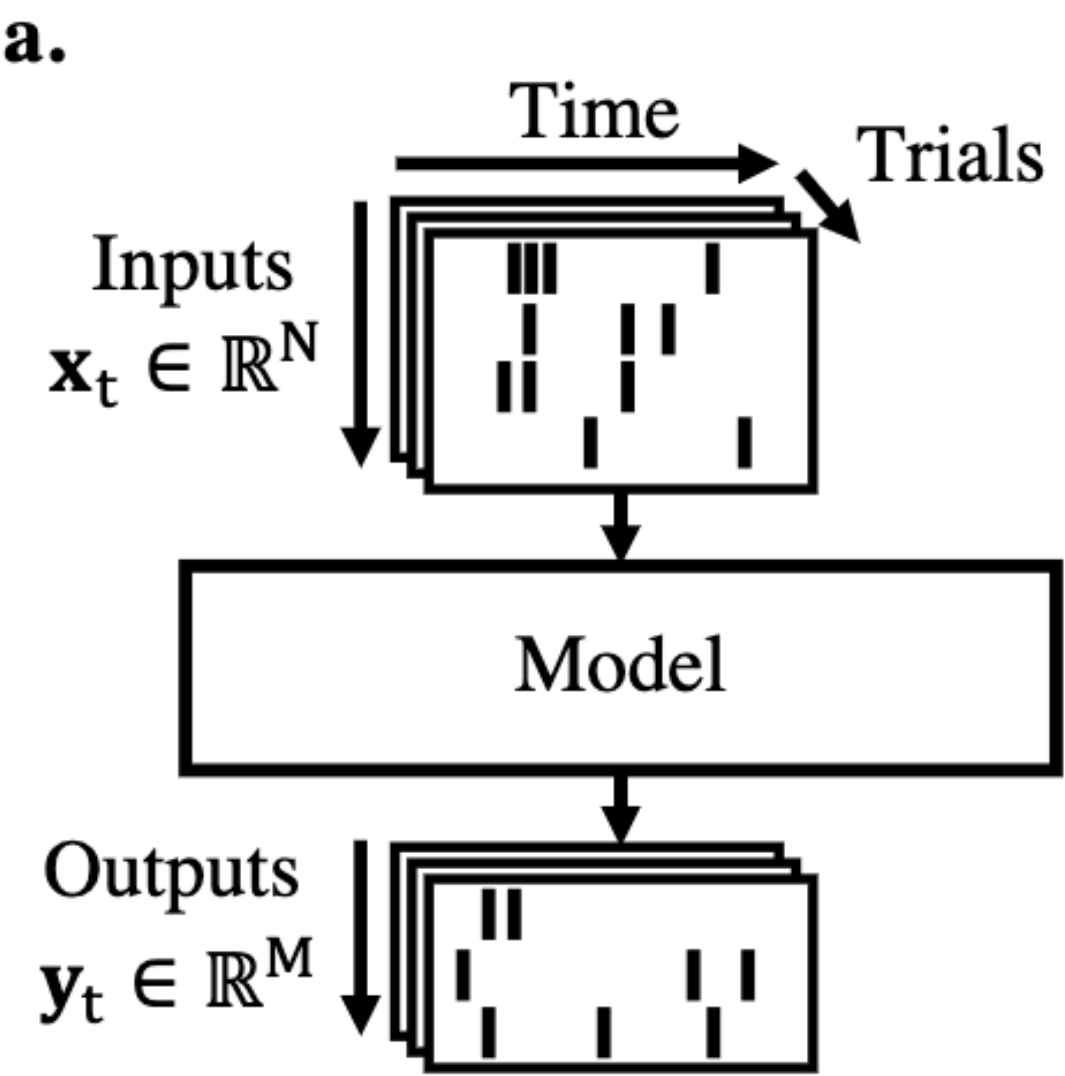


Table 1: Co-smoothing (in units of bits-per-spike) metric on MC Maze and DMFC RSG benchmarks [Pei et al., 2021] for S5 compared to SOTA methods. Note: we exclude ensemble methods and only consider single models.

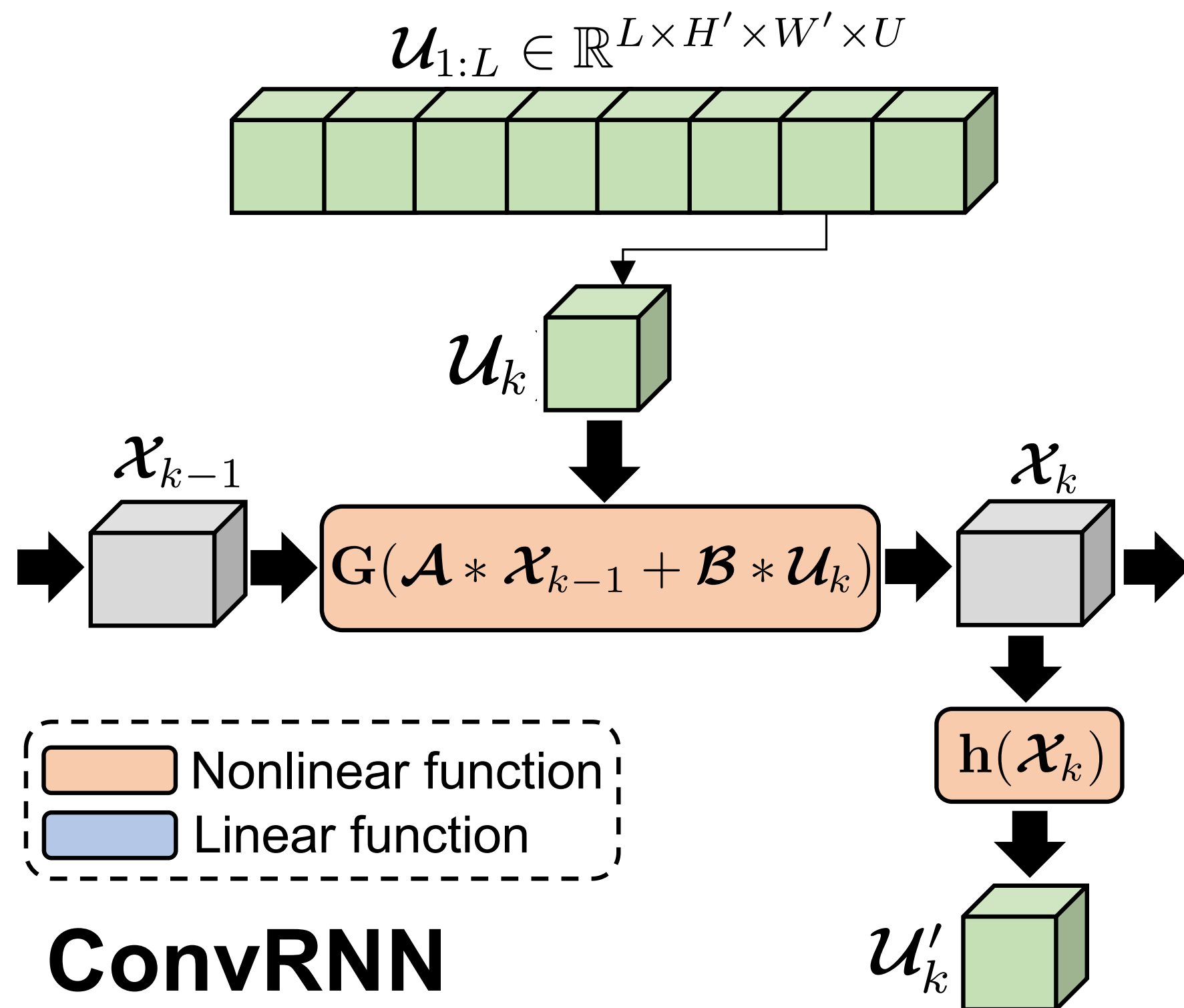
Method	MC Maze (\uparrow)	DMFC RSG (\uparrow)
S5 (Ours)	0.3826	0.1981
SSLFADS	0.3748	N/A
STNDT	0.3691	0.1859
iLQR-VAE	0.3559	N/A
Neural RoBERTa	0.3551	N/A
RNNf	0.3382	0.1781
AutoLFADS	0.3364	0.1829
MINT	0.3304	0.1821
NDT	0.3229	0.1720
SLDS	0.2249	0.1243

w/ Hyun Lee



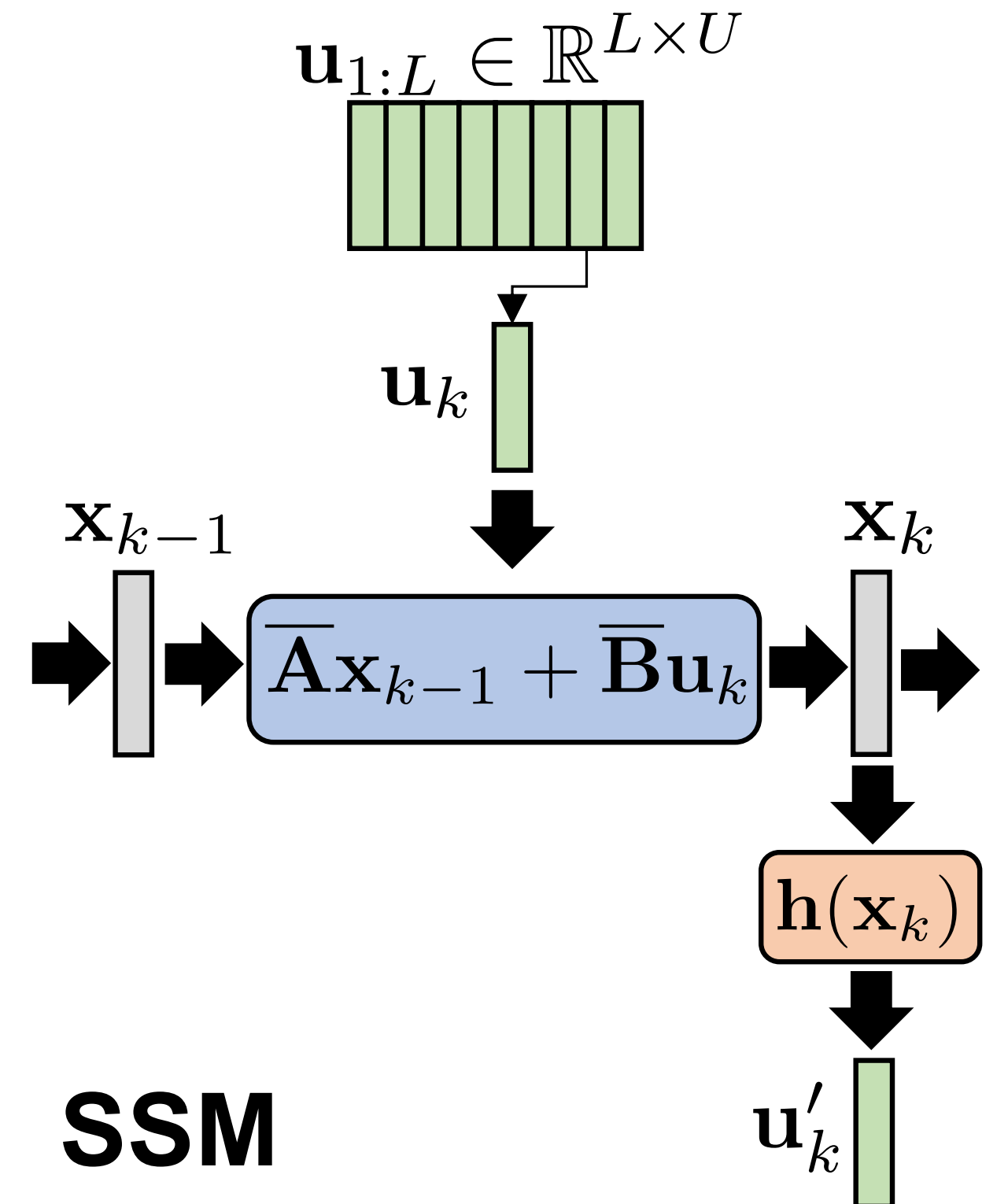
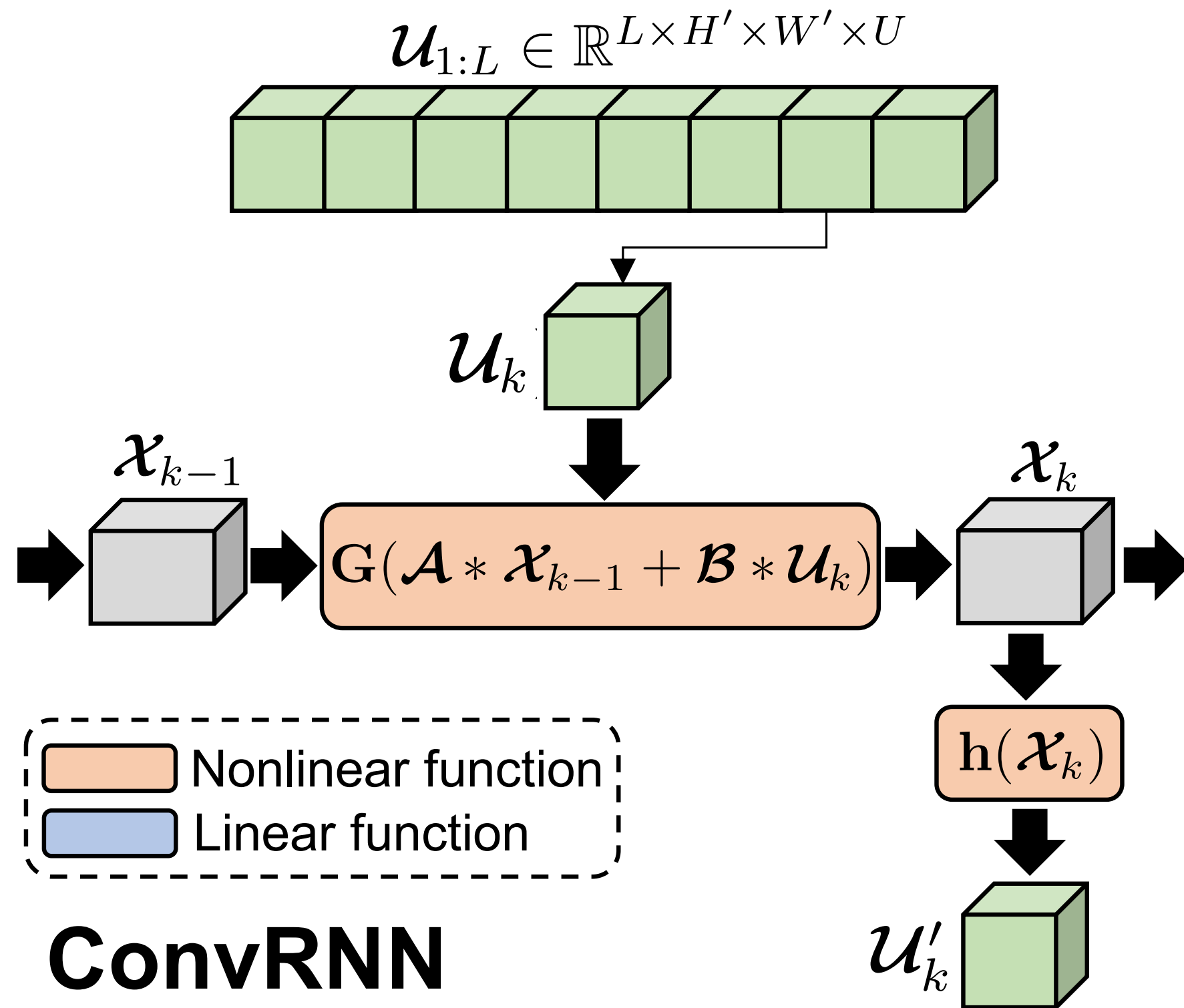
ConvSSMs: ConvRNN spatiotemporal inductive biases + SSM efficiency and long range modeling

ConvSSMs: ConvRNN spatiotemporal inductive biases + SSM efficiency and long range modeling



- ✓ Spatiotemporal modeling
- ✓ Fast autoregressive generation
- ✗ Slow, sequential training
- ✗ Vanishing gradients

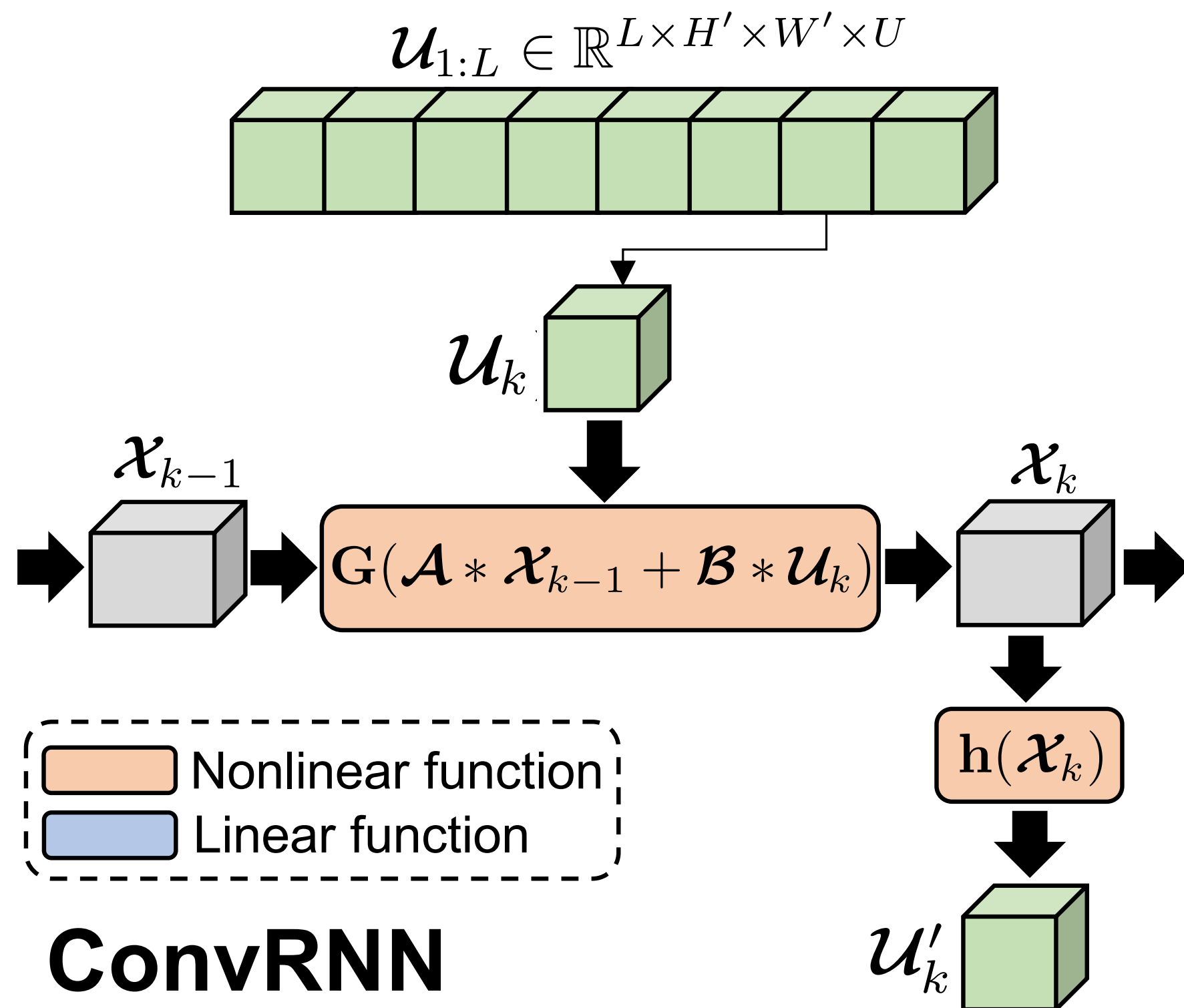
ConvSSMs: ConvRNN spatiotemporal inductive biases + SSM efficiency and long range modeling



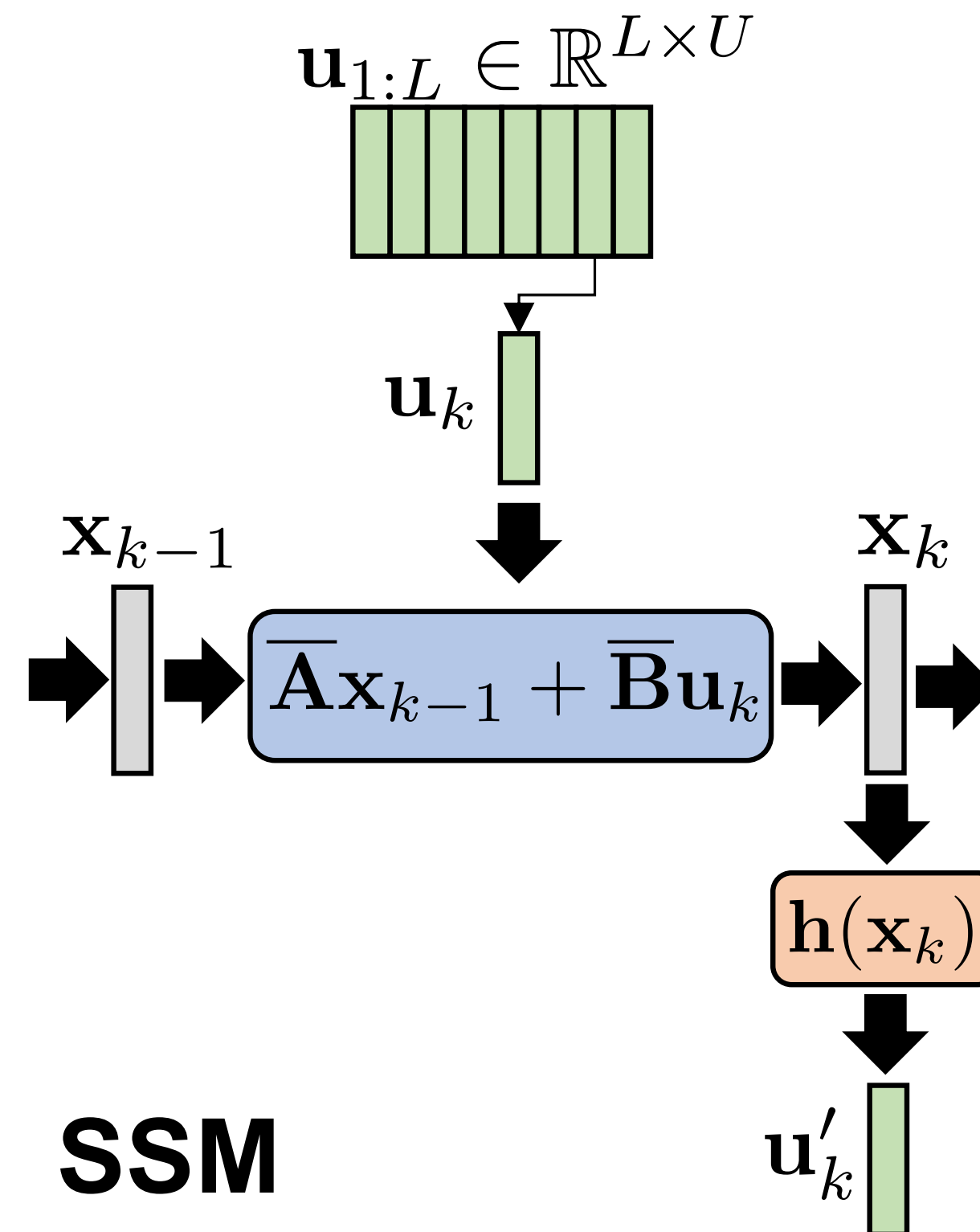
- ✓ Spatiotemporal modeling
- ✓ Fast autoregressive generation
- ✗ Slow, sequential training
- ✗ Vanishing gradients

- ✓ Efficient, parallel training
- ✓ Fast autoregressive generation
- ✓ Captures long-range dependencies
- ✗ Designed for vector-valued sequences

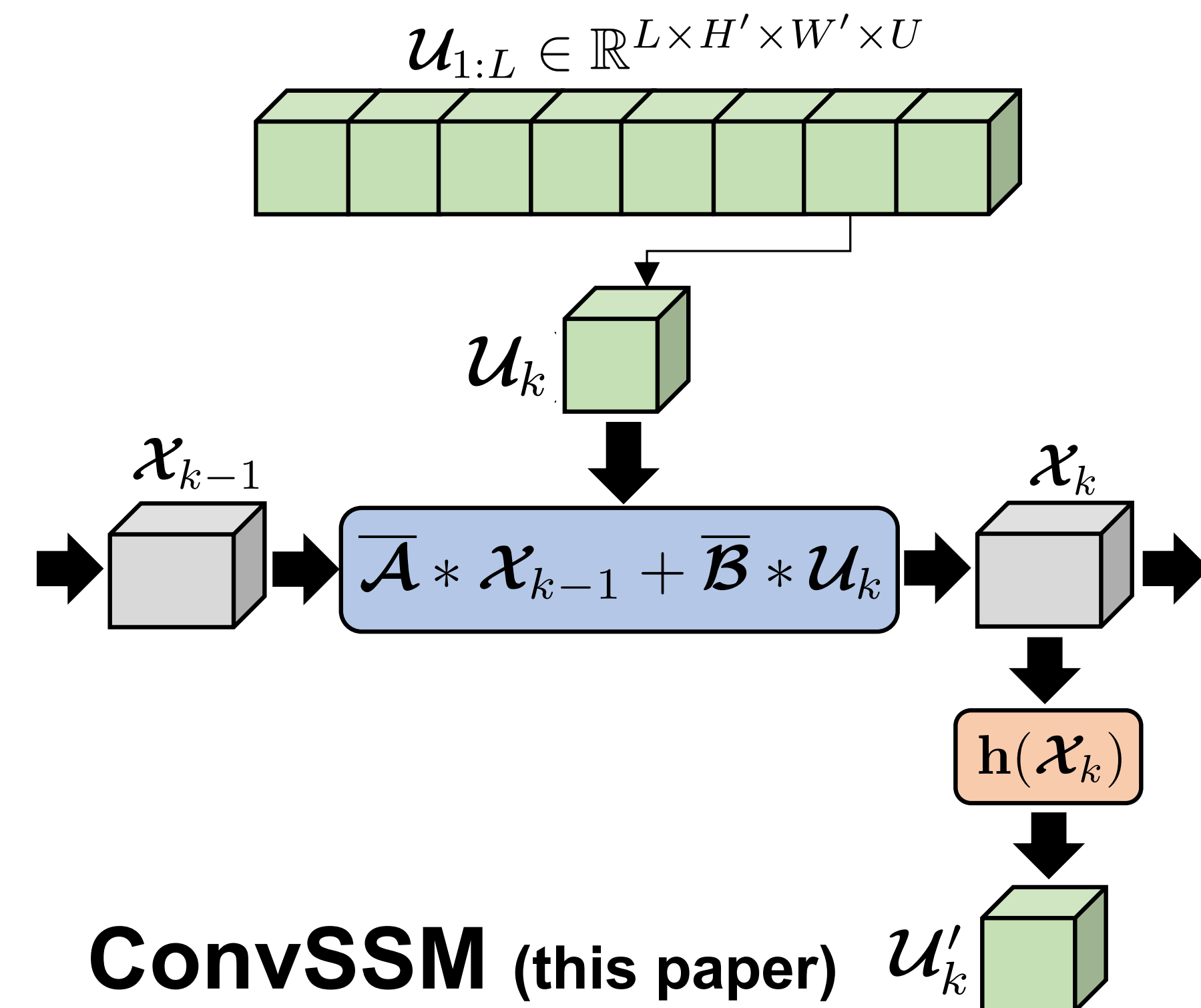
ConvSSMs: ConvRNN spatiotemporal inductive biases + SSM efficiency and long range modeling



- ✓ Spatiotemporal modeling
- ✓ Fast autoregressive generation
- ✗ Slow, sequential training
- ✗ Vanishing gradients

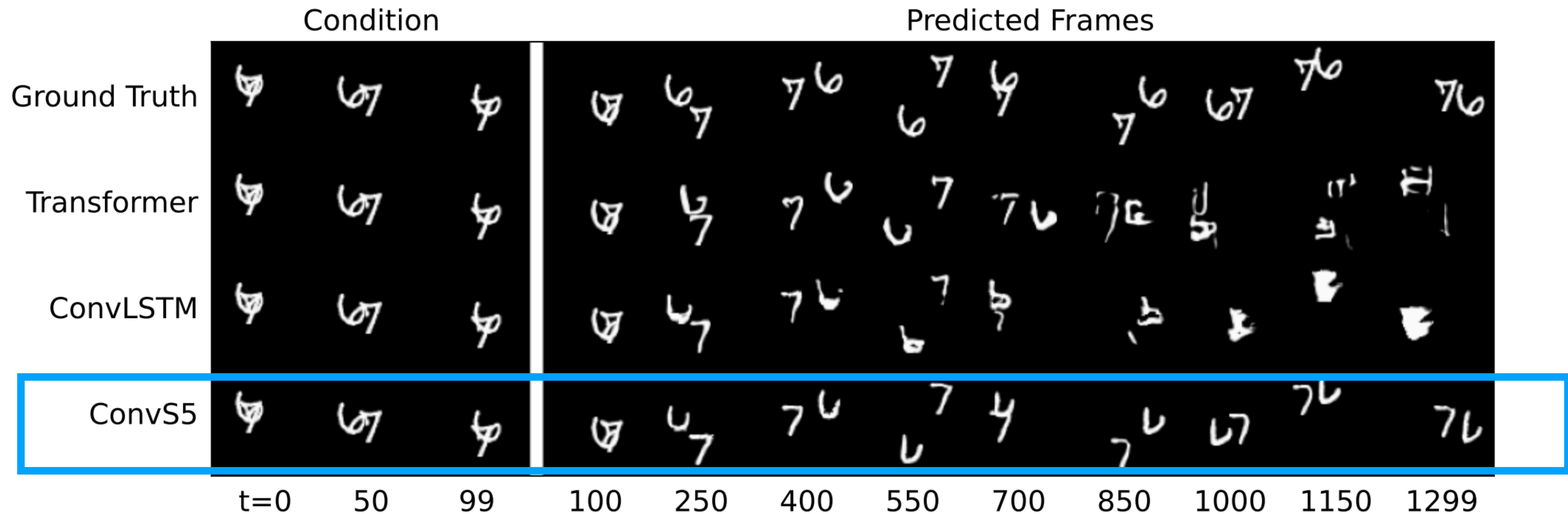


- ✓ Efficient, parallel training
- ✓ Fast autoregressive generation
- ✓ Captures long-range dependencies
- ✗ Designed for vector-valued sequences

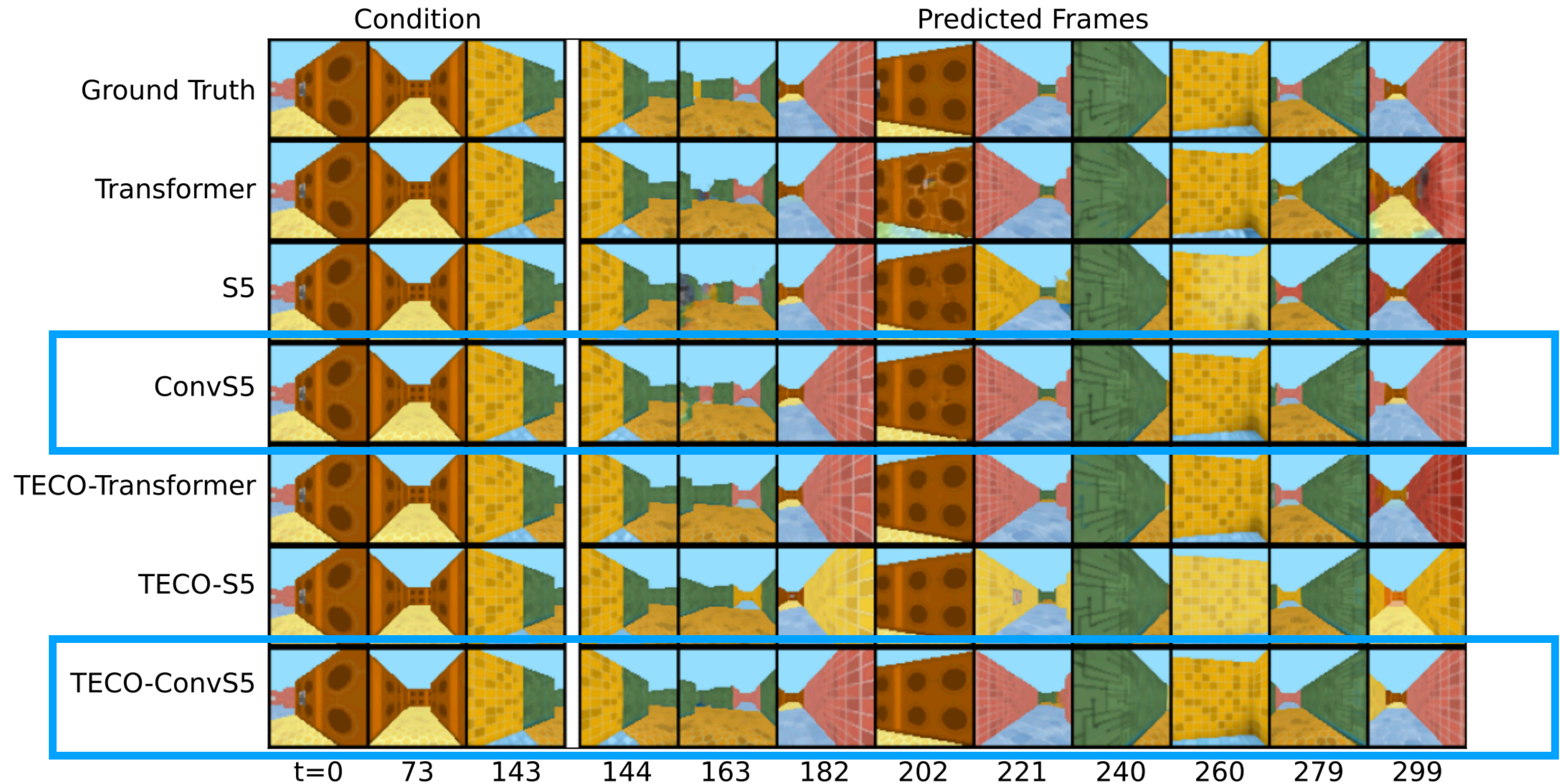


- ✓ Efficient, parallel training
- ✓ Fast autoregressive generation
- ✓ Captures long-range dependencies
- ✓ Spatiotemporal modeling

ConvSSMs: ConvRNN spatiotemporal inductive biases + SSM efficiency and long range modeling



ConvSSMs: ConvRNN spatiotemporal inductive biases + SSM efficiency and long range modeling



S5 enables new capabilities

S5 enables new capabilities

S5 can use linear time-varying (LTV) state space models:

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{B}(t)\mathbf{u}(t)$$

S5 enables new capabilities

S5 can use linear time-varying (LTV) state space models:

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{B}(t)\mathbf{u}(t)$$

- Context dependent dynamics

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_{1:k})\mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_{1:k})\mathbf{u}_k$$

S5 enables new capabilities

S5 can use linear time-varying (LTV) state space models:

$$\frac{dx(t)}{dt} = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{B}(t)\mathbf{u}(t)$$

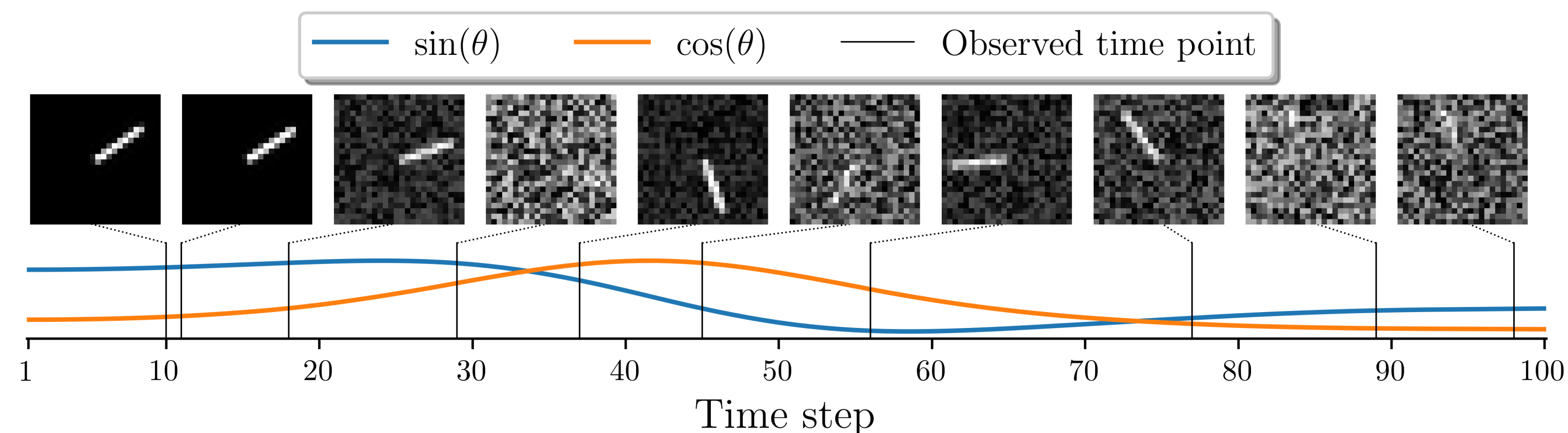
- **Context dependent dynamics**

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_{1:k})\mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_{1:k})\mathbf{u}_k$$

- **Irregularly sampled time-series**

$$\mathbf{x}_k = \overline{\mathbf{A}}(\Delta_k)\mathbf{x}_{k-1} + \overline{\mathbf{B}}(\Delta_k)\mathbf{u}_k$$

LTV example: Irregularly sampled time series



Model	Relative speed \uparrow	Regression MSE ($\times 10^{-3}$) \downarrow
mTAND	<u>12x</u>	65.64 (4.05)
RKN	1.9x	8.43 (0.61)
RKN- Δ_t	1.9x	5.09 (0.40)
GRU	3.0x	9.44 (1.00)
GRU- Δ_t	3.0x	5.44 (0.99)
Latent ODE	0.7x	15.70 (2.85)
ODE-RNN	1.0x	7.26 (0.41)
GRU-ODE-B	0.6x	9.78 (3.40)
f-CRU	1.2x	6.16 (0.88)
CRU	1.0x	4.63 (1.07)
CRU (our run)	1.0x	<u>3.94</u> (0.21)
S5	86x	3.41 (0.27)

Agenda

- Introduction, motivation, prior approaches
- Linear state space models (SSMs) overview
- S4, convolutions, parameterization
- S5, diagonalization, parallel scans
- **S6/Mamba/Hawk/Griffin, data-dependent dynamics**
- Challenges and opportunities

SSMs/RNNs vs Softmax Attention on Language

Deep SSMs, such as S4 and S5, mostly using LTI systems, have proven effective in a variety of data modalities: speech, image, video, reinforcement learning etc.

SSMs/RNNs vs Softmax Attention on Language

Deep SSMs, such as S4 and S5, mostly using LTI systems, have proven effective in a variety of data modalities: speech, image, video, reinforcement learning etc.

But language has proven troublesome compared to Softmax Attention.

SSMs/RNNs vs Softmax Attention on Language

Deep SSMs, such as S4 and S5, mostly using LTI systems, have proven effective in a variety of data modalities: speech, image, video, reinforcement learning etc.

But language has proven troublesome compared to Softmax Attention.

Several works, such as Zoology (Arora et al. 2023), suggests the ability to perform exact recall/retrieval/copying is extremely important for modeling language.

Hakuna Matata! It means no worries for the rest of your days! Hakuna Matata means no → worries
Key-Value Key-Value Query AR Hit Query AR Hit

SSMs/RNNs vs Softmax Attention on Language

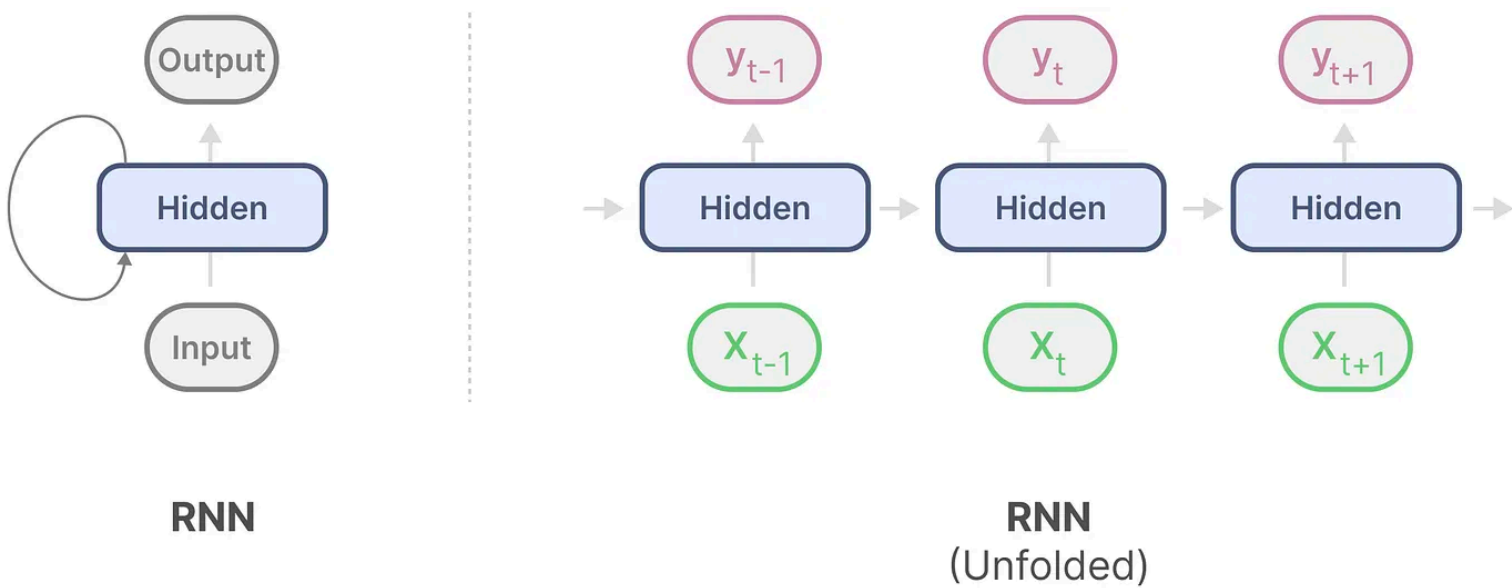
Deep SSMs, such as S4 and S5, mostly using LTI systems, have proven effective in a variety of data modalities: speech, image, video, reinforcement learning etc.

But language has proven troublesome compared to Softmax Attention.

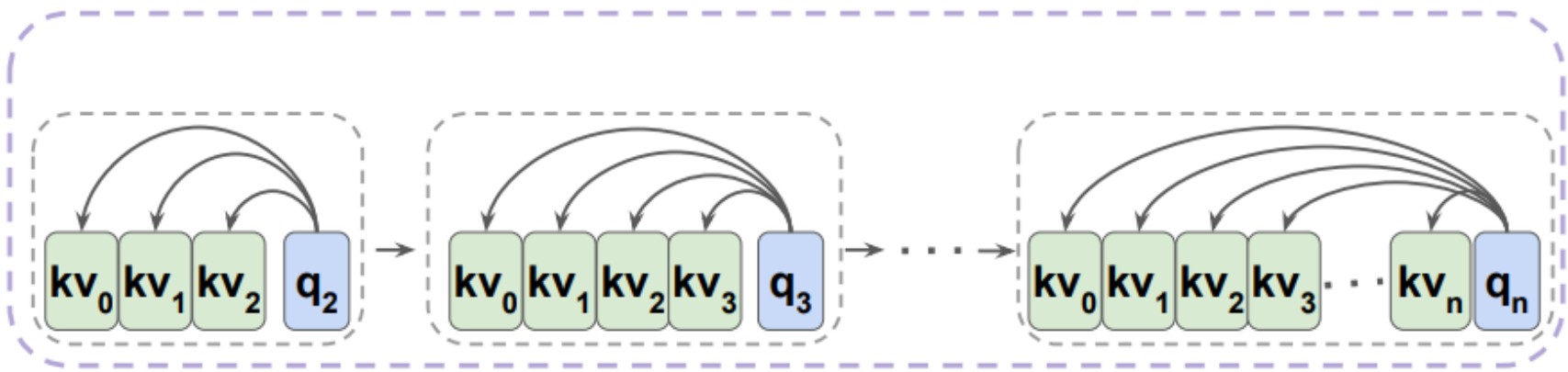
Several works, such as Zoology (Arora et al. 2023), suggests the ability to perform exact recall/retrieval/copying is extremely important for modeling language.

Hakuna Matata! It means no worries for the rest of your days! Hakuna Matata means no → worries
Key-Value Key-Value Query AR Hit Query AR Hit

But exact, lossless recall is difficult for fixed state models such as SSMs/RNNs compared to Softmax Attention.



Vs.



SSMs/RNNs vs Softmax Attention on Language

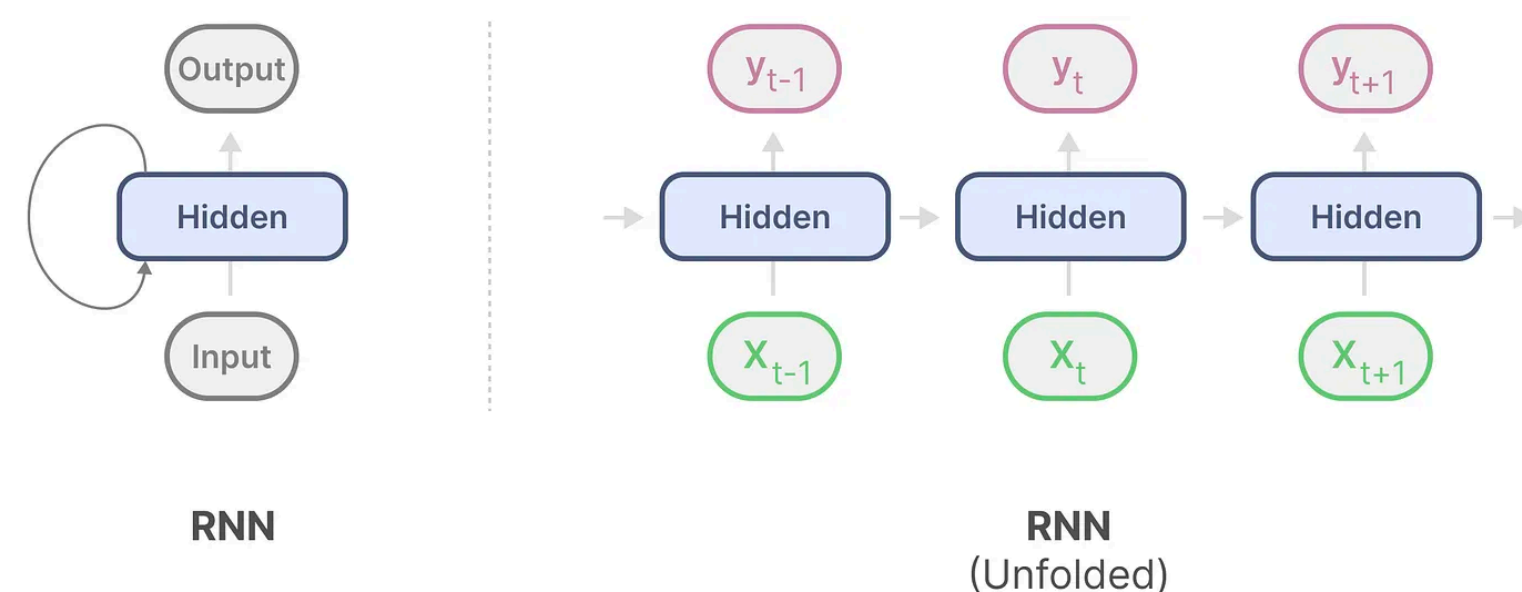
Deep SSMs, such as S4 and S5, mostly using LTI systems, have proven effective in a variety of data modalities: speech, image, video, reinforcement learning etc.

But language has proven troublesome compared to Softmax Attention.

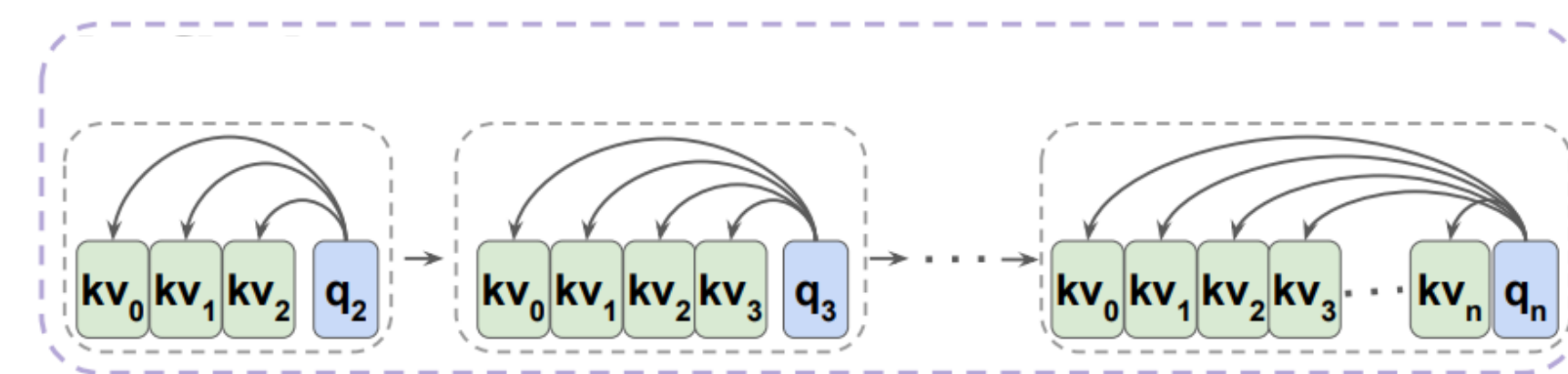
Several works, such as Zoology (Arora et al. 2023), suggests the ability to perform exact recall/retrieval/copying is extremely important for modeling language.

Hakuna Matata! It means no worries for the rest of your days! Hakuna Matata means no → worries
Key-Value Key-Value Query AR Hit Query AR Hit

But exact, lossless recall is difficult for fixed state models such as SSMs/RNNs compared to Softmax Attention.



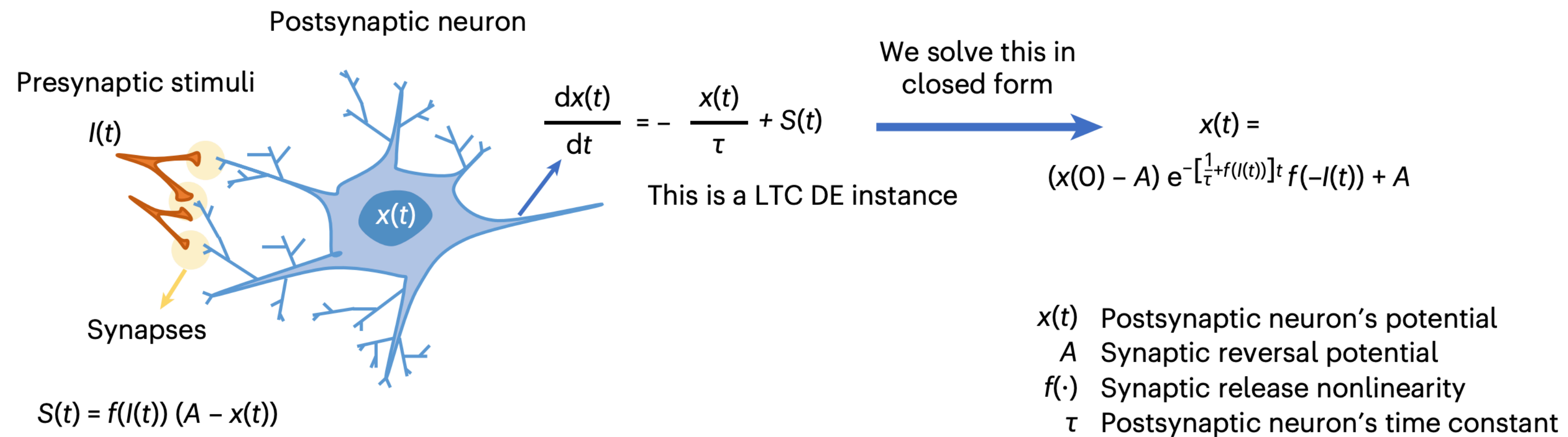
Vs.



Can we make better use of this fixed state with linear time-varying systems (LTV)?

Liquid Time-Constant Networks

$$\frac{d\mathbf{x}(t)}{dt} = - \left[\frac{1}{\tau} + f(\mathbf{x}(t), \mathbf{I}(t), t, \theta) \right] \mathbf{x}(t) + f(\mathbf{x}(t), \mathbf{I}(t), t, \theta) A.$$



LTV example: Liquid S4

$$x_k = (\bar{\mathbf{A}} + \bar{\mathbf{B}} u_k) x_{k-1} + \bar{\mathbf{B}} u_k, \quad y_k = \bar{\mathbf{C}} x_k$$

$$x_0 = \bar{\mathbf{B}} u_0, \quad y_0 = \bar{\mathbf{C}} \bar{\mathbf{B}} u_0$$

$$x_1 = \bar{\mathbf{A}} \bar{\mathbf{B}} u_0 + \bar{\mathbf{B}} u_1 + \bar{\mathbf{B}}^2 u_0 u_1, \quad y_1 = \bar{\mathbf{C}} \bar{\mathbf{A}} \bar{\mathbf{B}} u_0 + \bar{\mathbf{C}} \bar{\mathbf{B}} u_1 + \bar{\mathbf{C}} \bar{\mathbf{B}}^2 u_0 u_1$$

$$x_2 = \bar{\mathbf{A}}^2 \bar{\mathbf{B}} u_0 + \bar{\mathbf{A}} \bar{\mathbf{B}} u_1 + \bar{\mathbf{B}} u_2 + \bar{\mathbf{A}} \bar{\mathbf{B}}^2 u_0 u_1 + \bar{\mathbf{A}} \bar{\mathbf{B}}^2 u_0 u_2 + \bar{\mathbf{B}}^2 u_1 u_2 + \bar{\mathbf{B}}^3 u_0 u_1 u_2$$

$$y_2 = \bar{\mathbf{C}} \bar{\mathbf{A}}^2 \bar{\mathbf{B}} u_0 + \bar{\mathbf{C}} \bar{\mathbf{A}} \bar{\mathbf{B}} u_1 + \bar{\mathbf{C}} \bar{\mathbf{B}} u_2 + \bar{\mathbf{C}} \bar{\mathbf{A}} \bar{\mathbf{B}}^2 u_0 u_1 + \bar{\mathbf{C}} \bar{\mathbf{A}} \bar{\mathbf{B}}^2 u_0 u_2 + \bar{\mathbf{C}} \bar{\mathbf{B}}^2 u_1 u_2 + \bar{\mathbf{C}} \bar{\mathbf{B}}^3 u_0 u_1 u_2, \quad \dots$$

- Generally, LTV systems cannot be computed using convolutions.
- But Liquid-S4 work shows how this specific LTV form can be computed efficiently using convolutions.

Linear time-varying systems: S6/Mamba

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_{1:k})\mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_{1:k})\mathbf{u}_k$$

Linear time-varying systems: S6/Mamba

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_{1:k})\mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_{1:k})\mathbf{u}_k$$

S4 + S5 + Liquid S4 = S6:

Linear time-varying systems: S6/Mamba

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_{1:k})\mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_{1:k})\mathbf{u}_k$$

S4 + S5 + Liquid S4 = S6:

- Keeps the stack of SISO SSMs as in S4

Linear time-varying systems: S6/Mamba

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_{1:k})\mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_{1:k})\mathbf{u}_k$$

S4 + S5 + Liquid S4 = S6:

- Keeps the stack of SISO SSMs as in S4
- But uses a parallel scan like S5 (but with a clever hardware-aware algorithm) to allow LTV.

Linear time-varying systems: S6/Mamba

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_{1:k})\mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_{1:k})\mathbf{u}_k$$

S4 + S5 + Liquid S4 = S6:

- Keeps the stack of SISO SSMs as in S4
- But uses a parallel scan like S5 (but with a clever hardware-aware algorithm) to allow LTV.
- Time-varying, data-dependent SSM parameters, similar to Liquid-S4, but more general.

Linear time-varying systems: S6/Mamba

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_{1:k})\mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_{1:k})\mathbf{u}_k$$

S4 + S5 + Liquid S4 = S6:

- Keeps the stack of SISO SSMs as in S4
- But uses a parallel scan like S5 (but with a clever hardware-aware algorithm) to allow LTV.
- Time-varying, data-dependent SSM parameters, similar to Liquid-S4.

Algorithm 1 SSM (S4)	Algorithm 2 SSM + Selection (S6)
Input: $x : (B, L, D)$ Output: $y : (B, L, D)$ 1: $\underline{A} : (D, N) \leftarrow \text{Parameter}$ <div>▷ Represents structured $N \times N$ matrix</div> 2: $\underline{B} : (D, N) \leftarrow \text{Parameter}$ 3: $\underline{C} : (D, N) \leftarrow \text{Parameter}$ 4: $\underline{\Delta} : (D) \leftarrow \tau_{\Delta}(\text{Parameter})$ 5: $\underline{\underline{A}}, \underline{\underline{B}} : (D, \underline{N}) \leftarrow \text{discretize}(\underline{\Delta}, \underline{A}, \underline{B})$ 6: $y \leftarrow \text{SSM}(\underline{\underline{A}}, \underline{\underline{B}}, \underline{C})(x)$ <div>▷ Time-invariant: recurrence or convolution</div> 7: return y	Input: $x : (B, L, D)$ Output: $y : (B, L, D)$ 1: $\underline{A} : (D, N) \leftarrow \text{Parameter}$ <div>▷ Represents structured $N \times N$ matrix</div> 2: $\underline{B} : (B, L, N) \leftarrow s_B(x)$ 3: $\underline{C} : (B, L, N) \leftarrow s_C(x)$ 4: $\underline{\Delta} : (B, L, D) \leftarrow \tau_{\Delta}(\text{Parameter} + s_{\Delta}(x))$ 5: $\underline{\underline{A}}, \underline{\underline{B}} : (B, L, D, N) \leftarrow \text{discretize}(\underline{\Delta}, \underline{A}, \underline{B})$ 6: $y \leftarrow \text{SSM}(\underline{\underline{A}}, \underline{\underline{B}}, \underline{C})(x)$ <div>▷ Time-varying: recurrence (<i>scan</i>) only</div> 7: return y

Linear time-varying systems: S6/Mamba

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_{1:k})\mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_{1:k})\mathbf{u}_k$$

S4 + S5 + Liquid S4 = S6:

- Keeps the stack of SISO SSMs as in S4
- But uses a parallel scan like S5 (but with a clever hardware-aware algorithm) to allow LTV.
- Time-varying, data-dependent SSM parameters, similar to Liquid-S4, but more general.

Algorithm 1 SSM (S4)

Input: $x : (B, L, D)$

Output: $y : (B, L, D)$

```
1:  $\underline{A} : (D, N) \leftarrow \text{Parameter}$   
     $\triangleright$  Represents structured  $N \times N$  matrix  
2:  $\underline{B} : (D, N) \leftarrow \text{Parameter}$   
3:  $\underline{C} : (D, N) \leftarrow \text{Parameter}$   
4:  $\underline{\Delta} : (D) \leftarrow \tau_{\Delta}(\text{Parameter})$   
5:  $\underline{\underline{A}}, \underline{\underline{B}} : (D, \underline{N}) \leftarrow \text{discretize}(\underline{\Delta}, \underline{A}, \underline{B})$   
6:  $y \leftarrow \text{SSM}(\underline{\underline{A}}, \underline{\underline{B}}, \underline{C})(x)$   
     $\triangleright$  Time-invariant: recurrence or convolution  
7: return  $y$ 
```

Algorithm 2 SSM + Selection (S6)

Input: $x : (B, L, D)$

Output: $y : (B, L, D)$

```
1:  $\underline{A} : (D, N) \leftarrow \text{Parameter}$   
     $\triangleright$  Represents structured  $N \times N$  matrix  
2:  $\underline{B} : (B, L, N) \leftarrow s_B(x)$   
3:  $\underline{C} : (B, L, N) \leftarrow s_C(x)$   
4:  $\underline{\Delta} : (B, L, D) \leftarrow \tau_{\Delta}(\text{Parameter} + s_{\Delta}(x))$   
5:  $\underline{\underline{A}}, \underline{\underline{B}} : (B, L, D, N) \leftarrow \text{discretize}(\underline{\Delta}, \underline{A}, \underline{B})$   
6:  $y \leftarrow \text{SSM}(\underline{\underline{A}}, \underline{\underline{B}}, \underline{C})(x)$   
     $\triangleright$  Time-varying: recurrence (scan) only  
7: return  $y$ 
```

Time varying dynamics allows for ignoring irrelevant inputs, or forgetting information that is no longer important to remember.

Linear time-varying systems: S6/Mamba

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_k) \mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_k) \mathbf{u}_k$$

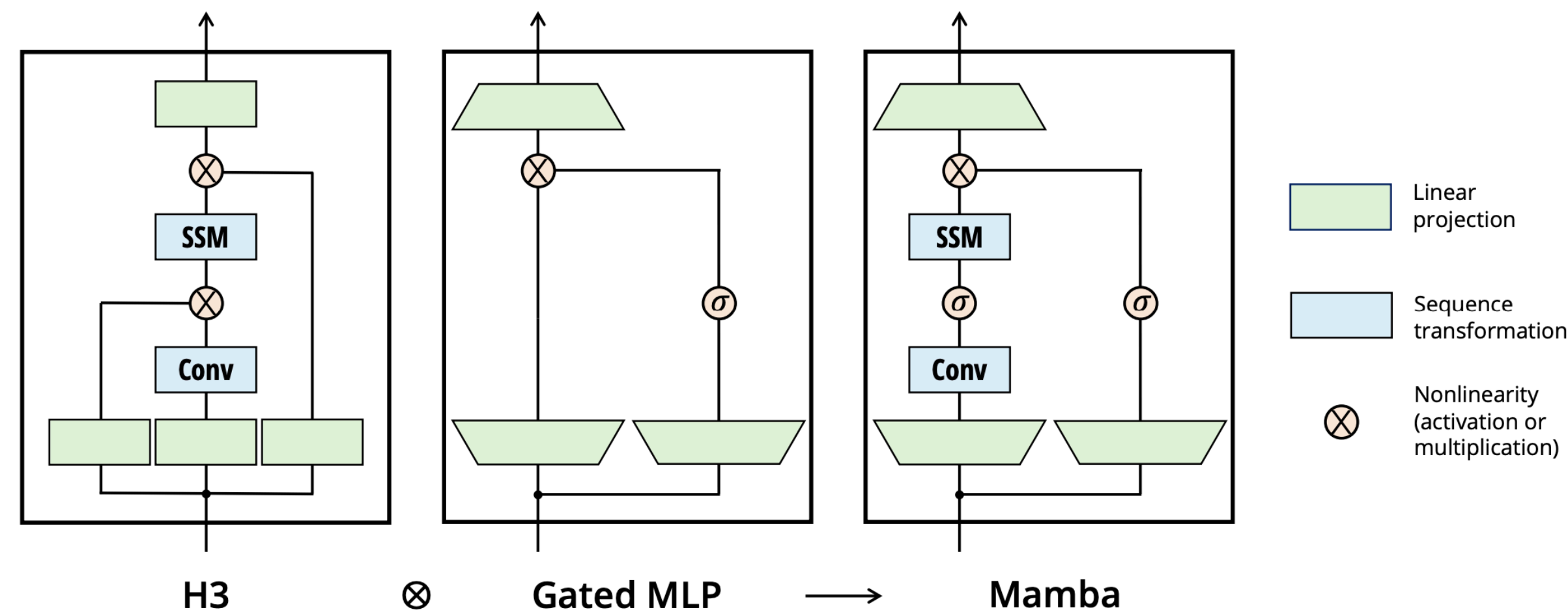
Algorithm 1 SSM (S4)	Algorithm 2 SSM + Selection (S6)
Input: $x : (B, L, D)$ Output: $y : (B, L, D)$ 1: $\mathbf{A} : (D, N) \leftarrow \text{Parameter}$ \triangleright Represents structured $N \times N$ matrix 2: $\mathbf{B} : (D, N) \leftarrow \text{Parameter}$ 3: $\mathbf{C} : (D, N) \leftarrow \text{Parameter}$ 4: $\Delta : (D) \leftarrow \tau_\Delta(\text{Parameter})$ 5: $\overline{\mathbf{A}}, \overline{\mathbf{B}} : (D, N) \leftarrow \text{discretize}(\Delta, \mathbf{A}, \mathbf{B})$ 6: $y \leftarrow \text{SSM}(\overline{\mathbf{A}}, \overline{\mathbf{B}}, \mathbf{C})(x)$ \triangleright Time-invariant: recurrence or convolution 7: return y	Input: $x : (B, L, D)$ Output: $y : (B, L, D)$ 1: $\mathbf{A} : (D, N) \leftarrow \text{Parameter}$ \triangleright Represents structured $N \times N$ matrix 2: $\mathbf{B} : (B, L, N) \leftarrow s_B(x)$ 3: $\mathbf{C} : (B, L, N) \leftarrow s_C(x)$ 4: $\Delta : (B, L, D) \leftarrow \tau_\Delta(\text{Parameter} + s_\Delta(x))$ 5: $\overline{\mathbf{A}}, \overline{\mathbf{B}} : (B, L, D, N) \leftarrow \text{discretize}(\Delta, \mathbf{A}, \mathbf{B})$ 6: $y \leftarrow \text{SSM}(\overline{\mathbf{A}}, \overline{\mathbf{B}}, \mathbf{C})(x)$ \triangleright Time-varying: recurrence (<i>scan</i>) only 7: return y

Linear time-varying systems: S6/Mamba

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_k) \mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_k) \mathbf{u}_k$$

Algorithm 1 SSM (S4)	Algorithm 2 SSM + Selection (S6)
Input: $x : (B, L, D)$ Output: $y : (B, L, D)$ 1: $\mathbf{A} : (D, N) \leftarrow \text{Parameter}$ ▷ Represents structured $N \times N$ matrix 2: $\mathbf{B} : (D, N) \leftarrow \text{Parameter}$ 3: $\mathbf{C} : (D, N) \leftarrow \text{Parameter}$ 4: $\Delta : (D) \leftarrow \tau_\Delta(\text{Parameter})$ 5: $\overline{\mathbf{A}}, \overline{\mathbf{B}} : (D, N) \leftarrow \text{discretize}(\Delta, \mathbf{A}, \mathbf{B})$ 6: $y \leftarrow \text{SSM}(\overline{\mathbf{A}}, \overline{\mathbf{B}}, \mathbf{C})(x)$ ▷ Time-invariant: recurrence or convolution 7: return y	Input: $x : (B, L, D)$ Output: $y : (B, L, D)$ 1: $\mathbf{A} : (D, N) \leftarrow \text{Parameter}$ ▷ Represents structured $N \times N$ matrix 2: $\mathbf{B} : (B, L, N) \leftarrow s_B(x)$ 3: $\mathbf{C} : (B, L, N) \leftarrow s_C(x)$ 4: $\Delta : (B, L, D) \leftarrow \tau_\Delta(\text{Parameter} + s_\Delta(x))$ 5: $\overline{\mathbf{A}}, \overline{\mathbf{B}} : (B, L, D, N) \leftarrow \text{discretize}(\Delta, \mathbf{A}, \mathbf{B})$ 6: $y \leftarrow \text{SSM}(\overline{\mathbf{A}}, \overline{\mathbf{B}}, \mathbf{C})(x)$ ▷ Time-varying: recurrence (<i>scan</i>) only 7: return y

Block design:

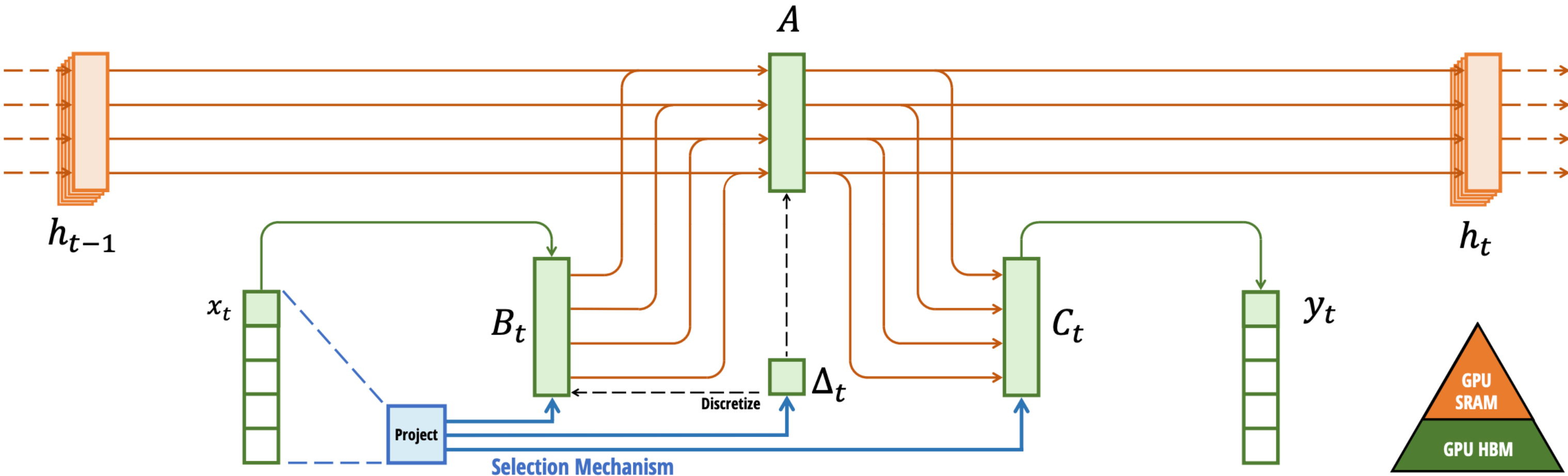


Linear time-varying systems: S6/Mamba

$$\mathbf{x}_k = \overline{\mathbf{A}}(\mathbf{u}_k) \mathbf{x}_{k-1} + \overline{\mathbf{B}}(\mathbf{u}_k) \mathbf{u}_k$$

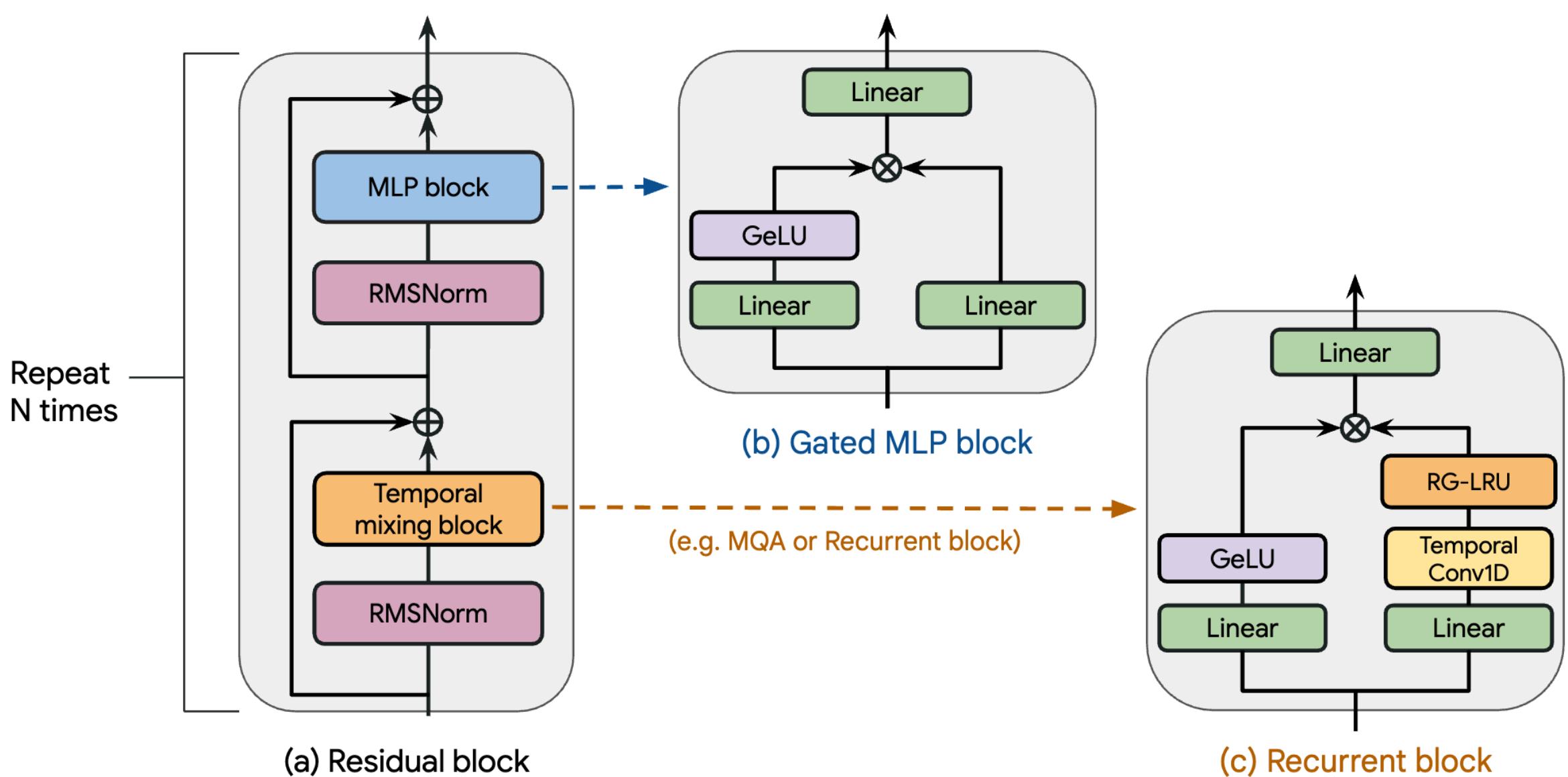
Algorithm 1 SSM (S4)	Algorithm 2 SSM + Selection (S6)
Input: $x : (B, L, D)$ Output: $y : (B, L, D)$ 1: $A : (D, N) \leftarrow \text{Parameter}$ ▷ Represents structured $N \times N$ matrix 2: $B : (D, N) \leftarrow \text{Parameter}$ 3: $C : (D, N) \leftarrow \text{Parameter}$ 4: $\Delta : (D) \leftarrow \tau_\Delta(\text{Parameter})$ 5: $\overline{A}, \overline{B} : (D, N) \leftarrow \text{discretize}(\Delta, A, B)$ 6: $y \leftarrow \text{SSM}(\overline{A}, \overline{B}, C)(x)$ ▷ Time-invariant: recurrence or convolution 7: return y	Input: $x : (B, L, D)$ Output: $y : (B, L, D)$ 1: $A : (D, N) \leftarrow \text{Parameter}$ ▷ Represents structured $N \times N$ matrix 2: $B : (B, L, N) \leftarrow s_B(x)$ 3: $C : (B, L, N) \leftarrow s_C(x)$ 4: $\Delta : (B, L, D) \leftarrow \tau_\Delta(\text{Parameter} + s_\Delta(x))$ 5: $\overline{A}, \overline{B} : (B, L, D, N) \leftarrow \text{discretize}(\Delta, A, B)$ 6: $y \leftarrow \text{SSM}(\overline{A}, \overline{B}, C)(x)$ ▷ Time-varying: recurrence (<i>scan</i>) only 7: return y

- Scans are limited by memory bandwidth
- This scan loads SSM params from slow HBM to fast SRAM, performs the discretization and recurrence in SRAM, and then writes outputs back to HBM

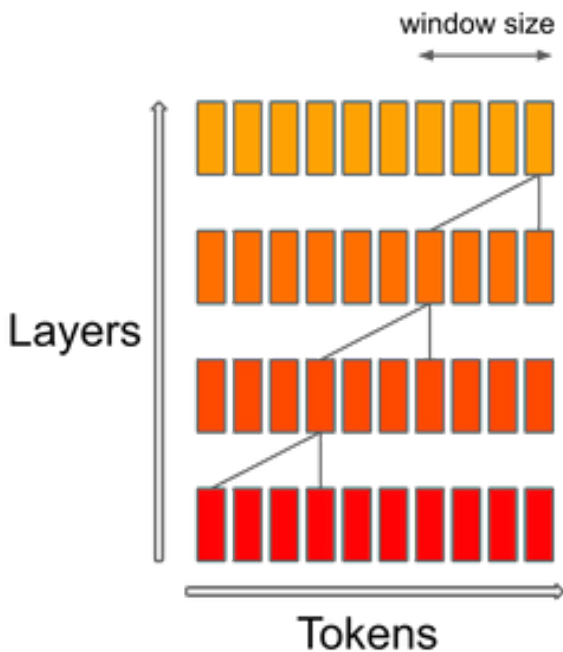


Linear time-varying systems: Hawk/Griffin

$$\begin{aligned} r_t &= \sigma(W_a x_t + b_a), && \text{recurrence gate} \\ i_t &= \sigma(W_x x_t + b_x), && \text{input gate} \\ a_t &= a^{cr_t}, \\ h_t &= a_t \odot h_{t-1} + \sqrt{1 - a_t^2} \odot (i_t \odot x_t). \end{aligned}$$



Griffin/Recurrent Gemma: Also include variant that includes sliding window attention in some layers



Linear time-varying systems: Language Results

Model Type	Model Size	Training Tokens	MMLU	HellaSwag	PIQA	WinoGrande	ARC-E	ARC-C	Average
Mamba	3B	600B	26.2	71.0	78.1	65.9	68.2	41.7	58.5
Llama-2	7B	2T	45.3	77.2	78.8	69.2	75.2	45.9	65.3
	13B	2T	54.8	80.7	80.5	72.8	77.3	49.4	69.3
MQA	1B	300B	28.9	64.8	75.0	62.0	60.2	35.4	54.4
Transformer (Baseline)	3B	300B	31.7	71.0	77.6	66.1	68.1	39.2	59.0
	6B	300B	38.9	77.0	79.5	70.4	74.1	45.2	64.2
Hawk	1B	300B	29.7	63.3	76.1	57.2	60.6	34.6	53.6
	3B	300B	31.3	71.7	78.8	66.5	68.4	40.2	59.5
	7B	300B	35.0	77.6	80.0	69.9	74.4	45.9	63.8
Griffin	1B	300B	29.5	67.2	77.4	65.2	67.0	36.9	57.2
	3B	300B	32.6	73.5	78.1	67.2	71.5	41.4	60.7
	7B	300B	39.3	78.6	81.0	72.6	75.4	47.9	65.8
	14B	300B	49.5	81.4	81.8	74.1	79.1	50.8	69.5

These models can sometimes match Transformer perplexity and standard academic benchmark results

Linear time-varying systems: Language Results

Model Type	Model Size	Training Tokens	MMLU	HellaSwag	PIQA	WinoGrande	ARC-E	ARC-C	Average
Mamba	3B	600B	26.2	71.0	78.1	65.9	68.2	41.7	58.5
Llama-2	7B	2T	45.3	77.2	78.8	69.2	75.2	45.9	65.3
	13B	2T	54.8	80.7	80.5	72.8	77.3	49.4	69.3
MQA	1B	300B	28.9	64.8	75.0	62.0	60.2	35.4	54.4
Transformer (Baseline)	3B	300B	31.7	71.0	77.6	66.1	68.1	39.2	59.0
	6B	300B	38.9	77.0	79.5	70.4	74.1	45.2	64.2
Hawk	1B	300B	29.7	63.3	76.1	57.2	60.6	34.6	53.6
	3B	300B	31.3	71.7	78.8	66.5	68.4	40.2	59.5
	7B	300B	35.0	77.6	80.0	69.9	74.4	45.9	63.8
Griffin	1B	300B	29.5	67.2	77.4	65.2	67.0	36.9	57.2
	3B	300B	32.6	73.5	78.1	67.2	71.5	41.4	60.7
	7B	300B	39.3	78.6	81.0	72.6	75.4	47.9	65.8
	14B	300B	49.5	81.4	81.8	74.1	79.1	50.8	69.5

These models can sometimes match Transformer perplexity and standard academic benchmark results

But do they have the same recall and in-context learning abilities expected of Transformers?

Agenda

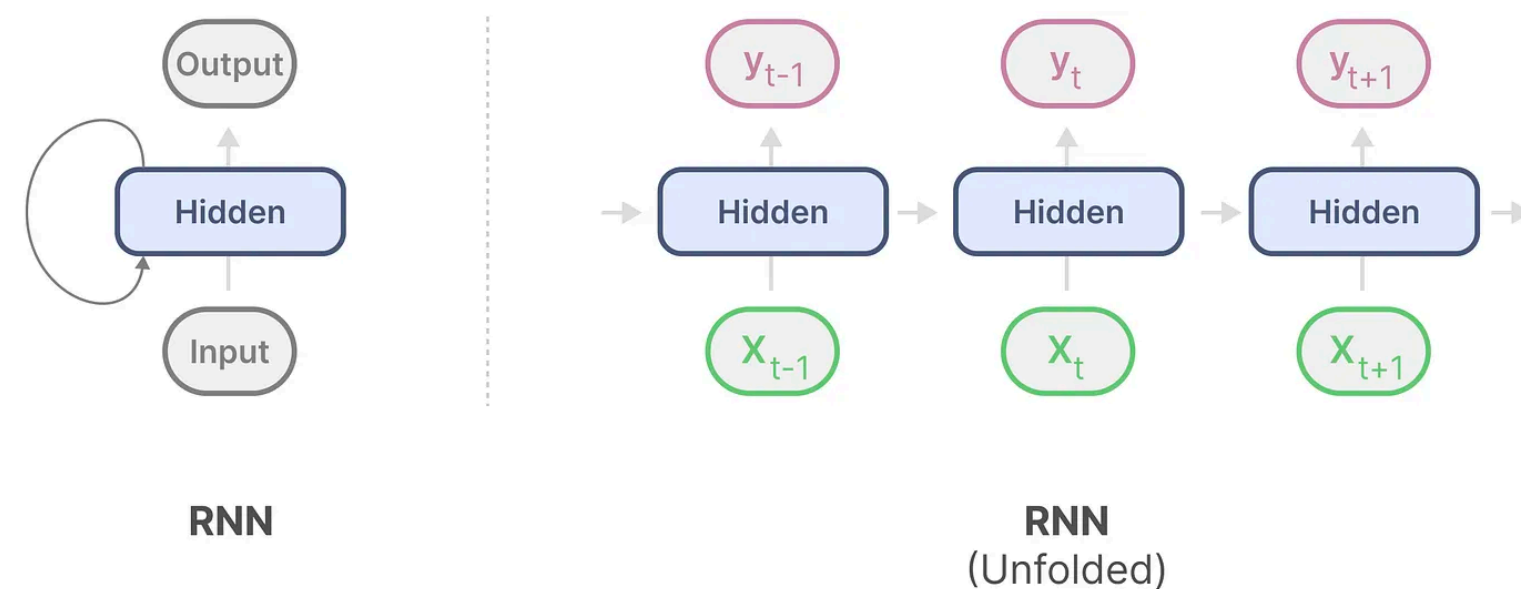
- Introduction, motivation, prior approaches
- Linear state space models (SSMs) overview
- S4, convolutions, parameterization
- S5, diagonalization, parallel scans
- S6/Mamba/Hawk/Griffin, data-dependent dynamics
- **Challenges and opportunities**

Recall challenges remain

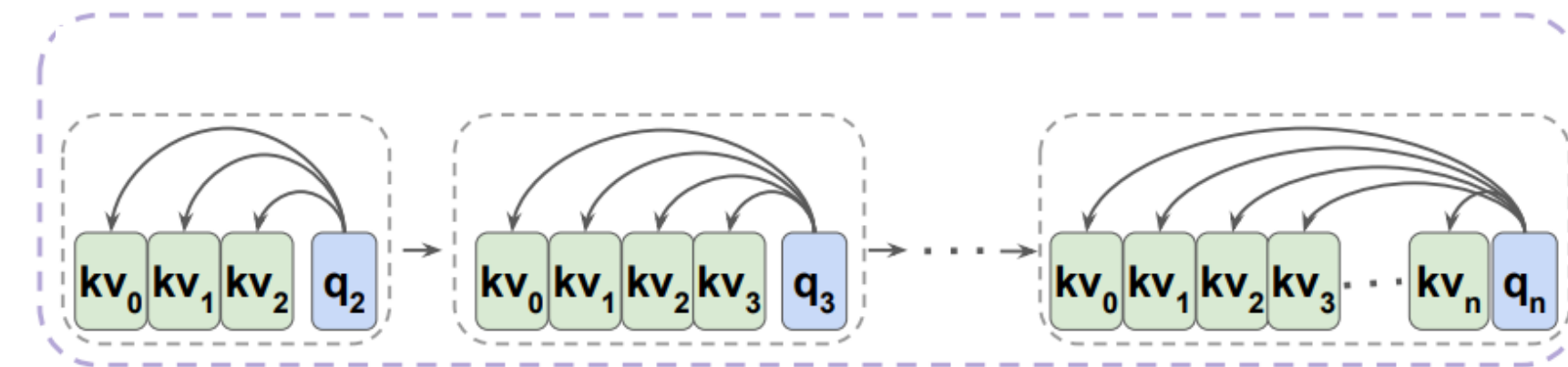
Many Mamba for X papers quickly followed showing strong results in vision, diffusion etc., suggesting these LTV systems can be very strong models.

But recall/copying problem in language seems to persist...:

- Repeat after me: Transformers are better than state space models at copying <https://arxiv.org/abs/2402.01032>
- Simple linear attention models balance recall-throughput tradeoff <https://arxiv.org/abs/2402.18668>
- Can Mamba learn how to learn? A comparative study on in-context learning tasks: <https://arxiv.org/abs/2402.04248>



Vs.



Recall challenges remain

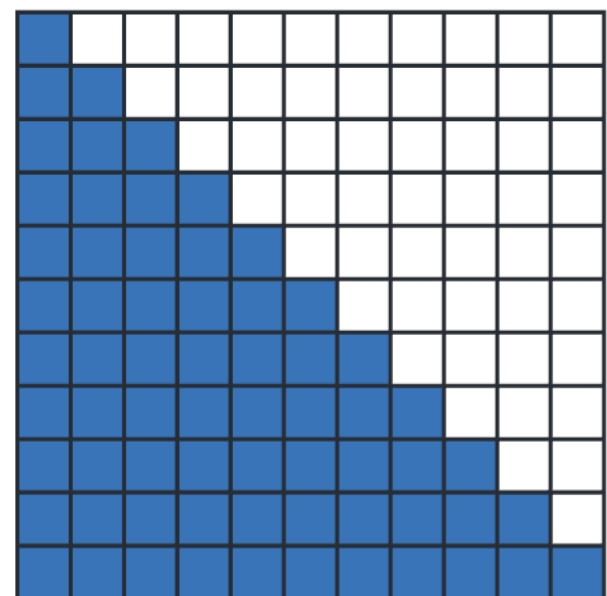
Griffin, Section 6.2 Copying and Retrieval capabilities:

*From Figure 6(c), we see that while **Hawk** can do reasonably well on the task for very short phonebook lengths, it **fails to memorize and retrieve** the correct phone number when the phonebook length grows, similar to the observation made by Jelassi et al. (2024)...Our **Transformer baseline** can **almost perfectly solve** this task up to the training sequence length...Interestingly, **Griffin** can perfectly solve this task up to a context length that matches its local attention **window size of 1024**...*

*Once the context length is long enough such that the local attention window **does not cover** the whole phonebook, **performance starts to degrade**...While the performance of Griffin is promising for the ability of models with fixed-size state to solve copying and retrieval tasks, **our results suggest more work is needed to improve these capabilities for such models.***

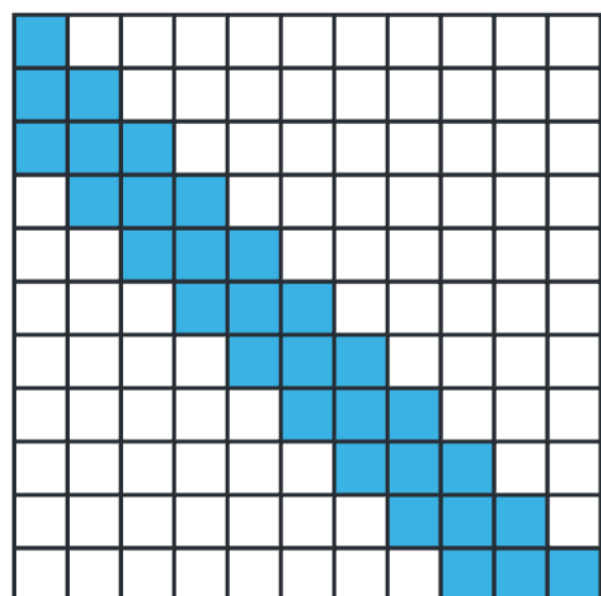
More sliding window hybrids

Linear Attention



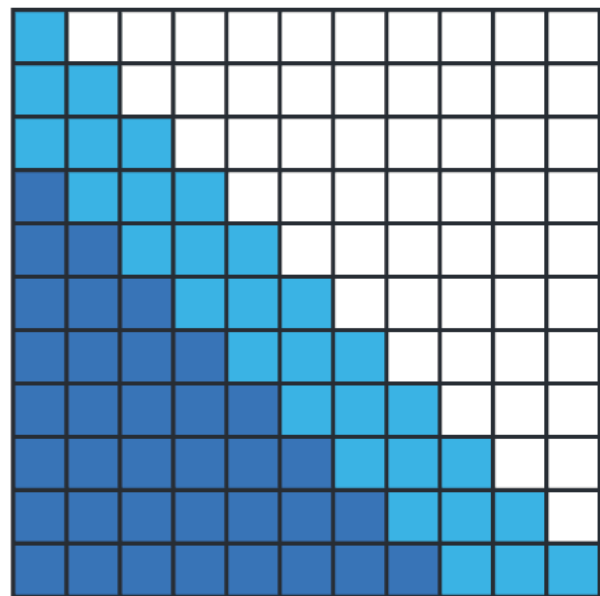
- ✓ Taylor approximation provides large memory for recall
- ✗ Precise local token shifts and comparison

Sliding Window

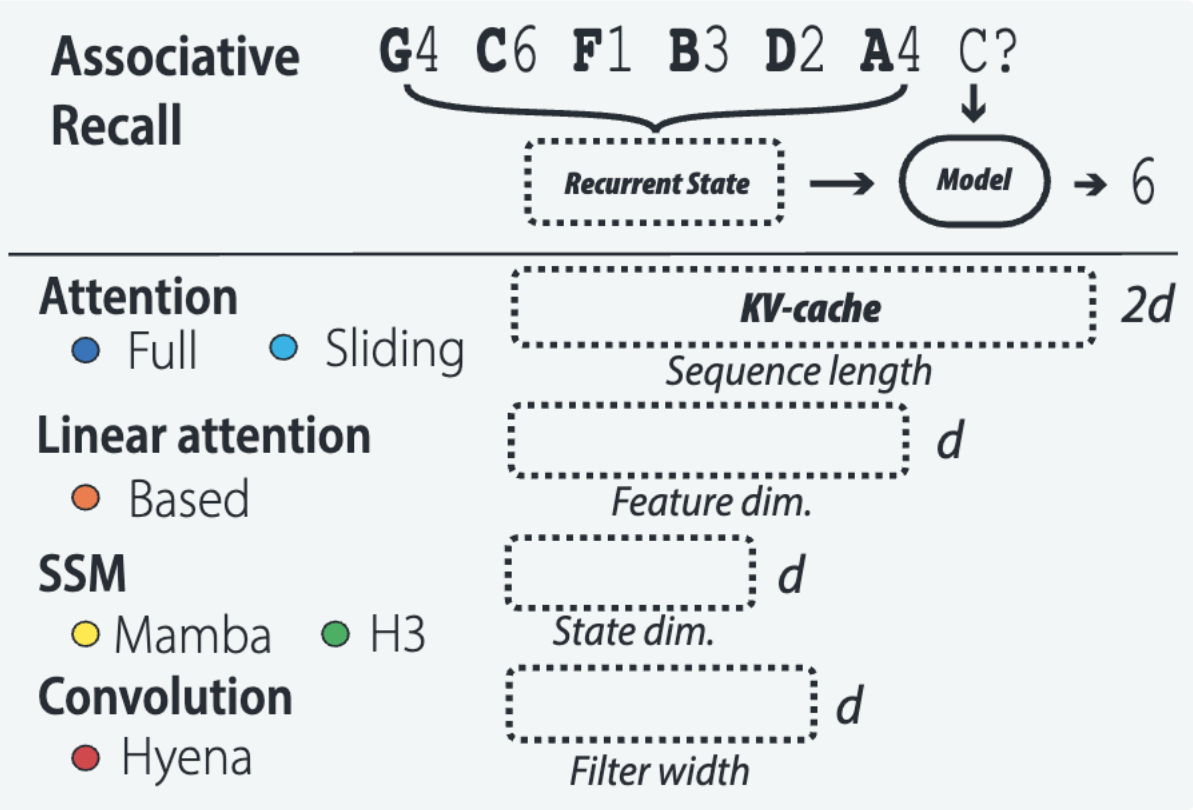
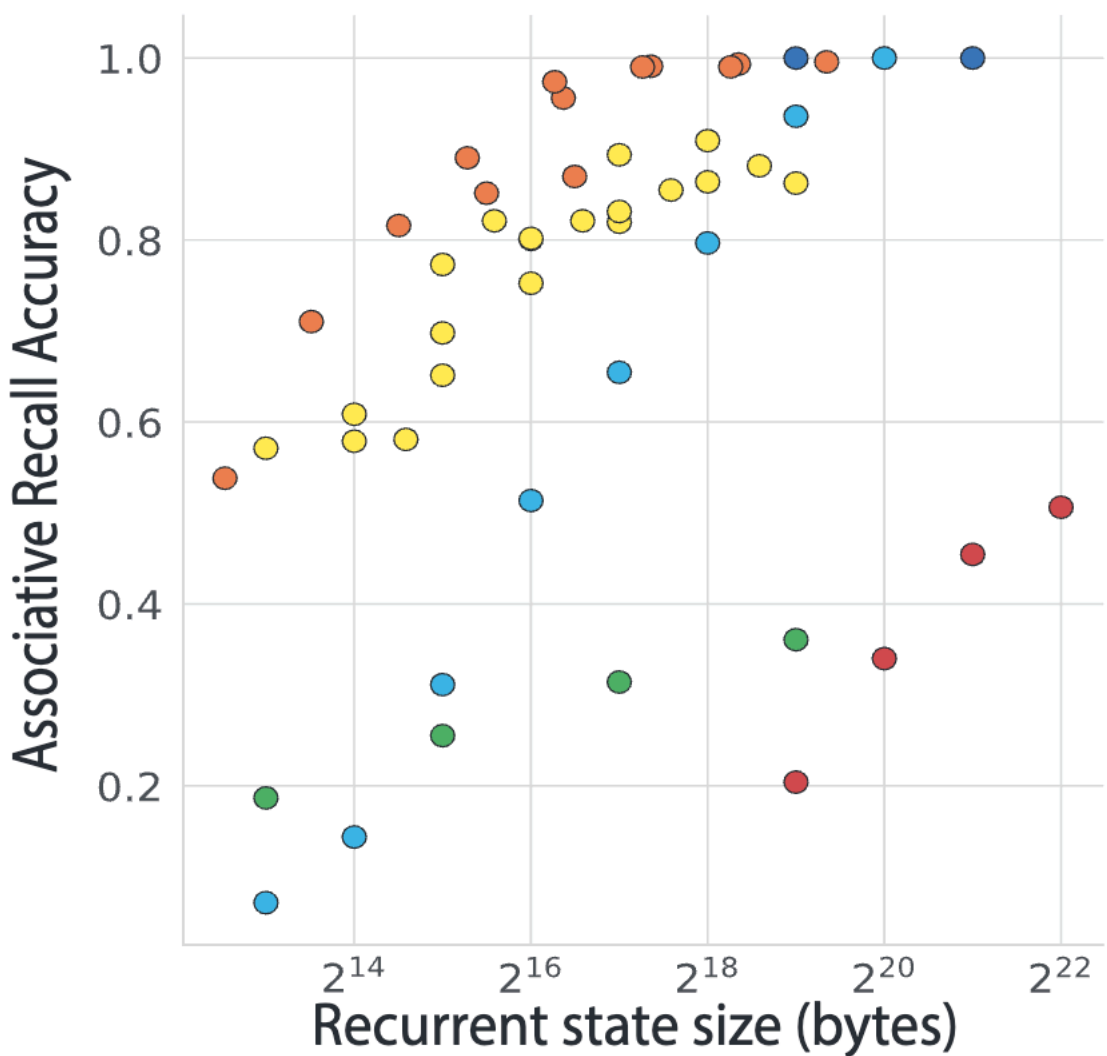


- ✗ Limited memory for long range recall
- ✓ Precise local token shifts and comparison

Based



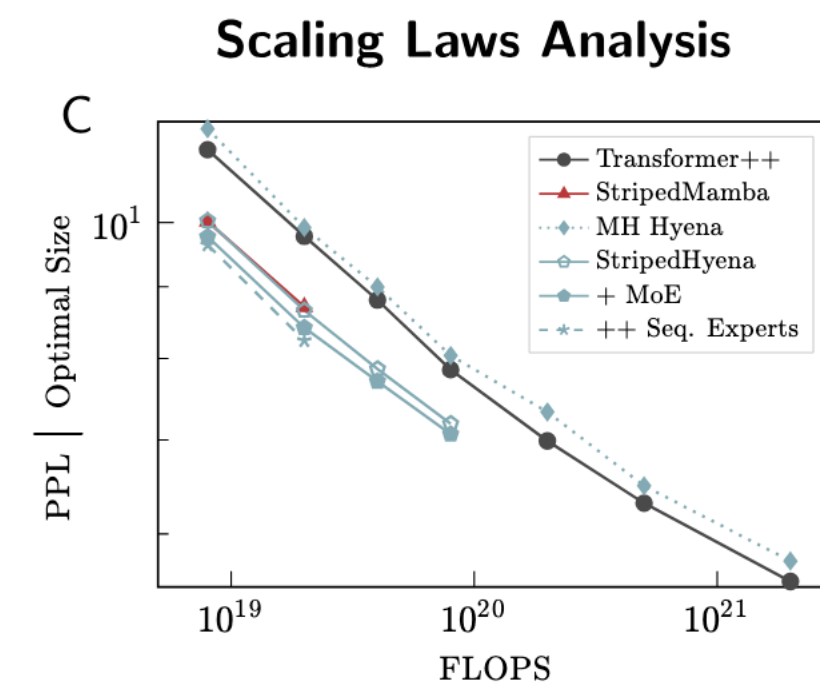
- ✓ Taylor approximation provides large memory for recall
- ✓ Precise local token shifts and comparison



Possible Solutions?

Possible Solutions?

- **Training flop/cost penalty: Can more tokens (or better data) help fixed state size methods make better use of state?**

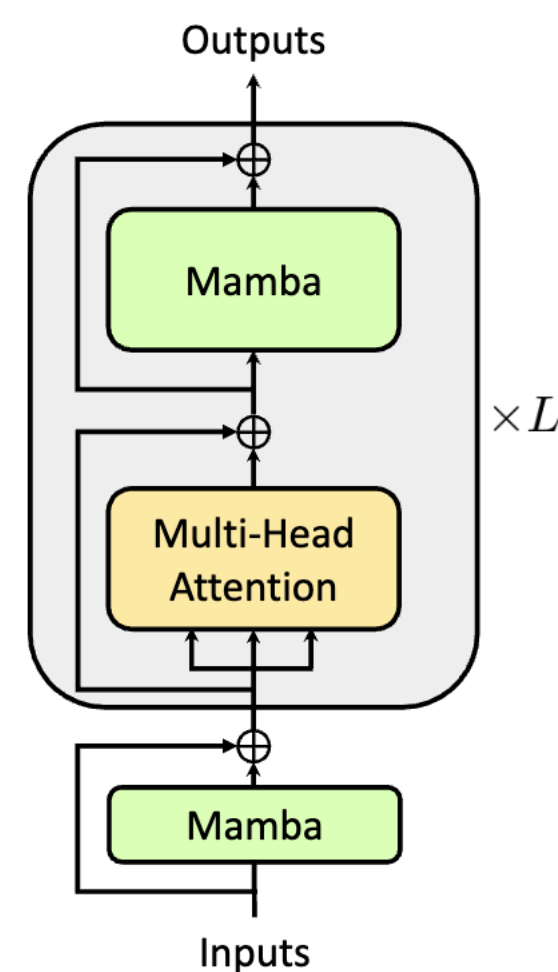


Poli et al 2024. Mechanistic Design and Scaling of Hybrid architectures

Possible Solutions?

- Training flop/cost penalty: Can more tokens (or better data) help fixed state size methods make better use of state?

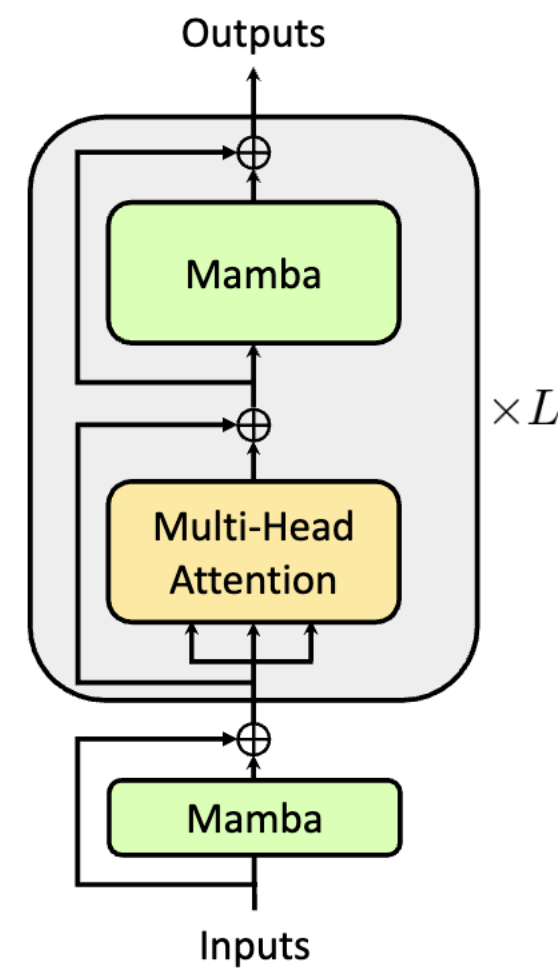
- Hybrid models?



Possible Solutions?

- **Training flop/cost penalty: Can more tokens (or better data) help fixed state size methods make better use of state?**

- **Hybrid models?**



- **Better dynamic storage/retrieval mechanisms?**

Taxonomy

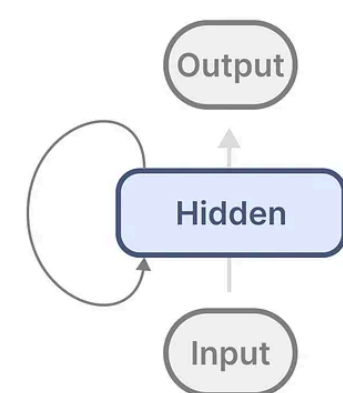
Fixed state size methods

Time-invariant

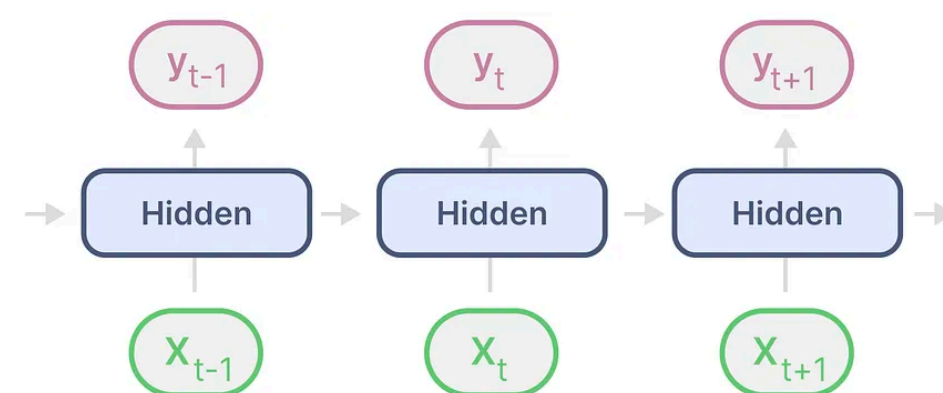
- LTI SSMs/RNNs (e.g. S4/S5/LRU)
- Linear Attention (e.g. Retnet)
- Long Convolutions (*Distilled)

Time-varying

- Nonlinear RNNs
- LTV SSMs/RNNs (e.g. Mamba/S6/Hawk/Griffin)
- Gated Linear Attention



RNN

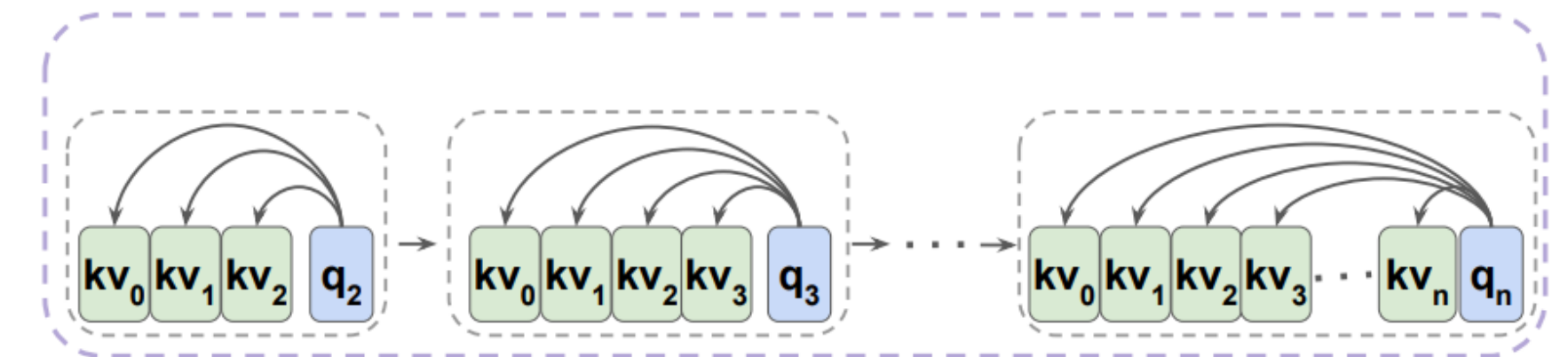


RNN
(Unfolded)

Vs.

Dynamic state size methods

- Softmax attention
- Other methods whose state grows with sequence



Thank you!