

SSMs: State Space Models

Prerequisites to Mamba

Plan

- Review: Sequence modelling
- Motivation for the new architecture
- Alternatives to attention:
 - Linear RNNs
 - State space models
 - Mathematical model
 - Discretization
 - Recurrent/Convolutional computation
 - Structure and Dimensions
- Prior work to Mamba: Architectures with SSMs

Review: Sequence modelling

- RNNs
- CNNs
- Transformer

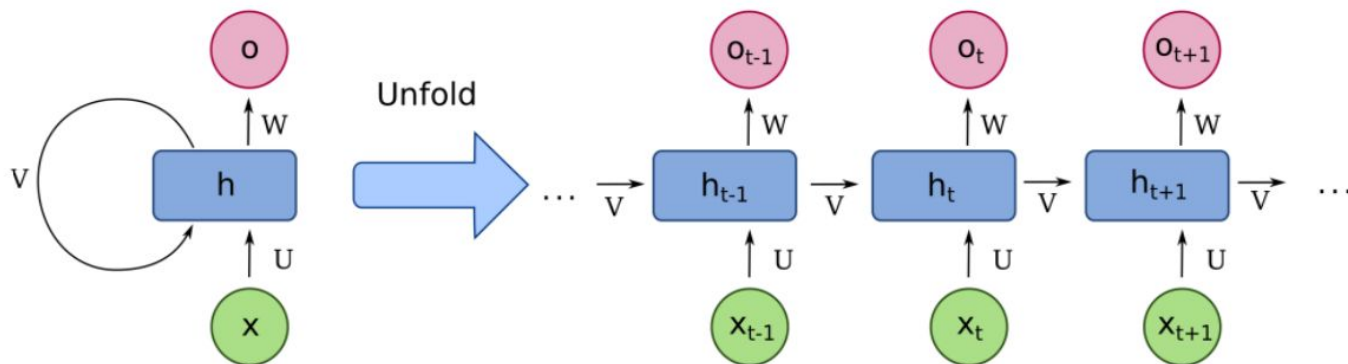
Review: RNNs

Theoretically with infinite context window

Practically suffer from vanishing/exploding gradient

Training: $O(N)$, not parallelizable

Inference: constant time for each token



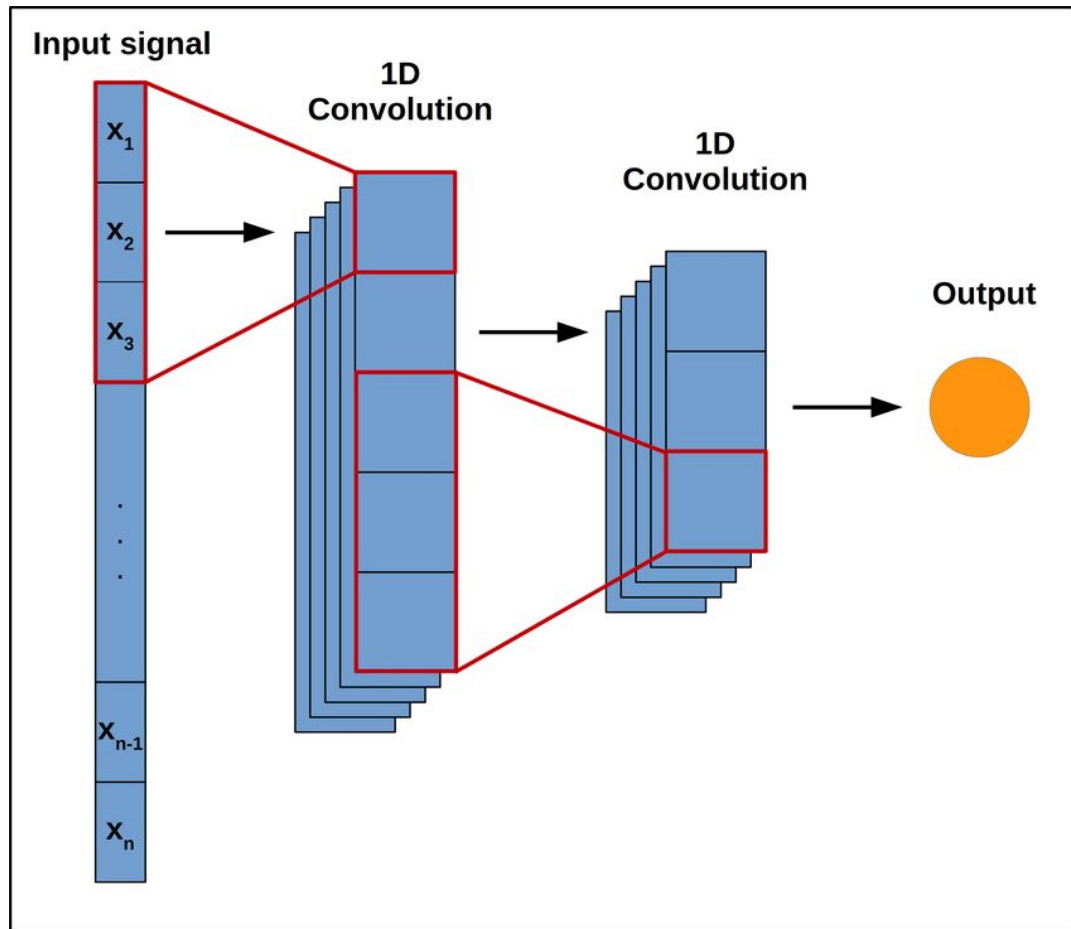
Review: CNNs

Finite context window
(depending on kernel size)

Need to materialize the kernel
before using it

Training and inference depend
on kernel size

Easily parallelizable

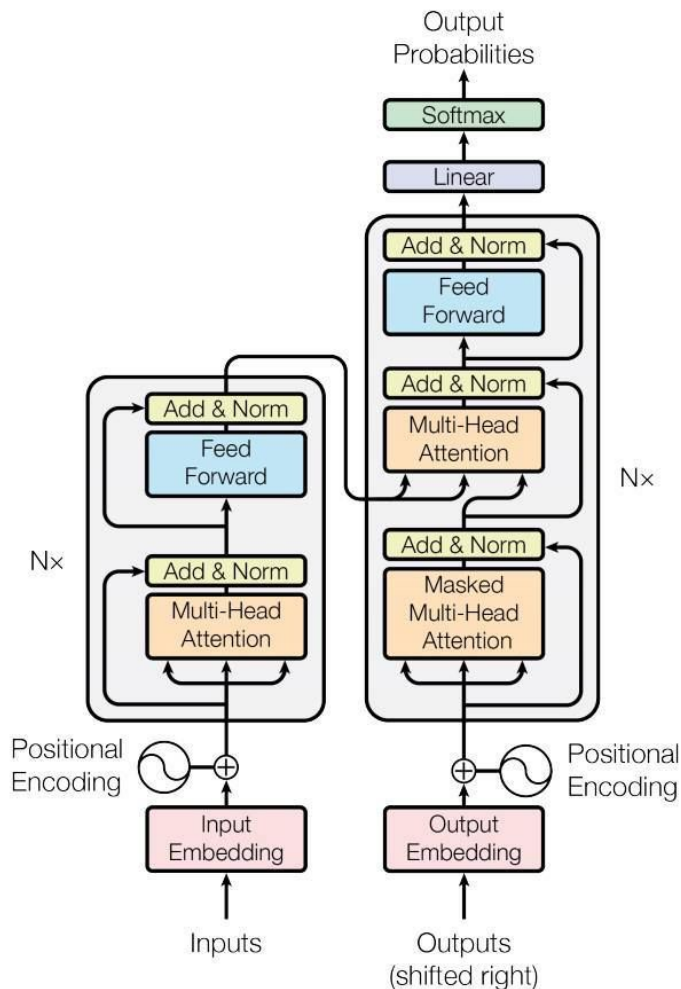


Review: Transformers

Finite context window

Training: $O(N^2)$, easily parallelizable

Inference: $O(N)$ when using KV-Cache,
for each token



Motivation for the new architecture

We want something as efficient as RNNs:

1. constant-time per step inference
2. linear training + parallelizable

... and as good at performance as Transformers

Motivation for the new architecture

We want something as efficient as RNNs:

1. constant-time per step inference
2. linear training + parallelizable

... and as good at performance as Transformers

SSMs is one the directions.

Linear RNNs

Let's get rid of non-linearity, the basic equation of is thus:

$$h_t = Ah_{t-1} + Bx_t$$
$$y_t = Ch_t$$

Linear RNNs

Let's get rid of non-linearity, the basic equation of is thus:

$$h_t = Ah_{t-1} + Bx_t$$

$$y_t = Ch_t$$

- Training Speed: Strong (Parallelizable convolution! – **later**)
- Generation Speed: Strong (constant-time per step)

Linear RNNs

Let's get rid of non-linearity, the basic equation of is thus:

$$h_t = Ah_{t-1} + Bx_t$$

$$y_t = Ch_t$$

- Training Speed: Strong (Parallelizable convolution! – **later**)
- Generation Speed: Strong (constant-time per step)
- **Accuracy: Extremely Poor... Barely learns.**

State space models

Unlike Linear RNNs,

- we look at **continuous-time** values
- instead of modelling hidden state in every position, we use **differential equation** to remodel the change in hidden state over time

$$\begin{aligned}h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) &= \mathbf{C}h(t) + \mathbf{D}x(t)\end{aligned}$$

State space models

A state space model allows us to map an input signal $x(t)$ to an output signal $y(t)$ by means of a state representation $h(t)$ as follows:

$$\begin{aligned}h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) &= \mathbf{C}h(t) + \mathbf{D}x(t)\end{aligned}$$

State space models

A state space model allows us to map an input signal $x(t)$ to an output signal $y(t)$ by means of a state representation $h(t)$ as follows:

$$\begin{aligned}h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) &= \mathbf{C}h(t) + \mathbf{D}x(t)\end{aligned}$$

Note:

$$x(t) \in \mathbb{R} \mapsto y(t) \in \mathbb{R}$$

$$h(t) \in \mathbb{R}^N$$

$$\mathbf{A} \in \mathbb{R}^{N \times N}, \mathbf{B} \in \mathbb{R}^{N \times 1}, \mathbf{C} \in \mathbb{R}^{1 \times N}$$

State space models

A state space model allows us to map an input signal $x(t)$ to an output signal $y(t)$ by means of a state representation $h(t)$ as follows:

$$\begin{aligned}h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) &= \mathbf{C}h(t) + \mathbf{D}x(t)\end{aligned}$$

Note:

$$x(t) \in \mathbb{R} \mapsto y(t) \in \mathbb{R}$$

$$h(t) \in \mathbb{R}^N$$

$$\mathbf{A} \in \mathbb{R}^{N \times N}, \mathbf{B} \in \mathbb{R}^{N \times 1}, \mathbf{C} \in \mathbb{R}^{1 \times N}$$

\mathbf{D} is set to zeros,
represents skip-connection

State space models

A state space model allows us to map an input signal $x(t)$ to an output signal $y(t)$ by means of a state representation $h(t)$ as follows:

$$\begin{aligned}h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) &= \mathbf{C}h(t)\end{aligned}$$

Note:

$$x(t) \in \mathbb{R} \mapsto y(t) \in \mathbb{R}$$

$$h(t) \in \mathbb{R}^N$$

$$\mathbf{A} \in \mathbb{R}^{N \times N}, \mathbf{B} \in \mathbb{R}^{N \times 1}, \mathbf{C} \in \mathbb{R}^{1 \times N}$$

\mathbf{D} is set to zeros,
represents skip-connection

State space models

A state space model allows us to map an input signal $x(t)$ to an output signal $y(t)$ by means of a state representation $h(t)$ as follows:

$$\begin{aligned}h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) &= \mathbf{C}h(t)\end{aligned}$$

This state space model is linear and time invariant.

State space models

A state space model allows us to map an input signal $x(t)$ to an output signal $y(t)$ by means of a state representation $h(t)$ as follows:

$$\begin{aligned}h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) &= \mathbf{C}h(t)\end{aligned}$$

This state space model is linear and time invariant.

To find the output signal $y(t)$ at time t , we first need to find a function $h(t)$ that describes the state of the system for all time steps.

We first need to discretize our system

Simple Discretization Example

Discretize state space model, so that we can calculate the evolution of the state over time by using a **recurrent formulation**:

Simple Discretization Example

Discretize state space model, so that we can calculate the evolution of the state over time by using a **recurrent formulation**:

1. Using the definition of derivative (fix small Δ): $h(t + \Delta) \approx \Delta h'(t) + h(t)$
2. This is the (continuous) state space model: $h'(t) = Ah(t) + Bx(t)$

Simple Discretization Example

Discretize state space model, so that we can calculate the evolution of the state over time by using a **recurrent formulation**:

1. Using the definition of derivative (fix small Δ): $h(t + \Delta) \cong \Delta h'(t) + h(t)$
2. This is the (continuous) state space model: $h'(t) = \mathbf{A}h(t) + \mathbf{B}x(t)$
3. Substitute the state space model into the first expression:

$$\begin{aligned}h(t + \Delta) &\cong \Delta(\mathbf{A}h(t) + \mathbf{B}x(t)) + h(t) \\&= \Delta\mathbf{A}h(t) + \Delta\mathbf{B}x(t) + h(t) \\&= (\mathbf{I} + \Delta\mathbf{A})h(t) + \Delta\mathbf{B}x(t) \\&= \bar{\mathbf{A}}h(t) + \bar{\mathbf{B}}x(t)\end{aligned}$$

Mamba Discretization

The matrices \bar{A} and \bar{B} are the discretized parameters of the model.

This allows us to also calculate the output $y(t)$ of the system for discrete time steps.

$$h'(t) = Ah(t) + Bx(t) \quad (1a)$$

$$h_t = \bar{A}h_{t-1} + \bar{B}x_t \quad (2a)$$

$$y(t) = Ch(t) \quad (1b)$$

$$y_t = Ch_t \quad (2b)$$

Mamba Discretization

The matrices $\bar{\mathbf{A}}$ and $\bar{\mathbf{B}}$ are the discretized parameters of the model.

This allows us to also calculate the output $y(t)$ of the system for discrete time steps.

$$h'(t) = \mathbf{A}h(t) + \mathbf{B}x(t) \quad (1a) \qquad h_t = \bar{\mathbf{A}}h_{t-1} + \bar{\mathbf{B}}x_t \quad (2a)$$

$$y(t) = \mathbf{C}h(t) \quad (1b) \qquad y_t = \mathbf{C}h_t \quad (2b)$$

For **Mamba**, **Zero-Order Hold (ZOH)** rule is used to discretize the system:

Discretization. The first stage transforms the “continuous parameters” $(\Delta, \mathbf{A}, \mathbf{B})$ to “discrete parameters” $(\bar{\mathbf{A}}, \bar{\mathbf{B}})$ through fixed formulas $\bar{\mathbf{A}} = f_A(\Delta, \mathbf{A})$ and $\bar{\mathbf{B}} = f_B(\Delta, \mathbf{A}, \mathbf{B})$, where the pair (f_A, f_B) is called a discretization rule. Various rules can be used such as the zero-order hold (ZOH) defined in equation (4).

$$\bar{\mathbf{A}} = \exp(\Delta \mathbf{A}) \quad \bar{\mathbf{B}} = (\Delta \mathbf{A})^{-1}(\exp(\Delta \mathbf{A}) - \mathbf{I}) \cdot \Delta \mathbf{B} \quad (4)$$

Mamba Discretization

$$h'(t) = Ah(t) + Bx(t) \quad (1a)$$

$$y(t) = Ch(t) \quad (1b)$$

$$h_t = \overline{A}h_{t-1} + \overline{B}x_t \quad (2a)$$

$$y_t = Ch_t \quad (2b)$$

For **Mamba**, **Zero-Order Hold (ZOH)** rule is used to discretize the system:

Discretization. The first stage transforms the “continuous parameters” (Δ, A, B) to “discrete parameters” $(\overline{A}, \overline{B})$ through fixed formulas $\overline{A} = f_A(\Delta, A)$ and $\overline{B} = f_B(\Delta, A, B)$, where the pair (f_A, f_B) is called a discretization rule. Various rules can be used such as the zero-order hold (ZOH) defined in equation (4).

$$\overline{A} = \exp(\Delta A) \quad \overline{B} = (\Delta A)^{-1}(\exp(\Delta A) - I) \cdot \Delta B \quad (4)$$

Note: Δ is learnable parameter

Recurrent computation mode

Discretization leads to recurrent formulation.

$$h_t = \bar{A}h_{t-1} + \bar{B}x_t$$

$$y_t = Ch_t$$

Suppose that, for simplicity, the initial state of the system is 0.

$$h_0 = \bar{B}x_0$$

$$y_0 = Ch_0$$

$$h_1 = \bar{A}h_0 + \bar{B}x_1$$

$$y_1 = Ch_1$$

$$h_2 = \bar{A}h_1 + \bar{B}x_2$$

$$y_2 = Ch_2$$

Recurrent computation: the problem

The recurrent formulation is great for inference: $O(1)$

The recurrent formulation is not good for training: $O(N)$

We want to parallelize the computation as much as possible, just like the Transformer does.

Recurrent computation: the problem

The recurrent formulation is great for inference: $O(1)$

The recurrent formulation is not good for training: $O(N)$

We want to parallelize the computation as much as possible, just like the Transformer does.

Thankfully, State Space Models provide a convolutional mode as well, let's see how it works!

Convolutional computation mode

$$h_0 = \bar{B}x_0$$

$$y_0 = Ch_0 = C\bar{B}x_0$$

$$h_t = \bar{A}h_{t-1} + \bar{B}x_t$$

$$y_t = Ch_t$$

Convolutional computation mode

$$h_0 = \bar{B}x_0$$

$$y_0 = Ch_0 = C\bar{B}x_0$$

$$h_1 = \bar{A}h_0 + \bar{B}x_1 = \bar{A}\bar{B}x_0 + \bar{B}x_1$$

$$y_1 = Ch_1 = C(\bar{A}\bar{B}x_0 + \bar{B}x_1) = C\bar{A}\bar{B}x_0 + C\bar{B}x_1$$

$$h_t = \bar{A}h_{t-1} + \bar{B}x_t$$

$$y_t = Ch_t$$

Convolutional computation mode

$$h_0 = \bar{B}x_0$$

$$y_0 = Ch_0 = C\bar{B}x_0$$

$$h_t = \bar{A}h_{t-1} + \bar{B}x_t$$

$$y_t = Ch_t$$

$$h_1 = \bar{A}h_0 + \bar{B}x_1 = \bar{A}\bar{B}x_0 + \bar{B}x_1$$

$$y_1 = Ch_1 = C(\bar{A}\bar{B}x_0 + \bar{B}x_1) = C\bar{A}\bar{B}x_0 + C\bar{B}x_1$$

$$h_2 = \bar{A}h_1 + \bar{B}x_2 = \bar{A}(\bar{A}\bar{B}x_0 + \bar{B}x_1) + \bar{B}x_2 = \bar{A}^2\bar{B}x_0 + \bar{A}\bar{B}x_1 + \bar{B}x_2$$

$$y_2 = Ch_2 = C(\bar{A}^2\bar{B}x_0 + \bar{A}\bar{B}x_1 + \bar{B}x_2) = C\bar{A}^2\bar{B}x_0 + C\bar{A}\bar{B}x_1 + C\bar{B}x_2$$

Convolutional computation mode

$$h_0 = \bar{B}x_0$$

$$y_0 = Ch_0 = C\bar{B}x_0$$

$$h_t = \bar{A}h_{t-1} + \bar{B}x_t$$

$$y_t = Ch_t$$

$$h_1 = \bar{A}h_0 + \bar{B}x_1 = \bar{A}\bar{B}x_0 + \bar{B}x_1$$

$$y_1 = Ch_1 = C(\bar{A}\bar{B}x_0 + \bar{B}x_1) = C\bar{A}\bar{B}x_0 + C\bar{B}x_1$$

$$h_2 = \bar{A}h_1 + \bar{B}x_2 = \bar{A}(\bar{A}\bar{B}x_0 + \bar{B}x_1) + \bar{B}x_2 = \bar{A}^2\bar{B}x_0 + \bar{A}\bar{B}x_1 + \bar{B}x_2$$

$$y_2 = Ch_2 = C(\bar{A}^2\bar{B}x_0 + \bar{A}\bar{B}x_1 + \bar{B}x_2) = C\bar{A}^2\bar{B}x_0 + C\bar{A}\bar{B}x_1 + C\bar{B}x_2$$

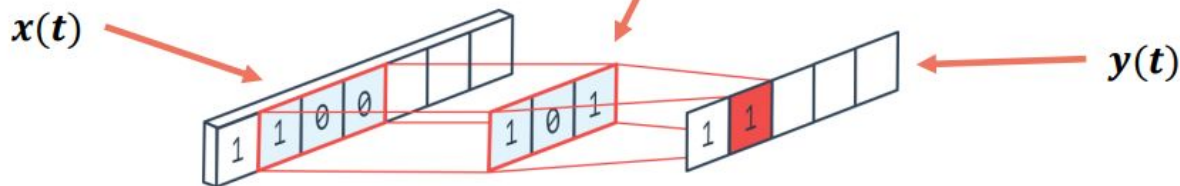
$$y_k = C\bar{A}^k\bar{B}x_0 + C\bar{A}^{k-1}\bar{B}x_1 + \cdots + C\bar{A}\bar{B}x_{k-1} + C\bar{B}x_k$$

Convolutional computation mode

The output of the system can be calculated using a convolution of a kernel \bar{K} with the input $x(t)$.

$$y_k = C\bar{A}^k\bar{B}x_0 + C\bar{A}^{k-1}\bar{B}x_1 + \dots + C\bar{A}\bar{B}x_{k-1} + C\bar{B}x_k \quad (3a)$$

$$y = x * \bar{K} \quad (3b)$$

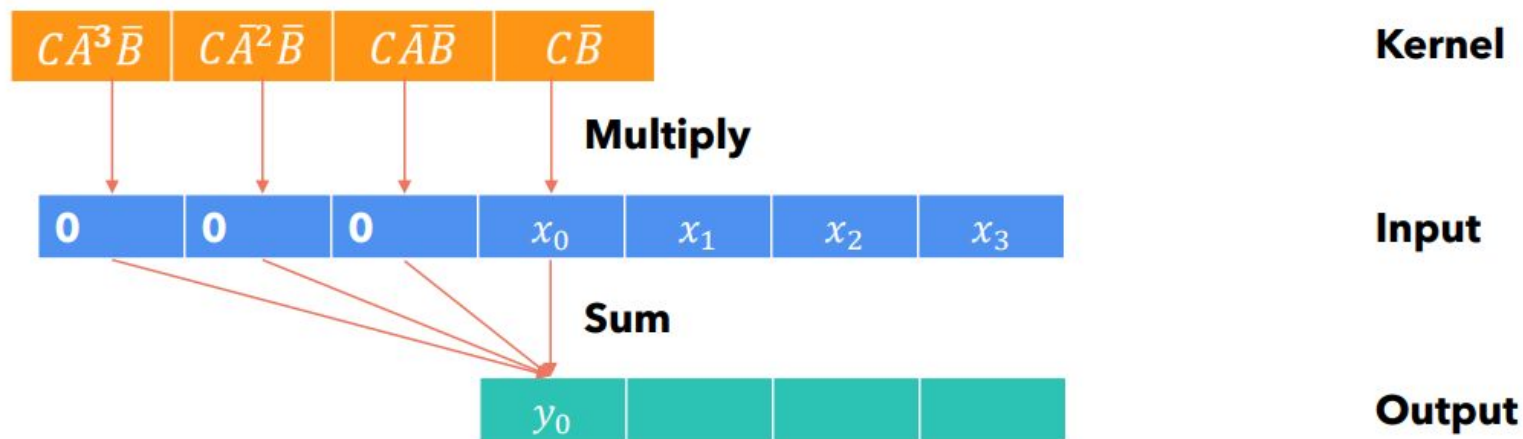


Convolutional computation: step 0

$$y_k = C\bar{A}^k\bar{B}x_0 + C\bar{A}^{k-1}\bar{B}x_1 + \dots + C\bar{A}\bar{B}x_{k-1} + C\bar{B}x_k$$

$$\bar{\mathbf{K}} = (C\bar{B}, C\bar{A}\bar{B}, \dots, C\bar{A}^k\bar{B}, \dots) \quad (3a)$$

$$y = x * \bar{\mathbf{K}} \quad (3b)$$



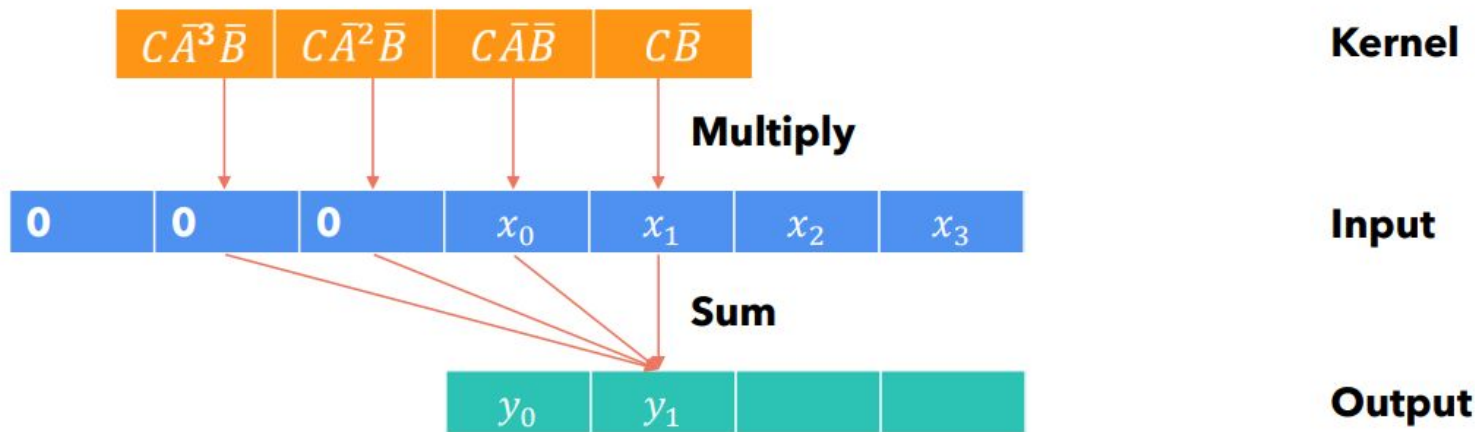
$$y_0 = C\bar{B}x_0$$

Convolutional computation: step 1

$$y_k = C\bar{A}^k\bar{B}x_0 + C\bar{A}^{k-1}\bar{B}x_1 + \dots + C\bar{A}\bar{B}x_{k-1} + C\bar{B}x_k$$

$$\bar{K} = (C\bar{B}, C\bar{A}\bar{B}, \dots, C\bar{A}^k\bar{B}, \dots) \quad (3a)$$

$$y = x * \bar{K} \quad (3b)$$



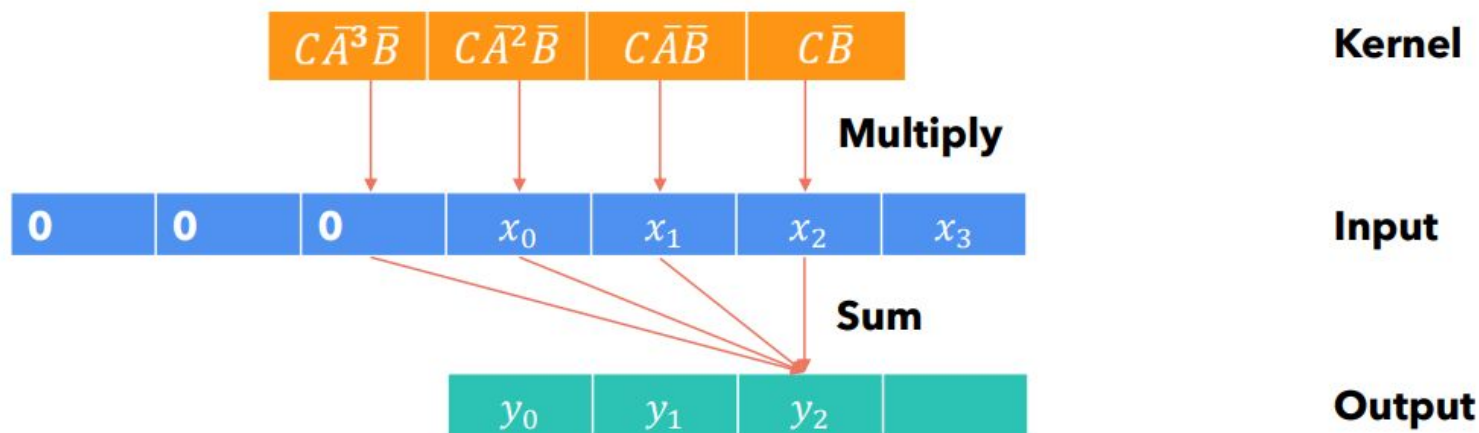
$$y_1 = C\bar{A}\bar{B}x_0 + C\bar{B}x_1$$

Convolutional computation: step 2

$$y_k = C\bar{A}^k\bar{B}x_0 + C\bar{A}^{k-1}\bar{B}x_1 + \dots + C\bar{A}\bar{B}x_{k-1} + C\bar{B}x_k$$

$$\bar{K} = (C\bar{B}, C\bar{A}\bar{B}, \dots, C\bar{A}^k\bar{B}, \dots) \quad (3a)$$

$$y = x * \bar{K} \quad (3b)$$



$$y_2 = C\bar{A}^2\bar{B}x_0 + C\bar{A}\bar{B}x_1 + C\bar{B}x_2$$

Recurrent/Convolutional mode

The best thing about the convolutional calculation is that it can be parallelized.

Recurrent/Convolutional mode

The best thing about the convolutional calculation is that it can be parallelized.

However, **building the kernel can be computationally expensive**, also from a memory point of view.

Recurrent/Convolutional mode

The best thing about the convolutional calculation is that it can be parallelized.

However, **building the kernel can be computationally expensive**, also from a memory point of view.

1. **Convolutional** mode can be used to perform **training**: we already have all the input sequence of tokens, and it can be easily parallelized.
2. **Recurrent** mode can be used to perform **inference**: one token at a time, using a constant amount of computation and memory for each token.

Structure and Dimensions

$$x(t) \in \mathbb{R} \mapsto y(t) \in \mathbb{R}$$

$$h(t) \in \mathbb{R}^N$$

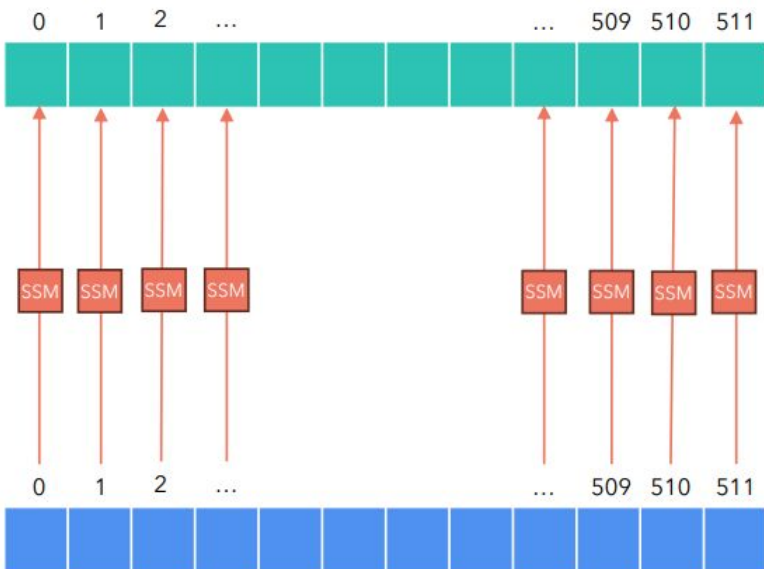
$$\mathbf{A} \in \mathbb{R}^{N \times N}, \mathbf{B} \in \mathbb{R}^{N \times 1}, \mathbf{C} \in \mathbb{R}^{1 \times N}$$

Finally, we note that structured SSMs are so named because computing them efficiently also requires imposing structure on the **A** matrix. The most popular form of structure is **diagonal** which we also use.

In this case, the **matrices can all be represented by N numbers.**

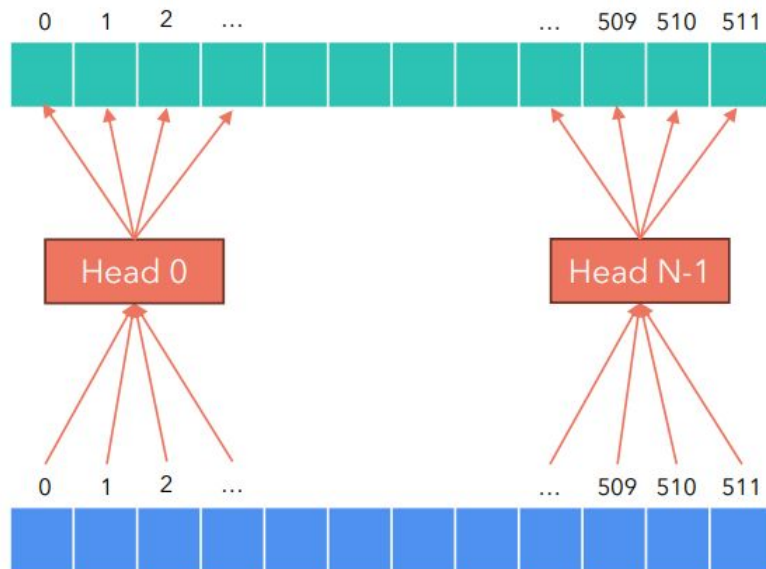
Multiple Dimensions

State Space Model: Each dimension is managed by an independent state space model!



Transformer: Each group of dimensions is managed by a different head of the multi-head attention!

Output



Input

SSM Algorithm

Algorithm 1 SSM (S4)

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

SSM Architectures

Recent research has made significant progress.

S4 [Gu et al., 2022a]

DSS [Gupta, 2022]

GSS [Mehta et al., 2022]

S4D [Gu et al., 2022b]

H3 [Dao et al., 2022]

S5 [Smith et al., 2022]

BiGS [Wang et al., 2022]

QRNN [Bradbury et al., 2016]

Mega [Ma et al., 2022]

SGConv [Li et al., 2022]

Hyena [Poli et al., 2023]

LRU [Orvieto et al., 2023]

RWKV [Peng et al., 2023]

MultiRes [Shi et al., 2023]

Mamba authors

Tri Dao



Tri Dao ✓

@tri_dao

Albert Gu



Albert Gu

@_albertgu

[FlashAttention](#) [Dao et al., 2022]

SSMs research:

[H3](#) [Dao et al., 2022]

[Hyena](#) [Poli et al., 2023]

SSMs research:

[S4](#) [Gu et al., 2022a]

[DSS](#) [Gupta, 2022]

[S4D](#) [Gu et al., 2022b]

[BiGS](#) [Wang et al., 2022]

[LRU](#) [Orvieto et al., 2023]

Sources

1. Mamba: Linear-Time Sequence Modeling with Selective State Spaces – <https://arxiv.org/abs/2312.00752>
2. Mamba and S4 Explained: Architecture, Parallel Scan, Kernel Fusion, Recurrent, Convolution, Math – [Youtube](#), [Slides](#)
3. Do we need Attention? - Linear RNNs and State Space Models (SSMs) for NLP – [Youtube](#), [Slides](#)
4. RNNs strike back – <https://adrian-valente.github.io/2023/10/03/linear-rnns.html#rnns-the-new-generation>