
State Space Models: Nuts and Bolts

Abhishek Aich
NEC Laboratories, America

Abstract

We present a theoretical support of State Space Models, in particular focusing on Linear State Space Layers, introduced in the paper ‘*Combining Recurrent, Convolutional, and Continuous-time Models with Linear State-Space Layers*’ [1]. This report disentangles the underlying concept of the objective function of these models to provide a comprehensive overview. Suggestions/corrections are welcome.

1 Introduction

State space models (SSMs) [1] have been proposed as a possible alternative to (now) ubiquitous transformers. SSMs extended the idea proposed in [2] that capturing long range dependencies is equivalent to optimally storing and incorporating information from previous inputs in *memory*. Transformers¹ were known to have difficulty scaling to long sequence lengths, primarily due to the quadratic complexity of their self-attention mechanism [3]. Hence, authors of [2] proposed to treat *memory* as a problem of online function approximation where a function (representing past memory) $f(t) : \mathbb{R}^+ \rightarrow \mathbb{R}$ is summarized by storing its optimal coefficients in terms of some basis functions. Following on [2] idea, [1] proposed *Linear State Space Layers* (LSSLs) for sequence modeling based on classical standard state space representations. Readers are recommended to go through this extremely well summarized video [4].

This report is dedicated to deciphering the minute details of LSSLs starting from the classical state space representations as the basics. Please do not hesitate to contact the author in case you find an error in the analysis.

2 The Classical (Continuous) "State Space Model"

In practice when we analyze any system, we start by developing a mathematical description of the system in the form of differential equations [5]. This is because differential equations provide a powerful and flexible framework for modeling, analyzing, and predicting the behavior of such systems.

A “state space” representation of a system is the mathematical model where the inputs, outputs and the system’s internal components (that mathematically describe the system, also called as “states”) [6] are related via *first-order differential equations*. Although there are other mathematical models [6, 7] for establishing the relationship between the aforementioned trio (e.g. using transfer functions), state space models leverage first-order differential equations for their conciseness (isolate each components or states and provide a representation) as well as the ease of efficient analysis (e.g. using linear algebra techniques for better system explainability). State-space models are one such method to convert a n -th order differential equation to *first-order matrix* differential equation [5] for their aforementioned useful properties.

Let’s discuss the origins of the state space representation of a system. Denote an input and output (both function of time t) with $u(t)$ and $y(t)$, respectively. Consider a n -th order system described by

¹this is in 2020.

a n -th order differential equation [6]

$$\frac{d^n \mathbf{y}(t)}{dt^n} + a_{n-1} \frac{d^{n-1} \mathbf{y}(t)}{dt^{n-1}} + \dots + a_1 \frac{d\mathbf{y}(t)}{dt} + a_0 \mathbf{y}(t) = b_0 \mathbf{u}(t) \quad (1)$$

Here, $a_i \forall i$ and b_0 represent the corresponding scalar coefficients. To understand how state-space models convert Eq. 1 into a batch of first-order differential equations, let's look at an example of a system defined by third-order differential equations (i.e. putting $n = 3$ in Eq. 1). We will drop the time function representation of all variables for simplicity.

$$\frac{d^3 \mathbf{y}}{dt^3} + a_2 \frac{d^2 \mathbf{y}}{dt^2} + a_1 \frac{d\mathbf{y}}{dt} + a_0 \mathbf{y} = b_0 \mathbf{u} \quad (2)$$

To convert this third-order differential equation to a set of first-order differential equations, we introduce the following three variables (all functions of time t) as follows.

$$\mathbf{h}_1 = \mathbf{y}, \mathbf{h}_2 = \frac{d\mathbf{y}}{dt}, \mathbf{h}_3 = \frac{d^2 \mathbf{y}}{dt^2} \quad (3)$$

For \mathbf{h}_2 and \mathbf{h}_3 , we can also observe that

$$\mathbf{h}_2 = \frac{d\mathbf{h}_1}{dt}, \mathbf{h}_3 = \frac{d\mathbf{h}_2}{dt}, \quad \text{along with, } \frac{d\mathbf{h}_3}{dt} = \frac{d^3 \mathbf{y}}{dt^3}$$

Using above in Eq. 2, we get

$$\frac{d\mathbf{h}_3}{dt} + a_2 \mathbf{h}_3 + a_1 \mathbf{h}_2 + a_0 \mathbf{h}_1 = b_0 \mathbf{u} \quad (4)$$

We summarize the obtained first-order differential equations in Eq. 3, needed to describe Eq. 1, as follows.

$$\frac{d\mathbf{h}_1}{dt} = \mathbf{h}_2, \frac{d\mathbf{h}_2}{dt} = \mathbf{h}_3, \frac{d\mathbf{h}_3}{dt} = -(a_2 \mathbf{h}_3 + a_1 \mathbf{h}_2 + a_0 \mathbf{h}_1) + b_0 \mathbf{u} \quad (5)$$

Writing the above equations together gives us the following matrix form.

$$\begin{bmatrix} \frac{d\mathbf{h}_1}{dt} \\ \frac{d\mathbf{h}_2}{dt} \\ \frac{d\mathbf{h}_3}{dt} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -a_0 & -a_1 & -a_2 \end{bmatrix} \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \\ \mathbf{h}_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ b_0 \end{bmatrix} \mathbf{u} \quad (6)$$

Setting $\mathbf{h} = \begin{bmatrix} \frac{d\mathbf{h}_1}{dt} \\ \frac{d\mathbf{h}_2}{dt} \\ \frac{d\mathbf{h}_3}{dt} \end{bmatrix}$, $\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -a_0 & -a_1 & -a_2 \end{bmatrix}$, $\mathbf{B} = \begin{bmatrix} 0 \\ 0 \\ b_0 \end{bmatrix}$, we get the *state* equations. After

computing the solution to this state equation, we can set $\mathbf{C} = [1 \ 0 \ 0]$ and arrive at the *output* equation.

$$\mathbf{y} = [1 \ 0 \ 0] \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \\ \mathbf{h}_3 \end{bmatrix} \quad (7)$$

Hence, the final SSM representation is given as

$$\dot{\mathbf{h}} = \mathbf{A}\mathbf{h} + \mathbf{B}\mathbf{u} \quad (8a)$$

$$\mathbf{y} = \mathbf{C}\mathbf{h} \quad (8b)$$

In general, for n -th order differential equation can be represented by (\mathbf{D} turned out to be zero for above example, but it's usually not the case)

$$\begin{aligned} \dot{\mathbf{h}}(t) &= \mathbf{A}\mathbf{h}(t) + \mathbf{B}\mathbf{u}(t) \\ \mathbf{y}(t) &= \mathbf{C}\mathbf{h}(t) + \mathbf{D}\mathbf{u}(t) \end{aligned} \quad (9)$$

In control system literature [7], \mathbf{A} is the *state* matrix, \mathbf{B} is the *input* matrix, \mathbf{C} is the *output* matrix, and \mathbf{D} is the *feedthrough* matrix. Further, we need to compute the solution of the state equation for using the output equation. We visualize the total representation in Fig. 1. Next, let's relate these fundamental SSM equations Eq. 9 to Linear State-Space Layers (LSSLs) introduced in [1]. For this, we will need to convert Eq. 9 to its discrete form.

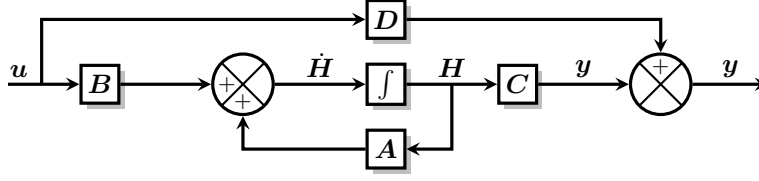


Figure 1: Classical SSMs

3 The Classical (Discrete) "State Space Model"

We aim to understand how [1] used SSMs (Eq. 9) as a black-box representation in a deep sequence model, where A , B , C , and D are parameters learned by gradient descent. Similar to [8], we will assume $D = 0$, interpreting it as an easily computable skip connection (can also be observed in Fig. 1).

To do this, we first convert the continuous-time equation group into discrete-time LSSLs (so that we can process discrete inputs such as image pixels). Essentially, we are estimating $\mathbf{y}(t)$ from $\mathbf{u}(t)$ using the system hidden state $\mathbf{h}(t)$. In particular, we leverage Generalized Bilinear Transform (GBT) [9]. GBT of Eq. 9 help us apply SSMs to operate on discrete input sequence $(\mathbf{u}_0, \mathbf{u}_1, \dots)$ instead of continuous function $\mathbf{u}(t)$. The inputs \mathbf{u}_k can be viewed as sampling the continuous signal $\mathbf{u}(t)$. We assume that the input $\mathbf{u}(t)$ is constant over the time interval Δ (also known as Zero-Order Hold assumption [10]), which is represented as follows.

$$\mathbf{u}_k = \mathbf{u}(k\Delta) := \mathbf{u}[k] \quad \forall \quad k\Delta \leq t \leq (k+1)\Delta \quad (10)$$

Step size Δ represents the gap between time points (i.e. time-period *or* time-period-scale *or* time-scale [1]) from where \mathbf{u}_k is sampled from $\mathbf{u}(t)$. This means that the length of dependencies (within the input sequence) that a (discrete) SSM can capture, is dependent on this value. This makes Δ an important factor in deep sequence modeling. So, let's first derive the discrete versions of Eq. 9 before computing the discrete versions of A , B , and C (using GBT).

There are two ways [10] to convert Eq. 9 to discrete form. One is an *Approximate discretisation* (which uses the Euler's approximation) and the other one, which we will focus on, is the *Exact discretisation*. The exact discretisation requires us first to calculate the solution to Eq. 9, in particular, solving for $\dot{\mathbf{h}}(t) = A\mathbf{h}(t) + B\mathbf{u}(t)$. For clarity, we will replace the variable t with τ while computing the solution of this state equation (as t will appear again later).

It is known that for a scalar equation (where $\mathbf{h}(\tau) \in \mathbb{R}$ and a is some constant), $\dot{\mathbf{h}}(\tau) = a\mathbf{h}(\tau)$ has a solution of the form $\mathbf{h}(\tau) = e^{A\tau}\mathbf{h}(0)$. Thus, in our case, where A is a matrix, we can start with the assumption that the solution is of the form $\mathbf{h}(\tau) = e^{A\tau}\mathbf{h}(\tau)$ [10]. With this, let's multiply $e^{-A\tau}$ on both sides of $\dot{\mathbf{h}}(\tau) = A\mathbf{h}(\tau) + B\mathbf{u}(\tau)$ giving us the following.

$$e^{-A\tau}\dot{\mathbf{h}}(\tau) = e^{-A\tau}(A\mathbf{h}(\tau) + B\mathbf{u}(\tau)) \quad (11)$$

$$\implies e^{-A\tau}\dot{\mathbf{h}}(\tau) - e^{-A\tau}A\mathbf{h}(\tau) = e^{-A\tau}B\mathbf{u}(\tau) \quad (12)$$

Noting that $\frac{d(e^{-A\tau}\mathbf{h}(\tau))}{d\tau} = e^{-A\tau}\dot{\mathbf{h}}(\tau) - e^{-A\tau}A\mathbf{h}(\tau)$, we get

$$\frac{d(e^{-A\tau}\mathbf{h}(\tau))}{d\tau} = e^{-A\tau}B\mathbf{u}(\tau) \quad (13)$$

Integrating both sides of Eq. 13 between 0 to t gives us

$$e^{-At}\mathbf{h}(t) \Big|_{\tau=0}^{\tau=t} = \int_{\tau=0}^{\tau=t} e^{-A\tau}B\mathbf{u}(\tau)d\tau \quad (14)$$

$$\implies e^{-At}\mathbf{h}(t) - e^0A\mathbf{h}(0) = \int_{\tau=0}^{\tau=t} e^{-A\tau}B\mathbf{u}(\tau)d\tau \quad (15)$$

Multiplying e^{At} to both sides of Eq. 15 and noting $e^0 = \mathbf{I}$ (identity matrix), we get the solution of the state equation in Eq. 9 as follows

$$\mathbf{h}(t) = e^{At}\mathbf{h}(0) + \int_{\tau=0}^{\tau=t} e^{A(t-\tau)}B\mathbf{u}(\tau)d\tau \quad (16)$$

Inserting this solution Eq. 16 in the output equation of Eq. 9, we get the output prediction as follows.

$$\mathbf{y}(t) = \mathbf{C} \left(e^{\mathbf{A}t} \mathbf{h}(0) + \int_{\tau=0}^{\tau=t} e^{\mathbf{A}(t-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau \right) + \mathbf{D} \mathbf{u}(t) \quad (17)$$

Now that we have computed the solution in Eq. 16, we can now derive the discrete versions of Eq. 9. With discrete input defined as Eq. 10, we evaluate the state equation at the sampling instant $t = (k+1)\Delta$ as follows [10].

$$\mathbf{h}((k+1)\Delta) = e^{\mathbf{A}(k+1)\Delta} \mathbf{h}(0) + \int_{\tau=0}^{\tau=(k+1)\Delta} e^{\mathbf{A}((k+1)\Delta-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau \quad (18)$$

$$\implies \mathbf{h}((k+1)\Delta) = e^{\mathbf{A}\Delta} \left(e^{\mathbf{A}k\Delta} \mathbf{h}(0) \right) + \underbrace{\left(\int_{\tau=0}^{\tau=(k+1)\Delta} e^{\mathbf{A}((k+1)\Delta-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau \right)}_{\mathcal{I}} \quad (19)$$

Let's focus on the term \mathcal{I} . We can break this term in two intervals as follows.

$$\mathcal{I} = \underbrace{\int_{\tau=0}^{\tau=k\Delta} e^{\mathbf{A}((k+1)\Delta-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau}_{\mathcal{I}.1} + \underbrace{\int_{\tau=k\Delta}^{\tau=(k+1)\Delta} e^{\mathbf{A}((k+1)\Delta-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau}_{\mathcal{I}.2} \quad (20)$$

Taking $e^{\mathbf{A}\Delta}$ outside the integral in $\mathcal{I}.1$, and noting that $\mathbf{u}(\tau) = \mathbf{u}(k\Delta)$ in the interval $k\Delta \leq \tau \leq (k+1)\Delta$ for $\mathcal{I}.2$, we get

$$\mathcal{I} = e^{\mathbf{A}\Delta} \left(\int_{\tau=0}^{\tau=k\Delta} e^{\mathbf{A}(k\Delta-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau \right) + \left(\int_{\tau=k\Delta}^{\tau=(k+1)\Delta} e^{\mathbf{A}((k+1)\Delta-\tau)} \mathbf{B} d\tau \right) \mathbf{u}(k\Delta) \quad (21)$$

Putting Eq. 21 back in Eq. 19, we get

$$\begin{aligned} \mathbf{h}((k+1)\Delta) &= e^{\mathbf{A}\Delta} \left(e^{\mathbf{A}k\Delta} \mathbf{h}(0) \right) + e^{\mathbf{A}\Delta} \left(\int_{\tau=0}^{\tau=k\Delta} e^{\mathbf{A}(k\Delta-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau \right) \\ &\quad + \left(\int_{\tau=k\Delta}^{\tau=(k+1)\Delta} e^{\mathbf{A}((k+1)\Delta-\tau)} \mathbf{B} d\tau \right) \mathbf{u}(k\Delta) \quad (22) \end{aligned}$$

$$\begin{aligned} \implies \mathbf{h}((k+1)\Delta) &= e^{\mathbf{A}\Delta} \left(\underbrace{e^{\mathbf{A}k\Delta} \mathbf{h}(0) + \int_{\tau=0}^{\tau=k\Delta} e^{\mathbf{A}(k\Delta-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau}_{\mathcal{II}} \right) \\ &\quad + \underbrace{\left(\int_{\tau=k\Delta}^{\tau=(k+1)\Delta} e^{\mathbf{A}((k+1)\Delta-\tau)} \mathbf{B} d\tau \right)}_{\mathcal{III}} \mathbf{u}(k\Delta) \quad (23) \end{aligned}$$

The term \mathcal{II} is clearly equivalent to $\mathbf{h}(k\Delta)$ if we set $t = k\Delta$ in Eq. 16. The term \mathcal{III} can be further simplified by setting $\epsilon = (k+1)\Delta - \tau$. This also means $d\tau = -d\epsilon$.

$$\mathcal{III} = \int_{\epsilon=\Delta}^{\epsilon=0} e^{\mathbf{A}(\epsilon)} \mathbf{B} (-d\epsilon) = \int_{\epsilon=0}^{\epsilon=\Delta} e^{\mathbf{A}\epsilon} \mathbf{B} d\epsilon \quad (24)$$

If we further simplify \mathcal{III} to remove the integral, we can use the following steps [11]. Assuming \mathbf{A} is invertible, we can write $\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}$ where \mathbf{I} is the identity matrix. Then, the term \mathcal{III} can be written as follows.

$$\mathcal{III} = \mathbf{A}^{-1} \left(\int_{\epsilon=0}^{\epsilon=\Delta} \mathbf{A} e^{\mathbf{A}\epsilon} d\epsilon \right) \mathbf{B} \quad (25)$$

With $\frac{d(e^{\mathbf{A}\epsilon})}{d\epsilon} = \mathbf{A} e^{\mathbf{A}\epsilon}$,

$$\mathcal{III} = \mathbf{A}^{-1} \left(\int_{\epsilon=0}^{\epsilon=\Delta} \frac{d(e^{\mathbf{A}\epsilon})}{d\epsilon} d\epsilon \right) \mathbf{B} = \mathbf{A}^{-1} \left(e^{\mathbf{A}\epsilon} \Big|_{\epsilon=0}^{\epsilon=\Delta} \right) \mathbf{B} = \mathbf{A}^{-1} (e^{\mathbf{A}\Delta} - \mathbf{I}) \mathbf{B} \quad (26)$$

Putting \mathbb{I} and \mathbb{III} back in Eq. 22, we obtain the discretized state space equation.

$$\mathbf{h}((k+1)\Delta) = \left(e^{A\Delta}\right)\mathbf{h}(k\Delta) + \left(\mathbf{A}^{-1}(e^{A\Delta} - \mathbf{I})\mathbf{B}\right)\mathbf{u}(k\Delta) \quad (27)$$

Similar to computing the solution (Eq. 16) of continuous-time state equation in Eq. 9, we now compute the solution of Eq. 27 using the Iterative Method. Let $\bar{\mathbf{A}} = e^{A\Delta}$, $\bar{\mathbf{B}} = \mathbf{A}^{-1}(e^{A\Delta} - \mathbf{I})\mathbf{B}$. Start by setting $k = 0$ and $k = 1$ (drop Δ in the notation for brevity).

$$\mathbf{h}(1) = \bar{\mathbf{A}}\mathbf{h}(0) + \bar{\mathbf{B}}\mathbf{u}(0) \quad (28)$$

$$\mathbf{h}(2) = \bar{\mathbf{A}}\mathbf{h}(1) + \bar{\mathbf{B}}\mathbf{u}(1) \quad (29)$$

Let's use Eq. 28 in Eq. 29.

$$\begin{aligned} \mathbf{h}(2) &= \bar{\mathbf{A}}(\bar{\mathbf{A}}\mathbf{h}(0) + \bar{\mathbf{B}}\mathbf{u}(0)) + \bar{\mathbf{B}}\mathbf{u}(1) \\ &= \bar{\mathbf{A}}^2\mathbf{h}(0) + \bar{\mathbf{A}}\bar{\mathbf{B}}\mathbf{u}(0) + \bar{\mathbf{B}}\mathbf{u}(1) \end{aligned} \quad (30)$$

We can keep performing the same steps for $k = 3$, which leads to the following.

$$\mathbf{h}(3) = \bar{\mathbf{A}}^3\mathbf{h}(0) + \bar{\mathbf{A}}^2\bar{\mathbf{B}}\mathbf{u}(0) + \bar{\mathbf{A}}\bar{\mathbf{B}}\mathbf{u}(1) + \bar{\mathbf{B}}\mathbf{u}(2) \quad (31)$$

Clearly, we can generalize this form to the following [10].

$$\mathbf{h}(k) = \bar{\mathbf{A}}^k\mathbf{h}(0) + \sum_{n=0}^{k-1} \bar{\mathbf{A}}^{(k-1-n)}\bar{\mathbf{B}}\mathbf{u}(n) \quad (32)$$

This state solution now can be used with the discretized output equation. To obtain the output equation, we set $t = k\Delta$ in Eq. 17 and noting the occurrence of the term \mathbb{II} .

$$\mathbf{y}(k\Delta) = \mathbf{C} \left(\underbrace{e^{A k \Delta} \mathbf{h}(0) + \int_{\tau=0}^{\tau=k\Delta} e^{A(k\Delta-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau}_{\mathbb{II}} \right) + \mathbf{D} \mathbf{u}(k\Delta) \quad (33)$$

$$\implies \mathbf{y}(k\Delta) = \mathbf{C} \mathbf{h}(k\Delta) + \mathbf{D} \mathbf{u}(k\Delta) \quad (34)$$

To summarize, we have the following discrete SSM of Eq. 9.

$$\begin{aligned} \mathbf{h}_{k+1} &= \bar{\mathbf{A}}\mathbf{h}_k + \bar{\mathbf{B}}\mathbf{u}_k \\ \mathbf{y}_k &= \bar{\mathbf{C}}\mathbf{h}_k + \bar{\mathbf{D}}\mathbf{u}_k \end{aligned} \quad (35)$$

where,

$$\bar{\mathbf{A}} = e^{A\Delta}, \quad \bar{\mathbf{B}} = \mathbf{A}^{-1}(e^{A\Delta} - \mathbf{I})\mathbf{B}, \quad \bar{\mathbf{C}} = \mathbf{C}, \quad \bar{\mathbf{D}} = \mathbf{D} \quad (36)$$

This discrete SSM in Eq. 35 is exactly the same as Eq. (4) in Mamba [12]. With the parameter transformation from $(\Delta, \mathbf{A}, \mathbf{B})$ to $(\bar{\mathbf{A}}, \bar{\mathbf{B}})$, the LSSL based sequence model can be represented in following two ways: either as a linear recurrence or as a global convolution.

4 State Space Sequence Models

Now, that we have arrived at the discrete form of SSMs in Eq. 35 that can handle discrete input sequences \mathbf{u}_k , lets formalize the notations and nomenclature before we discuss SSM's application for deep sequence modeling, following [1].

Given a fixed state space representations $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$, an LSSL is the sequence-to-sequence mapping defined by discretizing the linear state-space model Eq. 9 (which we derived in Eq. 35). Eq. 35 can be directly interpreted as a Recurrent Neural Network (RNN) representation. In particular, Eq. 35 maps the sequence \mathbf{u}_k to \mathbf{y}_k using the hidden state in \mathbf{h}_k . However to process the (new) sequence \mathbf{u}_{k+1} , we need to recompute this hidden state \mathbf{h}_{k+1} recursively using the transition matrix $\bar{\mathbf{A}}$. Hence, we will leverage the *convolutional representation*.

4.1 The Convolutional Representation

Let's go back to the solution of the discrete state equation in Eq. 32.

In Eq. 32, assume $\mathbf{h}(0) = 0$ and write $\mathbf{w}[k-1-n] = \overline{\mathbf{A}}^{(k-1-n)}\overline{\mathbf{B}}$ as follows.

$$\mathbf{h}[k] = \sum_{n=0}^{k-1} \mathbf{w}[k-1-n]\mathbf{u}[n] \quad (37)$$

We can note that Eq. 37 represents the discrete convolution operation between the interval $0 \leq n \leq k-1$. Insert Eq. 37 in output equation of Eq. 35 to get

$$\begin{aligned} \mathbf{y}[k] &= \mathbf{C} \left(\sum_{n=0}^{k-1} \mathbf{w}[k-1-n]\mathbf{u}[n] \right) + \mathbf{D}\mathbf{u}[k] \\ &= \overline{\mathbf{K}} * \mathbf{u}[k] + \mathbf{D}\mathbf{u}[k] \end{aligned} \quad (38)$$

With $\overline{\mathbf{K}}$, the convolution operation can be computed comparatively efficiently than the RNN form [1].

5 Conclusion

We provided an overview of Linear State Space Layers (LSSLs), which are the core components of State Space Models (SSMs). We traced their development from classical continuous-time models to discrete-time formulations, demonstrating their application in deep sequence modeling based on solid mathematical principles. Finally, we discussed their efficient convolutional representation, which makes LSSLs well-suited for modern deep sequence modeling.

References

- [1] Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Combining recurrent, convolutional, and continuous-time models with linear state space layers. *Advances in neural information processing systems*, 34:572–585, 2021.
- [2] Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory with optimal polynomial projections. *Advances in neural information processing systems*, 33: 1474–1487, 2020.
- [3] Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena: A benchmark for efficient transformers. *arXiv preprint arXiv:2011.04006*, 2020.
- [4] Samuel Albanie. Mamba - a replacement for transformers?, 2023. URL https://youtu.be/ouF-H35at0Y?si=mRJfTM_-sTBbdDny. Accessed: September 8, 2024.
- [5] Erik Cheever. Representing linear physical systems with differential equations. Online, 2022. URL <https://lpsa.swarthmore.edu/Representations/SysRepDiffEq.html>. Accessed: 2024-05-30.
- [6] Ramkrishna Pasumarthy and Viswanath Talasila. Introduction to state space systems, 2017. URL <https://youtu.be/xajgSUci9zs?feature=shared>. Accessed: 2024-05-29.
- [7] Richard Bishop and Robert H C Dorf. *Modern control systems*. Pearson, 2011.
- [8] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- [9] Guofeng Zhang, Tongwen Chen, and Xiang Chen. Performance recovery in digital implementation of analogue systems. *SIAM journal on control and optimization*, 45(6):2207–2223, 2007.
- [10] Julio H. Braslavsky. Lecture notes in control system design and management (elec4410), 2006. URL <https://www.eng.newcastle.edu.au/~jhb519/teaching/elec4410/Lectures/Lec11.pdf>.

- [11] Zhi Jane Li. Lecture notes in synergy of human and robotic systems (rbe59), 2017. URL https://users.wpi.edu/~zli11/teaching/rbe595_2017/LectureSlide_PDF/discretization.pdf.
- [12] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [13] Jeffrey L Elman. Finding structure in time. *Cognitive science*, 14(2):179–211, 1990.
- [14] Doron Levy. Lecture notes in introduction to numerical analysis i (amsc 466), 2016. URL <https://math.umd.edu/~dlevy/classes/amsc466/notes.pdf>.
- [15] Jeffrey Wong. Lecture notes in applied computational analysis (math 563), 2020. URL <https://services.math.duke.edu/~jtwong/math563-2020/lectures/Lec4-orth.pdf>.