

# AUTOCORRELATION MATTERS: UNDERSTANDING THE ROLE OF INITIALIZATION SCHEMES FOR STATE SPACE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Current methods for initializing state space model (SSM) parameters primarily rely on the HiPPO framework (Gu et al., 2023), which is based on online function approximation with the SSM kernel basis. However, the HiPPO framework does not explicitly account for the effects of the temporal structures of input sequences on the optimization of SSMs. In this paper, we take a further step to investigate the roles of SSM initialization schemes by considering the autocorrelation of input sequences. Specifically, we: (1) rigorously characterize the dependency of the SSM timescale on sequence length based on sequence autocorrelation; (2) find that with a proper timescale, allowing a zero real part for the eigenvalues of the SSM state matrix mitigates the curse of memory while still maintaining stability at initialization; (3) show that the imaginary part of the eigenvalues of the SSM state matrix determines the conditioning of SSM optimization problems, and uncover an approximation-estimation tradeoff when training SSMs with a specific class of target functions.

## 1 INTRODUCTION

The state space model (SSM) is a sequence model that has recently shown great potential in long sequence modeling across various applications, including computer vision (Zhu et al., 2024; Liu et al., 2024), time series forecasting (Rangapuram et al., 2018; Zhang et al., 2023) and natural language processing (Gu & Dao, 2023; Dao & Gu, 2024). In mathematics, a SSM layer is defined by a continuous-time ordinary differential equation  $h'(t) = Wh(t) + Bx(t)$ ,  $y(t) = Ch(t) + Dx(t)$ , where  $W, B, C, D$  are trainable parameters,  $x(t)$  is the input sequence, and  $y(t)$  is the output sequence. For discrete input sequences, a timescale  $\Delta > 0$  will be introduced as a hyperparameter to discretize the model. Different from the attention mechanism (Vaswani et al., 2017), SSMs are recurrent-based architectures that treat the input sequence token by token, yet can achieve first-order time complexity on the sequence length through parallelization (Gu et al., 2022b). There are two well known issues for training recurrent-based architectures, the vanishing and the exploding gradient problems (Pascanu et al., 2013). By introducing complex-valued initialization schemes, proper parameterization methods and regularization techniques, recent works demonstrate that SSMs can achieve performance comparable to attention-based architectures in terms of both computational cost and sample efficiency (Gu & Dao, 2023; Dao & Gu, 2024; Zhu et al., 2024; Yu et al., 2024; Wang & Li, 2024; Liu & Li, 2024; Yu et al., 2024; Bick et al., 2024; Hwang et al., 2024; Wang et al., 2024; Waleffe et al., 2024). However, the theoretical understanding on the roles of the initialization schemes is still lacking and needs to further explored. In this paper, we particularly look into the timescale  $\Delta$  and the state matrix  $W$ , and we aim to study the following fundamental question

*Given a sequential dataset with length  $L$ , how should the timescale  $\Delta$  depend on  $L$  and what is the role of  $W$  on training SSMs?*

Based on the analysis of continuous-time SSMs, previous works (Gu et al., 2022b;c; 2023) propose the HiPPO framework where  $W, B$  are initialized such that the SSM basis kernels  $\{e_n^\top e^{Wt} B\}_{n=1}^\infty$  are orthogonal in  $L^2[0, \infty)$  with some measure  $\omega(t)$ , and the timescale  $\Delta$  scales as  $1/L$  to capture

054 long range dependencies of sequences with length  $L$ . Common HiPPO-based initialization methods  
 055 such as S4D-Legs and S4D-Lin typically presume that the measure  $\omega(t)$  is exponential decay and  
 056 the discrete input sequences  $x$  have an inherent timescale  $\Delta$  that is shared with the model. However,  
 057 these assumptions are restrictive because exponential decay measures weaken the effects of temporal  
 058 dependencies in input sequences, and in practice, we usually lack prior information about the data’s  
 059 timescale. To address this concern, we take an initial step towards understanding the relationship  
 060 between the autocorrelation of input sequences and the SSM initialization schemes. Specifically,  
 061 we focus on the diagonal SSM<sup>1</sup> (Gu et al., 2022c) where the state matrix  $W$  is a complex-valued  
 062 diagonal matrix. By studying the stability condition for given input sequences  $x \in \mathbb{R}^L$ , we find that  
 063 the connection of the timescale  $\Delta$  and the sequence length  $L$  is highly related with the spectrum of  
 064 the data autocorrelation matrix  $\mathbb{E}[xx^\top]$ . Different temporal dependencies in the input sequences can  
 065 cause significant variations in the spectrum of the autocorrelation matrix. For example, when  $x$  is  
 066 sampled from a standard normal distribution,  $x$  has zero temporal dependencies, and the autocorrelation  
 067 matrix becomes an identity matrix. On the other hand, if  $x$  consists of constant values, the  
 068 input sequence exhibits full temporal dependencies, and the autocorrelation matrix is low rank. For  
 069 the state matrix  $W$ , our stability analysis shows that even with a zero real part, i.e.  $\Re(W) = 0$ , the  
 070 diagonal SSM can still be stable at initialization if  $\Delta$  is properly set. However, during training, it  
 071 is worth noting that stability is not guaranteed because  $\Re(W) = 0$  places the SSM on the edge of  
 072 training stability. In this paper, we find that, at least for simple tasks, initializing  $\Re(W) = 0$  helps  
 073 improve training performance for fixed-length tasks. One benefit for setting the real part to zero is  
 074 that the learned SSM kernel functions at initialization do not exponentially decay, which helps to  
 075 mitigate the curse of memory (Li et al., 2022). Our convergence analysis indicates that the imaginary  
 076 part  $\Im(W)$  plays a crucial role in the convergence rate and explains the benefits for complex-valued  
 077 SSMs compared to real-valued SSMs in terms of the optimization. In particular, the more separated  
 078 the imaginary parts  $\Im(w)$  are, the faster the convergence. When considering both approximation  
 079 and optimization, we characterize an approximation-estimation tradeoff when the target function  
 080 has closely spaced dominant frequencies. Then well separated  $\Im(w)$  values lead to fast conver-  
 081 gence, while achieving a good approximation requires close imaginary parts. To summarize, our  
 082 contributions are as follows:

- 082 • In section 4.1, we characterize the dependency between the timescale  $\Delta$  and the sequence  
 083 length  $L$  by taking into account the autocorrelation of the input sequences. Even if the  
 084 eigenvalues of the state matrix  $W$  have zero real part, the stability condition on the magni-  
 085 tude of the output value at initialization can still hold with an appropriate setting of  $\Delta$ .
- 086 • In section 4.2, we show that the real part of the eigenvalues of the state matrix  $W$  determines  
 087 the decay rate of the SSM kernel functions. Allowing the eigenvalues of  $W$  to have zero  
 088 real part at initialization can significantly increase the model’s effective memory and help  
 089 mitigate the curse of memory for fixed-length tasks that require long-term memory.
- 090 • In section 4.3, we prove that the conditioning of SSM optimization problems is determined  
 091 by the separation distance of the imaginary parts of the eigenvalues of the state matrix.  
 092 Well-separated imaginary parts induce faster convergence, whereas closely spaced ones  
 093 lead to slower convergence. This explains the benefits of complex-valued SSMs over real-  
 094 valued SSMs. Furthermore, it uncovers an approximation-estimation tradeoff when the  
 095 target function has close dominant frequencies in the frequency domain.

## 097 2 RELATED WORKS

099 **Optimization of SSMs.** Recurrent-based architectures are known for two issues: training stability  
 100 and computational cost (Pascanu et al., 2013). To mitigate these challenges and capture long range  
 101 dependencies more effectively in sequence modeling, the S4 model was introduced with novel pa-  
 102 rameterization, initialization, and discretization techniques (Gu et al., 2022b). Recent updates to  
 103 the S4 model have further simplified the hidden state matrix by using a diagonal matrix, thereby  
 104 improving computational efficiency (Gu et al., 2022c; Gupta et al., 2022; Orvieto et al., 2023). Ad-  
 105 ditionally, regularization methods such as dropout, weight decay, and data-dependent regularizers  
 106 (Liu & Li, 2024) are employed with SSMs to prevent overfitting. In this study, we explore how  
 107

<sup>1</sup>To simplify the analysis, we omit the skip connection by letting  $D = 0$ .

temporal dependencies in input sequences impact initialization schemes in terms of optimization, with a particular focus on the timescale and state matrix.

**Curse of memory in SSMs.** The ‘‘curse of memory’’ is a newly introduced concept that highlights the difficulty recurrent-based models face in capturing long-term memory (Li et al., 2021; 2022), and has been discussed in recent works (Cirone et al., 2024; Sieber et al., 2024; Zucchet & Orvieto, 2024). This issue arises due to the exponential decay property of the model’s kernel basis functions. A common strategy to parameterize the real part of the state matrix’s eigenvalues involves stable parameterization (Gu et al., 2022c; Wang & Li, 2024), ensuring stable training dynamics even if the input sequence is infinitely long. However, this stable parameterization constrains the real part of the state matrix’s eigenvalues to be strictly negative, thereby limiting the model’s ability to capture long-term memory. In this paper, we argue that if input sequences have fixed lengths, it is reasonable to set the real part of the eigenvalues to zero by appropriately setting the timescale. This relaxation allows the model to capture long-term memory while still maintaining training stability.

### 3 PRELIMINARIES

In this section, we briefly introduce the diagonal SSM and the problem setting we consider throughout this paper. Specifically, we consider the following single-input single-output diagonal-SSM built in the complex number field  $\mathbb{C}$  and then cast into the real number field  $\mathbb{R}$ ,

$$\frac{d}{dt}h(t) = Wh(t) + bx(t), \quad y(t) = \Re(c^\top h(t)), \quad t \geq 0, \quad (1)$$

where  $\Re(\cdot)$  represents the real part;  $x(t)$  is input sequence from an input space<sup>2</sup>  $\mathcal{X} := C_0(\mathbb{R}_{\geq 0}, \mathbb{R})$ ;  $y(t) \in \mathbb{R}$  is the output sequence at time  $t$ ;  $h(t) \in \mathbb{C}^m$  is the hidden state with  $h(0) = 0$ ;  $W \in \mathbb{C}^{m \times m}$ ,  $b, c \in \mathbb{C}^m$  are trainable parameters. In particular, the state matrix  $W = \text{diag}(w_1, \dots, w_m)$  is a diagonal matrix. To simplify the analysis, we omit the skip connection matrix  $D$ . Following the training setup in Gu et al. (2022c), the read-out vector  $c$  follows standard normal distribution and the read-in vector  $b$  in (1) is fixed as an all-one vector at initialization without training. Under these settings, the input-output relation in (1) is explicitly given by the integral

$$y(t) = \int_0^t \Re(c^\top e^{ws})x(t-s)ds, \quad (2)$$

where  $w \in \mathbb{C}^m$  is the state vector that contains all the diagonal entries of the state matrix  $W$ , and the function  $\Re(c^\top e^{ws})$  is called the memory function or the kernel function.

**Discretization.** To handle discrete input sequences, we follow (Gu et al., 2022c) to use the zero-order (ZOH) hold method for discretization. Then given a timescale  $\Delta > 0$  and any discrete sequence  $(x_0, \dots, x_{L-1}) \subset \mathbb{R}$  with length  $L$ , the ZOH method induces a model output

$$y_\ell = \Re \left( \sum_{j=1}^m \frac{e^{\Delta w_j} - 1}{w_j} c_j e^{\Delta w_j (\ell-1)} \right) x_0 + \dots + \Re \left( \sum_{j=1}^m \frac{e^{\Delta w_j} - 1}{w_j} c_j e^{\Delta w_j 0} \right) x_{\ell-1}, \quad (3)$$

for  $\ell = 1, 2, \dots, L$ .

In the following section, we tackle the problems related to the initialization schemes of State Space Models (SSMs) that were introduced in the Introduction. Specifically, we will explore the following questions:

1. **Timescale Initialization:** How should we correctly initialize the model timescale  $\Delta$  for fixed-length tasks to enhance the training of SSMs? Is the previously used scaling  $\Delta = 1/L$  a universal approach?
2. **Real Part of the State Vector:** What role does  $\Re(w)$  play? Can we initialize  $\Re(w)$  to be zero, and what benefits might arise from a zero real part?
3. **Imaginary Part of the State Vector:** What role does  $\Im(w)$  play? What advantages do complex-valued SSMs offer compared to real-valued SSMs?

By probing these questions, we aim to deepen our understanding of effective initialization practices for SSMs, thereby improving their training performance.

<sup>2</sup>A linear space of continuous functions from  $\mathbb{R}_{\geq 0}$  to  $\mathbb{R}$  that vanishes at infinity.

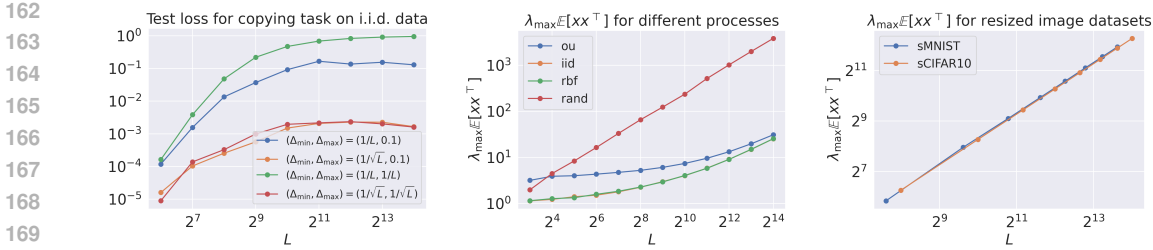


Figure 1: (Left) Training a diagonal SSM (3) on a copying task using i.i.d. data with a dimension of 128. We vary the minimal timescale  $\Delta_{\min} = 1/L, 1/\sqrt{L}$  and the maximal timescale  $\Delta_{\max} = 1/L, 1/\sqrt{L}, 0.1$  w.r.t. sequence length  $L$ . (Middle) The maximal eigenvalue of the autocorrelation matrix  $\mathbb{E}[xx^\top]$  on different random processes of  $x$ . (Right) The maximal eigenvalue of  $\mathbb{E}[xx^\top]$  on sequential image datasets sMNIST and sCIFAR10 with different resize rates varied from 0.5 to 4.

## 4 MAIN RESULTS

In this section, we present our main results by focusing on three initialization parameters  $\Delta$ ,  $\Re(w)$  and  $\Im(w)$  respectively. Specifically, in section 4.1, we rigorously characterize the relationship between the timescale  $\Delta$  and the sequence length  $L$  in terms of training stability at initialization by taking into account data autocorrelation. In section 4.2, we demonstrate that allowing the state vector’s real part to be zero can prevent exponential decay in the SSM kernel function, thereby mitigating the curse of memory in certain scenarios. In section 4.3, we explore the relationship between the convergence rate and the separation distance of the state vector’s imaginary part. In particular, we uncover an approximation-estimation tradeoff for a class of target functions.

### 4.1 RELATIONSHIP BETWEEN $\Delta$ AND $L$

In this subsection, we derive a stability condition for the ZOH-discretized diagonal SSM (3) when the state vector’s real part  $\Re(w)$  is non-positive. From both theoretical and numerical perspectives, we demonstrate that the dependency of the model timescale  $\Delta$  on the sequence length  $L$  is strongly influenced by the data autocorrelation. To start with, we prove the following theorem that provides an upper bound on the magnitude of the model output value.

**Theorem 4.1.** *Consider a ZOH discretized SSM (3) with timescale  $\Delta > 0$  and  $\Re(w_j) \leq 0$  for  $j = 1, \dots, m$ . Suppose that the input sequence  $(x_0, \dots, x_{L-1})$  is sampled from a unknown distribution in  $\mathbb{R}^L$ , and the read-out vector  $c$  is from i.i.d. standard normal distribution. Then we have*

$$\mathbb{E}_{c,x}[y_L^2] \leq \Delta^2 m^2 L \cdot \lambda_{\max}(\mathbb{E}[xx^\top]),$$

where  $\lambda_{\max}(\cdot)$  represents the maximal eigenvalue.

The proof is provided in Section C. In practice, the hidden state size  $m$  is often much smaller than the sequence length  $L$  (Gu et al., 2023). Given this, we focus on fixing the hidden size  $m$  and investigating the relationship between the model timescale  $\Delta$  and the sequence length  $L$ . We see that Theorem 4.1 connects the model timescale  $\Delta$  with the sequence length  $L$  in terms of the data autocorrelation matrix  $\mathbb{E}[xx^\top]$ . If we have normalized the sequences such that  $\mathbb{E}[\|x\|^2] = 1$ , then a simple observation is that  $1 \leq \lambda_{\max}(\mathbb{E}[xx^\top]) \leq L$  because  $\text{Tr}(\mathbb{E}[xx^\top]) = L$ . This indicates that the maximal eigenvalue of the autocorrelation matrix can have different dependencies on  $L$  based on the temporal dependencies. For example, when the elements in the sequence are uncorrelated with each other,  $x$  exhibits zero temporal dependencies, and the autocorrelation matrix is an identity matrix with  $\lambda_{\max}(\mathbb{E}[xx^\top]) = 1$ . In this case,  $\Delta$  should scale as  $1/\sqrt{L}$  to ensure training stability. On the other hand, when  $x$  is a constant sequence  $(1, 1, \dots, 1)$ , then  $x$  exhibits full temporal dependencies. The autocorrelation matrix then becomes a rank-1 matrix with  $\lambda_{\max}(\mathbb{E}[xx^\top]) = L$ , implying that  $\Delta$  should scale as  $1/L$ . Additionally, Theorem 4.1 includes the case for  $\Re(w) = 0$ . This deviates from the common practice, as noted in Gu et al. (2022c; 2023), where exponential parameterization is applied to the real part to ensure  $\Re(w)$  is strictly negative for training stability. We emphasize that for fixed-length sequences, it is also reasonable to have a zero real part, provided there is an estimate of the spectrum of the data autocorrelation matrix.

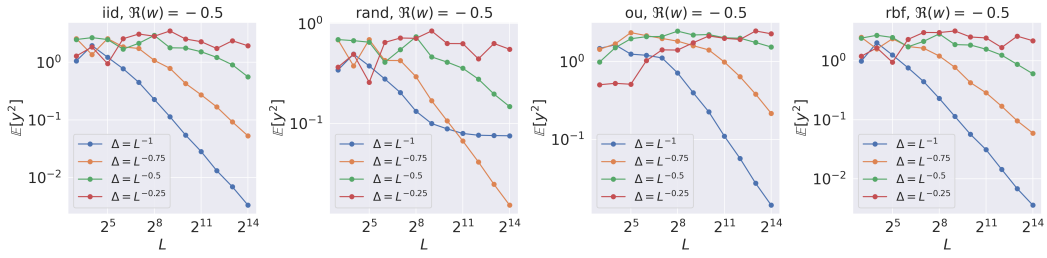


Figure 2: The expected magnitude of the SSM output value on synthetic sequences with different autocorrelation. The real part  $\Re(w) = -0.5$  follows the common practice and we consider four dependencies between the timescale  $\Delta$  and the sequence length  $L$ .

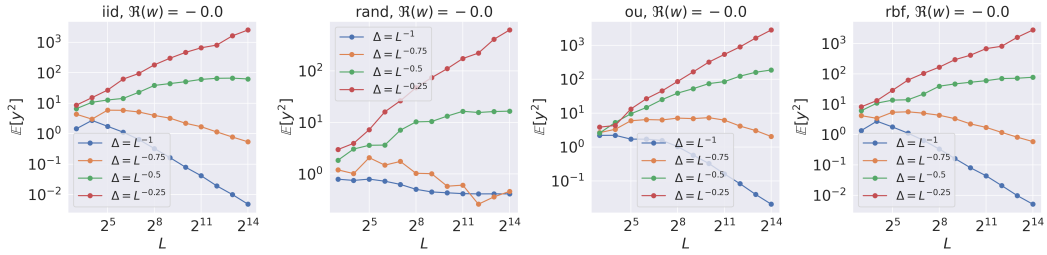


Figure 3: The expected magnitude of the SSM output value on synthetic sequences with different autocorrelation and different dependencies between  $\Delta$  and  $L$ . The real part  $\Re(w)$  is set to be zero.

**Numerical experiments on  $\lambda_{\max}(\mathbb{E}[xx^\top])$  and  $\mathbb{E}[y_L^2]$ .** To validate our theory, we conduct experiments on the exact values of the magnitude of the model output and  $\mathbb{E}[xx^\top]$ . Specifically, we consider both synthetic and real sequential datasets in both negative and zero real part cases. For synthetic datasets, we consider Gaussian process with mean 0 and autocovariance function  $\mathbb{E}[x_i x_j] = K(i, j)$ . By restricting  $K(i, i) = 1$  then the autocovariance matrix is exactly the same as the autocorrelation matrix. In this paper, we choose 4 Gaussian processes with different autocovariance functions and plot their maximal eigenvalues. The autocovariance functions for ‘ou, iid, rbf’ are  $K(i, j) = \exp(-|i - j|/\ell), \delta_{i-j}, \exp(-|i - j|^2/\ell)$  respectively. The autocovariance matrix for ‘rand’ is given by  $\Sigma\Sigma^\top$  where  $\Sigma$  is a random matrix with i.i.d. uniform distributed entries in  $[0, 1]$ . As Figure 1 (Middle) shows, different processes have varying dependencies of  $\lambda_{\max}(\mathbb{E}[xx^\top])$  on  $L$  ranging from  $\mathcal{O}(1)$  to  $\mathcal{O}(L)$ . For the i.i.d. case,  $\lambda_{\max}(\mathbb{E}[xx^\top])$  is not always 1 in Figure 1 (Middle), which is because we use the sample autocorrelation matrix to replace the expected autocorrelation matrix. For real sequential datasets, we choose to resize the MNIST dataset (LeCun et al., 2010) and the gray CIFAR10 dataset (Krizhevsky et al., 2009; Tay et al., 2021) with resize rates  $[0.5, 1, 1.5, 2, 2.5, 3, 3.5, 4]$  and the flatten the images to sequences. More experiment details are provided in Appendix A. We record  $\lambda_{\max}(\mathbb{E}[xx^\top])$  based on the entire training dataset. As shown in Figure 1 (Left), the maximal eigenvalue scales (almost) linearly with sequence length across the resize rate for both sequential MNIST (sMNIST) and sequential CIFAR10 (sCIFAR10) datasets. Additionally, we plot the relationship between the magnitude of the model output value and sequence length by varying the timescale  $\Delta = [L^{-1}, L^{-0.75}, L^{-0.5}, L^{-0.25}]$ . In Figures 2 and 4, when  $\Re(w) = -0.5$  (following the setup in Gu et al. (2022c; 2023)), the magnitude  $\mathbb{E}[y_L^2]$  remains stable for both synthetic and resized image datasets for all decay rates of  $\Delta$ . When  $\Re(w) = 0$ , Figures 3 and 4 demonstrate that for the ‘rand’ process,  $\Delta = L^{-1}$  is stable. For the ‘iid,’ ‘ou,’ and ‘rbf’ processes,  $\Delta = L^{-0.75}$  is stable. This indicates that our bound in Theorem 4.1 effectively characterizes the relation between  $\Delta$  and  $L$  for  $\Re(w) = 0$ . Moreover, as shown in Figure 4, for the sequential image datasets,  $\Delta$  should scale as  $1/L$  to ensure stability when  $\Re(w) = 0$ ; otherwise, the magnitude increases with sequence length. This finding aligns with the empirical results in (Gu et al., 2022c; 2023) that  $\Delta$  should scale as  $1/L$  to effectively capture the range of dependencies for length  $L$  for real-world tasks. But their theoretical reasons are based on Fourier analysis of continuous-time SSMs and do not explicitly account for the data autocorrelation.



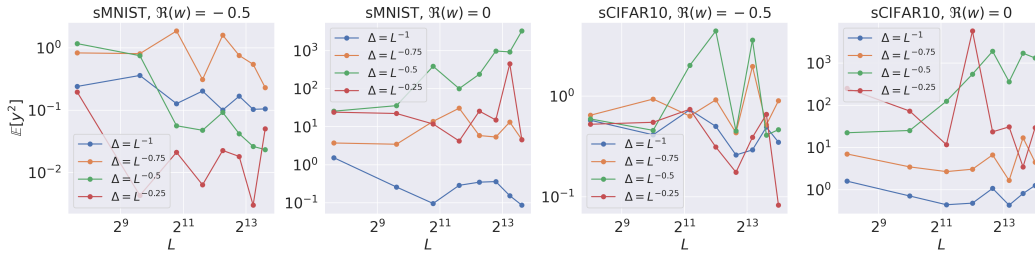


Figure 4: The expected magnitude of the SSM output value on sequential image datasets with different resize rates (ranging from 0.5 to 4) and different dependencies between  $\Delta$  and  $L$ .

**Experiments on copying task with different timescales.** We also tested the performance of the diagonal state-space model (SSM) (3) on a copying task with various dependencies of  $\Delta$  on  $L$ . It is worth noting that, as discussed in Jelassi et al. (2024), SSMs struggle with the copying task because the model’s state dimension needs to scale linearly with the sequence length to memorize all the input tokens. However, the limitation highlighted in Jelassi et al. (2024) pertains to the length generalization task—i.e., training an SSM with short sequences and then testing it on longer sequences will fail if the hidden size  $m$  does not grow linearly with  $L$ . Here, we focus on a fixed-length task, where both training and test sequences have the same length. We find that, with an appropriately initialized timescale, SSMs can effectively handle the copying task even with a small state size. In this paper, we use a diagonal SSM with a fixed state size  $m = 32$  to learn a copying task on i.i.d. data with a dimension of 128, and the timescale  $\Delta \in \mathbb{R}^{128}$ . We vary the minimal and maximal timescales ( $\Delta_{\min}, \Delta_{\max}$ ) with different dependencies on  $L$ . From Figure 1 (Left), we see that the combination  $(\Delta_{\min}, \Delta_{\max}) = (1/L, 0.1)$ , which is commonly used in practice (Gu et al., 2022c; 2023) to train real datasets, consistently performs worse than setting  $\Delta_{\min} = 1/\sqrt{L}$ . This stable scaling is in line with our theoretical suggestions for i.i.d. data. Therefore, the data autocorrelation is very crucial for us to get a good initialization scale on the timescale. More experiment details are provided in Appendix A.

#### 4.2 BENEFITS OF ZERO REAL PART

In this subsection, we investigate the benefits of initializing  $\Re(w) = 0$  for tasks that require long-term memory. In previous works (Li et al., 2021; 2022), it is shown that recurrent-based models suffer from the curse of memory in both approximation and optimization when there is long-term memory in the target. For example, we consider using a diagonal SSM (3) to learn an input-output relationship given by a real-valued target function  $\rho^*$  such that

$$y_\ell^* = \rho_{\ell-1}^* x_0 + \dots + \rho_0^* x_{\ell-1}, \quad \ell = 1, 2, \dots, L.$$

The objective function is given by the squared difference between the model output  $y_L$  and the corresponding label  $y_L^*$ . Then in a special case when the input sequences have zero temporal dependencies with  $\mathbb{E}[xx^\top] = \mathbb{I}_L$ , the expected mean squared error is given by

$$\mathbb{E}[|y_L - y_L^*|^2] = \|\tilde{\rho} - \rho^*\|^2,$$

where  $\tilde{\rho}$  is a vector  $(\Re(\sum_{j=1}^m \frac{e^{\Delta w_j} - 1}{w_j} c_j e^{\Delta w_j 0}), \dots, \Re(\sum_{j=1}^m \frac{e^{\Delta w_j} - 1}{w_j} c_j e^{\Delta w_j (L-1)}))$  that represents the model’s memory function, and  $\rho^* = (\rho_0^*, \dots, \rho_{L-1}^*)$ . Therefore, a well-trained SSM means that the model memory function matches with the target function, i.e.,

$$\Re\left(\sum_{j=1}^m \frac{e^{\Delta w_j} - 1}{w_j} c_j e^{\Delta w_j \ell}\right) = \rho_\ell^*, \quad \ell = 0, \dots, L-1.$$

Then we can see that the curse of memory happens when the target function  $\rho^*$  has a sudden spike in a very long distance. For instance, consider a shifting task that requires mapping an input sequence  $(x_0, \dots, x_{L-1})$  to a shifted sequence  $(0, \dots, 0, x_0)$ . In this task, the target  $\rho^*$  is  $(0, \dots, 0, 1)$ , which is challenging for an exponentially decaying SSM kernel  $\tilde{\rho}$  to capture long-term memory when

$\Re(w) < 0$ . However, if we allow the real part to be zero at initialization, then  $\tilde{\rho}$  does not undergo exponential decay. As a result, we can potentially avoid the curse of memory, even for long sequences, in this scenario. It is worth noting that in this paper, we do not consider a stable parameterization to ensure  $\Re(w) \leq 0$  strictly during training, which means there may be some optimization stability issues since  $\Re(w) = 0$  is on the stability boundary. However, our experiments on simple tasks demonstrate that initializing with a zero real part still helps enhance training, even without a stable parameterization. This suggests that, despite the potential optimization stability challenges during training, a zero real part can be beneficial for training on certain tasks.

**Experiments on the benefits of zero real part.** To validate the effectiveness of having a zero real part, we conduct experiments on both synthetic and real datasets that require long-term memory. For the synthetic task, we use i.i.d. sequential data to easily visualize the expected error via the memory function. The goal is to learn an input-output mapping from  $(x_0, \dots, x_{L-1})$  to  $x_0 + x_{L-1}$ , which requires the model to memorize both the first and last token. The target memory function  $\rho^*$  is  $(1, 0, \dots, 0, 1)$ . In our setting, the sequence length  $L$  is 128, and the hidden state size  $m$  is 32. As shown in Figure 5 (Left) and (Middle), the SSM with a zero real part outperforms the case with a negative real part. It is evident that by initializing  $\Re(w) = 0$ , the learned memory function is able to capture long range dependencies. For the real-world task, we utilize the sequential MNIST (sMNIST) dataset. Before training, we preprocess the entire dataset with a linear transformation to decorrelate the training sequences, resulting in an autocorrelation matrix that is an identity matrix. We recover the underlying target memory function by solving a least square problem  $\min_{\rho} \|X * \rho - Y\|_F^2$  where  $X \in \mathbb{R}^{50000 \times 784}$  is the collected sequence matrix,  $Y \in \mathbb{R}^{50000 \times 10}$  is the one-hot label matrix, and  $*$  denotes the convolution operator. The recovered target memory function  $\rho \in \mathbb{R}^{784 \times 10}$  has 10 channels. To illustrate the underlying memory patterns, we plot  $\sqrt{L}\rho$  for each channel in Figure 6. We observe that for the decorrelated sMNIST dataset, the underlying memory function exhibits a sudden spike at a long distance, implying the curse of memory when  $\Re(w) < 0$ . This observation is confirmed in Figure 5 (Right), which shows that initializing  $\Re(w) = 0$  outperforms the case with a negative real part. More experimental details are provided in Appendix A.

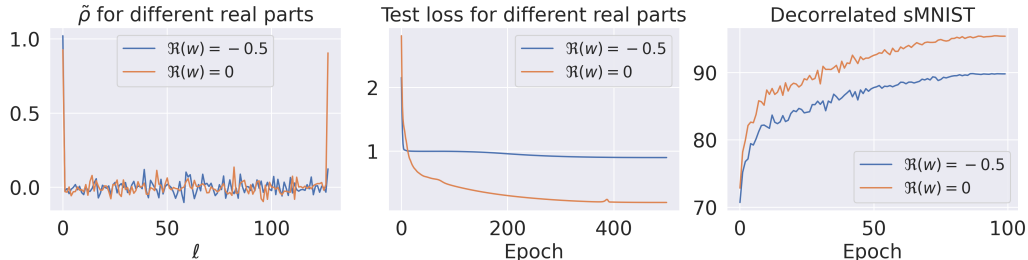
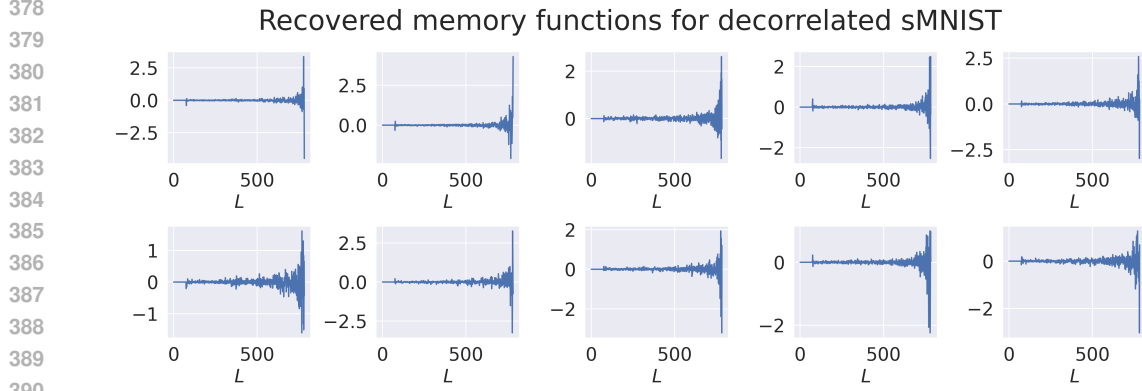


Figure 5: (Left) Training a diagonal SSM (3) on a task that requires long-term memory. The learned memory function  $\tilde{\rho}$  effectively captures the spike in long-range dependencies. However, it struggles to do so when the real part is negative. (Middle) Test loss on the long-term memory task when initializing  $\Re(w) = 0$  and  $\Re(w) = -0.5$ . (Right) Test accuracy for training a diagonal SSM on decorrelated sequential MNIST dataset with different real parts at initialization.

#### 4.3 IMAGINARY PART INDUCES AN APPROXIMATION-ESTIMATION TRADEOFF

In the previous subsection, we show that the real part  $\Re(w)$  is related with the long-term memory when training SSMs. In this subsection, we focus on the imaginary part  $\Im(w)$ . We will demonstrate how  $\Im(w)$  influences the conditioning of the SSM optimization problem within a convex framework. Additionally, from an approximation standpoint, we reveal an approximation-estimation tradeoff that arises when training SSMs with a particular class of target functions.

**Convergence analysis.** Here we consider the continuous-time SSM (2) and assume that the read-out vector  $c$  is in  $\mathbb{R}^m$ . Now we define the loss function. Suppose the ground truth input-output relation is given by some real-valued target function  $\rho^*(s) \in L^1[0, \infty) \wedge L^2[0, \infty)$  with  $y^*(t) = \int_0^t \rho^*(s)x(t-s)ds$ . We use the squared difference between the SSM output  $y(t)$  and the target



391  
392  
393  
394  
395  
396

Figure 6: Recovering the memory function  $\rho$  on the decorrelated sequential MNIST dataset by solving a linear equation  $X * \rho = Y$ , where  $X \in \mathbb{R}^{N \times L}$  is the collected sequence matrix,  $Y \in \mathbb{R}^{N \times 10}$  is the one-hot label matrix, and  $*$  is the convolution operator. Then  $\rho \in \mathbb{R}^{L \times 10}$  has 10 channels and we plot the scaled function  $\sqrt{L}\rho$  each channel to show the underlying memory patterns.

397  
398

output  $y^*(t)$  at some terminal time  $T > 0$  averaged over input distributions, which can be written as

399  
400

$$\mathcal{L}(c, a) := \mathbb{E}_x (y(T) - y^*(T))^2. \quad (4)$$

401  
402  
403

To make the theoretical analysis amenable, we make the simplification that the input sequence  $x(t)$  is sampled from white noise, i.e.,  $x(T-s)ds = dW_s$  where  $W_s$  is the canonical real-valued Wiener process. Then by Itô's isometry (Proposition B.2), the expected risk (4) can be rewritten as

404  
405  
406

$$\mathcal{L}(c, w) = \int_0^T (c^\top \Re(e^{ws}) - \rho^*(s))^2 ds.$$

407  
408  
409

In the practical training, the sequence length is very long and thus we take  $T \rightarrow \infty$  to investigate the effect of long-term memory. To study the effects of the state vector initialization, we consider the following convex optimization problem where  $w$  is fixed.

410  
411  
412  
413

$$\arg \min_{c \in \mathbb{R}^m} \mathcal{L}_c := \int_0^\infty \left( \sum_{j=1}^m c_j \Re(e^{w_j s}) - \rho^*(s) \right)^2 ds. \quad (5)$$

414  
415  
416  
417  
418  
419

From the perspective of function approximation, the HiPPO framework (Gu et al., 2020) initializes  $w$  such that the SSM basis kernel functions  $\{\Re(e^{w_j s})\}_{j=1}^\infty$  are orthogonal in  $L^2[0, \infty)$  w.r.t. some measure  $\omega(s)$ . In this paper, we discover the effects of the state initialization on the optimization problem (5). Let  $c^*$  be one of the solution of the convex problem (5), then  $c^*$  is a stationary point that satisfies

420  
421

$$Gc^* = \int_0^\infty \Re(e^{ws})\rho^*(s)ds,$$

422

where  $G \in \mathbb{R}^{m \times m}$  is a Gram matrix with

423  
424  
425

$$[G]_{j,k} = \int_0^\infty \Re(e^{w_j s})\Re(e^{w_k s})ds. \quad (6)$$

426  
427  
428  
429

Therefore, the spectrum of the Gram matrix  $G$  determines the numerical stability and convergence rate of optimization algorithms for solving the convex problem (5). We show in the following proposition that when  $w \in \mathbb{R}^m$  and all  $w_j$  are distinct, or when  $w \in \mathbb{C}^m$  and all the imaginary parts  $\Im(w)$  are non-zero and distinct, then  $G$  is positive definite.

430  
431

**Proposition 4.2.** *Let  $w_j = a_j + i \cdot v_j$  with  $a_j, v_j \in \mathbb{R}$  for  $j = 1, \dots, m$ . If all  $v_j = 0$ , i.e.,  $w \in \mathbb{R}^m$ , then  $G$  is positive definite given that all  $a_j$  are distinct. If  $v_j$  are all non-zero, i.e.,  $w \in \mathbb{C}^m$ , then  $G$  is positive definite given that all  $v_j$  are distinct.*



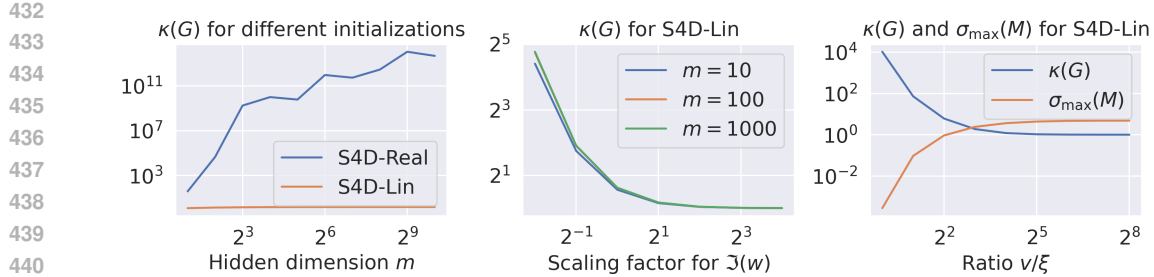


Figure 7: (Left) Condition number  $\kappa(G) := \frac{\lambda_{\max}(G)}{\lambda_{\min}(G)}$  for S4D-Real and S4D-Lin with different hidden size  $m$ . (Middle)  $\kappa(G)$  for S4D-Lin with different  $m$  by varying scaling factors of the imaginary part  $\Im(a)$ . (Right)  $\kappa(G)$  and approximation measure  $\sigma_{\max}(M)$  (in the approximation-estimation tradeoff part) for S4D-Lin by different ratios of model frequencies  $v$  and target frequencies  $\xi$ .

The proof is based on the argument of Vandermonde matrix, and we provide details in Appendix D. Given that the gram matrix  $G$  is positive-definite, we are ready to study its spectrum. In the following theorem, we show that for complex-valued SSMs, the gram matrix  $G$  can be well-conditioned provided that the imaginary parts  $\Im(w)$  are well separated.

**Theorem 4.3.** *Let  $\lambda_{\min}(G), \lambda_{\max}(G)$  be the extreme eigenvalues of  $G$  defined in (6), and let  $\coth(x) = \frac{e^{2x}+1}{e^{2x}-1}$ . Suppose that  $w_j = -0.5 + i \cdot v_j$  for  $v_j \in \mathbb{R}$ , and we define the separation distance  $\delta := \min_{j \neq k} |v_j - v_k|$ . Then if  $\delta > 0$ , we have*

$$1.19 - \frac{3\pi}{4\delta} \coth\left(\frac{\pi}{\delta}\right) < \lambda_{\min}(G) \leq \lambda_{\max}(G) < \frac{5}{12} + \frac{3\pi}{4\delta} \coth\left(\frac{\pi}{\delta}\right).$$

The proof is based on the Gershgorin circle theorem, with details provided in Appendix E. The setup  $w_j = -0.5 + i \cdot v_j$  follows the configurations in Gu et al. (2022a;c; 2023). This theorem shows that the Gram matrix  $G$  can be well-conditioned when the separation distance  $\delta$  is large. One example is that for the commonly used S4D-Lin initialization (Gu et al., 2022c),  $v_j = \pi \cdot j$ . Then the separation distance  $\delta = \pi$ . Numerical calculations show that  $0.2 < \lambda_{\min}(G) \leq \lambda_{\max}(G) < \sqrt{2}$ , meaning that  $G$  is well-conditioned for any hidden size  $m$ , and its condition number has a uniform upper bound w.r.t.  $m$ . Note that  $x \coth(x) \geq 1$  and is increasing on  $[0, \infty)$ , which implies that the bound for  $\lambda_{\min}(G)$  is non-trivial when  $\frac{3\pi}{4\delta} \coth\left(\frac{\pi}{\delta}\right) < 1.19$ . By numerically solving this inequality, it is sufficient to have  $\delta > 2.3$ . However, Proposition 4.2 suggests that as long as  $\delta > 0$ , the positive-definiteness of  $G$  is guaranteed. This indicates a gap between the lower bound and the actual minimal eigenvalue, which we leave for future research.

**Real vs complex.** We can now compare real-valued SSMs and complex-valued SSMs in terms of the conditioning of the convex optimization problem (5), which is determined by the condition number of  $G$ . For real-valued SSMs with the S4D-Real initialization (Gu et al., 2022c), where  $w_j = -j$ , we have  $G_{j,k} = \frac{1}{j+k}$ . In this case,  $G$  is a Hilbert matrix, whose condition number grows exponentially with respect to its size  $m$  (Todd, 1953). For complex-valued SSMs with  $w_j = -0.5 + iv_j$ , Theorem 4.3 indicates that if the separation distance  $\delta$  remains uniformly large with respect to  $m$ , then  $G$  can be well-conditioned even for larger values of  $m$ . For S4D-Lin initialization, we already know that  $0.2 < \lambda_{\min}(G) \leq \lambda_{\max}(G) < \sqrt{2}$  by the above argument. Therefore, unlike real-valued SSMs, the condition number of  $G$  in the complex-valued case can remain well-conditioned even for large  $m$ , given that the imaginary parts are well separated. This difference is illustrated in Figure 7 (Left), where we compare the exact condition numbers for S4D-Real and S4D-Lin. As the scaling factor of the imaginary part increases, the separation distance also increases. Figure 7 (Middle) shows that the Gram matrix  $G$  for S4D-Lin becomes better conditioned, validating Theorem 4.3.

**Approximation-estimation tradeoff.** Despite the fact that complex-valued SSMs with adequately separated imaginary parts  $\Im(w)$  enhance the conditioning of  $G$ , we cannot simply initialize  $w$  with widely separated  $\Im(w)$ . This is because  $\Im(w)$  determines the frequencies that the SSM can capture, and misaligned frequencies relative to the target  $\rho^*$  lead to a large approximation error  $\mathcal{L}_{c^*}$ . For example, suppose that the target memory function  $\rho^*(s) = e^{-s/2} \hat{c}^\top \cos(\xi s)$  with  $\hat{c}, \xi \in \mathbb{R}^m$ . Let

486  $w = -0.5 + iv$  for  $v \in \mathbb{R}^m$ , then we have

487  
488  
489  
490  
491

$$\begin{aligned} \mathcal{L}_{c^*} &= \int_0^\infty \rho^{*2}(s) ds - \left( \int_0^\infty e^{-\frac{s}{2}} \cos(vs) \rho^*(s) ds \right)^\top G^{-1} \left( \int_0^\infty e^{-\frac{s}{2}} \cos(vs) \rho^*(s) ds \right) \\ &= \hat{c}^\top M c, \end{aligned}$$

492 where  $M \in \mathbb{R}^{m \times m}$  is given by

493  
494  
495

$$\int_0^\infty e^{-s} \cos(\xi s) \cos(\xi s)^\top ds - \left( \int_0^\infty e^{-s} \cos(\xi s) \cos(v s)^\top ds \right) G^{-1} \left( \int_0^\infty e^{-s} \cos(v s) \cos(\xi s)^\top ds \right).$$

496 We can see that the maximum singular value  $\sigma_{\max}(M)$  of  $M$  determines the approximation error.  
497 Now, let's consider a limiting case when  $v_j = \mu j$  with  $\mu \rightarrow \infty$ . According to Lemma B.5, we know  
498 that  $G = \frac{1}{2} \mathbb{I}_m$ , a scaled identity matrix, possesses the best possible conditioning. Furthermore, if  
499  $\xi$  is finite, then as  $\mu \rightarrow \infty$ ,  $\int_0^\infty e^{-s} \cos(v_j s) \cos(\xi_k s) ds = 0$ , indicating that the worst approxi-  
500 mation error  $\int_0^\infty \rho^{*2}(s) ds$ . On the other hand, if we aim to minimize the approximation error, we  
501 might align the frequencies such that  $v = \xi$ . However, when the target function comprises closely  
502 spaced frequencies  $\xi_1, \dots, \xi_m$ , such alignment may cause  $G$  to have a large condition number (as  
503 per Theorem 4.3). Balancing these two aspects reveals an approximation-estimation tradeoff, which  
504 is crucial when selecting an SSM initialization. Numerical evidence for this tradeoff is illustrated  
505 in Figure 7 (Right). In this figure, we set  $\xi_j = 0.1\pi j$  with a relatively small separation distance  
506  $\delta = 0.1\pi$ , and we vary the ratio  $v_j/\xi_j$  from  $2^0$  to  $2^8$ . As the ratio increases, the optimization is  
507 expected to improve, while the approximation deteriorates. This trend is shown in Figure 7 (Right),  
508 where the induced Gram matrix  $G$  becomes better-conditioned, whereas the approximation measure  
509  $\sigma_{\max}(M)$  increases. In practice, the approximation-estimation tradeoff indicates that selecting a  
510 data-dependent initialization for  $\Im(w)$ , based on the dominant frequencies of the target function,  
511 can strike a balance that optimizes performance for a given training budget, such as the number of  
512 training epochs.

## 513 5 CONCLUSION

514  
515  
516 In this paper, we study the question proposed in the Introduction section, focusing on two initializa-  
517 tion schemes for state space models (SSMs): the timescale  $\Delta$  and the state matrix  $W$ . Regarding the  
518 timescale  $\Delta$ , we investigate it from the perspective of training stability at initialization. Our findings  
519 indicate that its dependency on sequence length is determined by data autocorrelation. By analyzing  
520 data autocorrelation, we can initialize  $\Delta$  to enhance SSM training for tasks involving fixed-length  
521 sequences. For the state matrix  $W$ , we differentiate between the real part  $\Re(W)$  and the imagi-  
522 nary part  $\Im(W)$ . The real part  $\Re(W)$  is crucial for capturing long-term memory in temporal data.  
523 Allowing for a zero real part can effectively mitigate the curse of memory while maintaining train-  
524 ing stability at initialization, provided the timescale is appropriately initialized. The imaginary part  
525  $\Im(W)$  affects the conditioning of the SSM optimization problem. A well-separated  $\Im(W)$  leads to  
526 a well-conditioned Gram matrix, improving the convergence rate. However, from an approxima-  
527 tion standpoint, excessively increasing the separation distance can result in a frequency mismatch  
528 between the SSM and the target function, leading to an approximation-estimation tradeoff. There  
529 are several potential future interesting directions. For instance, we have not discussed the effects of  
530 gating (Mehta et al., 2023) and model depth on the approximation and optimization of SSMs, which  
531 we leave for future research.

532  
533  
534  
535  
536  
537  
538  
539

## REFERENCES

- 540  
541  
542 Aviv Bick, Kevin Y Li, Eric P Xing, J Zico Kolter, and Albert Gu. Transformers to ssms: Distilling  
543 quadratic knowledge to subquadratic models. *arXiv preprint arXiv:2408.10189*, 2024.
- 544 Nicola Muca Cirone, Antonio Orvieto, Benjamin Walker, Cristopher Salvi, and Terry Lyons. Theo-  
545 retical foundations of deep selective state-space models. *arXiv preprint arXiv:2402.19047*, 2024.
- 546  
547 Tri Dao and Albert Gu. Transformers are SSMs: Generalized models and efficient algorithms  
548 through structured state space duality. In *Forty-first International Conference on Machine Learn-*  
549 *ing*, 2024.
- 550 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*  
551 *preprint arXiv:2312.00752*, 2023.
- 552 Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory  
553 with optimal polynomial projections. *Advances in neural information processing systems*, 33:  
554 1474–1487, 2020.
- 555  
556 Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization  
557 of diagonal state space models. *Advances in Neural Information Processing Systems*, 35:35971–  
558 35983, 2022a.
- 559 Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured  
560 state spaces. In *International Conference on Learning Representations*, 2022b.
- 561  
562 Albert Gu, Ankit Gupta, Karan Goel, and Christopher Ré. On the parameterization and initialization  
563 of diagonal state space models. *Advances in Neural Information Processing Systems*, 35, 2022c.
- 564 Albert Gu, Isys Johnson, Aman Timalsina, Atri Rudra, and Christopher Re. How to train your  
565 HIPPO: State space models with generalized orthogonal basis projections. In *International Con-*  
566 *ference on Learning Representations*, 2023.
- 567  
568 Ankit Gupta, Albert Gu, and Jonathan Berant. Diagonal state spaces are as effective as structured  
569 state spaces. In *Advances in Neural Information Processing Systems*, 2022.
- 570 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint*  
571 *arXiv:1606.08415*, 2016.
- 572  
573 Sukjun Hwang, Aakash Lahoti, Tri Dao, and Albert Gu. Hydra: Bidirectional state space models  
574 through generalized matrix mixers. *arXiv preprint arXiv:2407.09941*, 2024.
- 575  
576 Samy Jelassi, David Brandfonbrener, Sham M. Kakade, and eran malach. Repeat after me: Trans-  
577 formers are better than state space models at copying. In *Forty-first International Conference on*  
*Machine Learning*, 2024.
- 578  
579 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,  
580 2014.
- 581 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
582 2009.
- 583  
584 Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online].*,  
585 2, 2010.
- 586  
587 Zhong Li, Jiequn Han, Weinan E, and Qianxiao Li. On the curse of memory in recurrent neural  
588 networks: Approximation and optimization analysis. In *International Conference on Learning*  
*Representations*, 2021.
- 589  
590 Zhong Li, Jiequn Han, E Weinan, and Qianxiao Li. Approximation and optimization theory for  
591 linear continuous-time recurrent neural networks. *Journal of Machine Learning Research*, 23  
592 (42):1–85, 2022.
- 593  
Fusheng Liu and Qianxiao Li. From generalization analysis to optimization designs for state space  
models. In *Forty-first International Conference on Machine Learning*, 2024.

- 594 Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and  
595 Yunfan Liu. Vmamba: Visual state space model. *arXiv preprint arXiv:2401.10166*, 2024.
- 596
- 597 Harsh Mehta, Ankit Gupta, Ashok Cutkosky, and Behnam Neyshabur. Long range language model-  
598 ing via gated state spaces. In *The Eleventh International Conference on Learning Representations*,  
599 2023.
- 600 Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pas-  
601 canu, and Soham De. Resurrecting recurrent neural networks for long sequences. In *International  
602 Conference on Machine Learning*, pp. 26670–26698. PMLR, 2023.
- 603
- 604 Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural  
605 networks. In *Proceedings of the 30th International Conference on Machine Learning*, 2013.
- 606 Syama Sundar Rangapuram, Matthias W Seeger, Jan Gasthaus, Lorenzo Stella, Yuyang Wang, and  
607 Tim Januschowski. Deep state space models for time series forecasting. In *Advances in Neural  
608 Information Processing Systems*. Curran Associates, Inc., 2018.
- 609 Jerome Sieber, Carmen Amo Alonso, Alexandre Didier, Melanie N Zeilinger, and Antonio Orvieto.  
610 Understanding the differences in foundation models: Attention, state space models, and recurrent  
611 neural networks. *arXiv preprint arXiv:2405.15731*, 2024.
- 612
- 613 Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao,  
614 Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena : A benchmark for efficient  
615 transformers. In *International Conference on Learning Representations*, 2021.
- 616 J. Todd. *The condition of the finite segments of the Hilbert matrix*. National Bureau of Standards  
617 Applied Mathematics Series, 1953.
- 618
- 619 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
620 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-  
621 tion processing systems*, 30, 2017.
- 622 Roger Waleffe, Wonmin Byeon, Duncan Riach, Brandon Norick, Vijay Korthikanti, Tri Dao, Albert  
623 Gu, Ali Hatamizadeh, Sudhakar Singh, Deepak Narayanan, et al. An empirical study of mamba-  
624 based language models. *arXiv preprint arXiv:2406.07887*, 2024.
- 625
- 626 Junxiong Wang, Daniele Paliotta, Avner May, Alexander M Rush, and Tri Dao. The mamba in the  
627 llama: Distilling and accelerating hybrid models. *arXiv preprint arXiv:2408.15237*, 2024.
- 628 Shida Wang and Qianxiao Li. StableSSM: Alleviating the curse of memory in state-space models  
629 through stable reparameterization. In *Forty-first International Conference on Machine Learning*,  
630 2024.
- 631 Annan Yu, Michael W Mahoney, and N Benjamin Erichson. There is hope to avoid hippos for  
632 long-memory state space models. *arXiv preprint arXiv:2405.13975*, 2024.
- 633
- 634 Michael Zhang, Khaled Kamal Saab, Michael Poli, Tri Dao, Karan Goel, and Christopher Re. Ef-  
635 fectively modeling time series with simple discrete state spaces. In *The Eleventh International  
636 Conference on Learning Representations*, 2023.
- 637 Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision  
638 mamba: Efficient visual representation learning with bidirectional state space model. In *Forty-first  
639 International Conference on Machine Learning*, 2024.
- 640
- 641 Nicolas Zucchet and Antonio Orvieto. Recurrent neural networks: vanishing and exploding gradi-  
642 ents are not the end of the story. *arXiv preprint arXiv:2405.21064*, 2024.
- 643
- 644
- 645
- 646
- 647

## A EXPERIMENTS DETAILS

In this section, we provide more experiment details that produce Figure 1, 2, 3, 4, 5, 6 in section 4.

**Figure 1 (Left).** The synthetic dataset that we use to produce Figure 1 (Left) is i.i.d. sampled from standard normal distribution with dimension 128, i.e., each input sequence is of shape  $(1, L, 128)$  where  $L$  is its sequence length. We use a ZOH discretized diagonal SSM layer (3) with hidden size  $m = 32$ , model dimension  $d = 128$  to handle the 128 dimensional dataset. We initialize the state vector  $w$  by S4D-Lin with real part  $-0.5$ . The read-out vector  $c$  is initialized as i.i.d. standard normal distribution. We vary  $\Delta_{\min}$  and  $\Delta_{\max}$  in the SSM layer and use the Adam optimizer (Kingma, 2014) to train the hyperparameters  $\Delta, \Re(w), \Im(w), C$  without weight decay. The learning rate for  $\Delta, \Re(w), \Im(w)$  is 0.001 and the learning rate for  $c$  is 0.1.

**Figure 1 (Middle), 2, 3.** The synthetic datasets that we use to produce these figures are Gaussian processes with mean zero and varied autocovariance functions  $\mathbb{E}[x_i x_j] = K(i, j)$  for  $i, j = 1, 2, \dots, L$ . Specifically, the ‘iid’ dataset refers to  $K(i, j) = \delta_{i-j}$ ; the ‘ou’ dataset refers to  $K(i, j) = \exp(-|i - j|/2)$ ; the ‘rbf’ dataset refers to  $K(i, j) = \exp(-\pi|i - j|^2)$ ; and the autocovariance matrix for the ‘rand’ dataset is given by  $\Sigma \Sigma^\top / L$  where  $\Sigma \in \mathbb{R}^{L \times L}$  is a random matrix with i.i.d. entries sampled from a uniform distribution  $U[0, \sqrt{3}]$ . For all the four synthetic datasets, we have  $K(i, i) = 1$ . The plot for Figure 1 (Middle) records the maximal eigenvalue of the sample matrix that we fix the data size to be 1000 and vary the sequence length  $L$  as plotted. So we can see some deviations between theory and practice. For Figure 2 & 3, we also use the 1-dimensional SSM layer (3) with S4D-Lin initialization on  $\Im(w)$  and vary the real part  $\Re(w)$  to be  $-0.5$  or 0.

**Figure 1 (Right), 4, 5 (Right), 6.** For the resized sequential image datasets, we choose to resize the original images with resize rates  $[0.5, 1, 1.5, 2, 2.5, 3, 3.5, 4]$ . Then we standardize the whole images and flatten them into 1-d sequence. For sequential MNIST (sMNIST) dataset, the sequence length is  $784r^2$  and for sequential CIFAR10 (sCIFAR10), the sequence length is  $1024r^2$  where  $r$  is the resize rate. The plot for the maximal eigenvalue of the autocorrelation matrix and the output value are based on the whole training set. We use the 1-dimensional SSM layer (3) with S4D-Lin initialization and vary the real part  $\Re(w)$  to be  $-0.5$  or 0 to calculate the output value magnitude. For the decorrelated sMNIST dataset, we choose the original MNIST dataset and the decorrelation transformation is given by a centered matrix with a whitening matrix after flattening images. The centered matrix is the mean of the sequential data along the batch dimension, and the whitening matrix has shape  $L \times L$ . The whitening matrix can be obtained by SVD on the data matrix. To train the decorrelated sMNIST dataset, we use a 128-dimensional SSM layer (3) with  $m = 32$  and GELU activation (Hendrycks & Gimpel, 2016) on the model output, and also apply a gated linear unit after the GELU activation. We use dropout with rate 0.1 and apply a decoder layer for classification. We use Adam optimizer with learning rate 0.001 on  $\Delta, \Re(w), \Im(w)$  and AdamW optimizer with weight decay 0.01 on the rest hyperparameters. For the plot of the memory function in Figure 6, we solve a least square problem by taking the pseudo inverse of the sequence matrix  $X \in \mathbb{R}^{50000 \times 784}$  and then get the recovered memory function  $\rho$ .

**Figure 5 (Left), (Middle).** The comparisons on zero real part and negative real part in Figure 5 (Left) & (Middle) are conducted on a 1-dimensional synthetic dataset. We sample the training and test dataset from i.i.d. standard normal distribution with length 128. The training sample size and the test sample size are both 1000. We use the SSM layer (3) with  $m = 32$ , S4D-Lin initialization on  $\Im(w)$  and initialize the timescale  $\Delta = 1/\sqrt{128}$ . We use Adam optimizer with learning rate 0.001 on  $\Delta, \Re(w), \Im(w)$  and learning rate 0.01 on  $c$ .

### A.1 ADDITIONAL EXPERIMENTS FOR S4D-LEGS INITIALIZATION

In this subsection, we include more experiment results in Figure 8, 9, 10, 11 for SSMs with S4D-Legs (Gu et al., 2022c) initialization on the imaginary part  $\Im(w)$ . The S4D-Legs initialization is an approximation on the original S4-Legs initialization (Gu et al., 2022b) by taking diagonal part of the diagonal plus low-rank HiPPO-Legs matrix. In Figure 8, 9, 10, we plot the magnitude of the SSM output value given the S4D-Legs initialization for both zero real part and negative real part cases. The experiment settings follow the guidelines we introduce before with only a change on the initialization of  $\Im(w)$ . We can see that for S4D-Legs initialization, our conclusion still holds in the sense that negative real part is stable at initialization for all all the scaling that we



considered in this paper, while for zero real part, the dependencies of  $\Delta$  on  $L$  varies for different sequence autocorrelation. We also compare the effects of real parts on optimization with S4D-Legs initialization. The results are shown in Figure 11 and we obtain consistent results as the S4D-Lin initialization. One interesting finding is that on the decorrelated sMNIST dataset, the comparison between Figure 5 (Right) and Figure 11 (Right) shows that the S4D-Lin initialization outperforms the S4D-Legs initialization in both zero real part and negative real part cases.

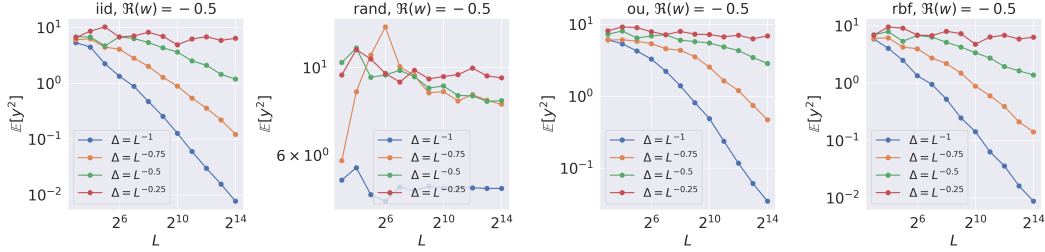


Figure 8: The expected magnitude of the SSM output value on synthetic sequences with S4D-Legs initialization and different autocorrelation. The real part  $\Re(w) = -0.5$  follows the common practice and we consider four dependencies between the timescale  $\Delta$  and the sequence length  $L$ .

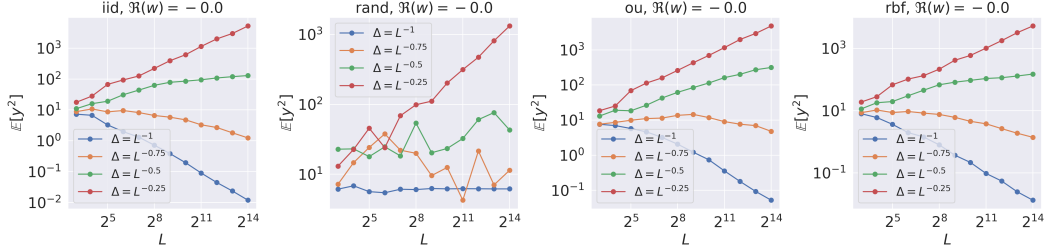


Figure 9: The expected magnitude of the SSM output value on synthetic sequences with S4D-Legs initialization and different autocorrelation and different dependencies between  $\Delta$  and  $L$ . The real part  $\Re(w)$  is set to be zero.

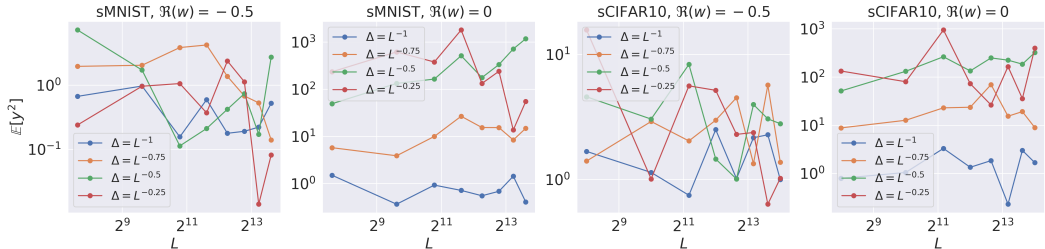


Figure 10: The expected magnitude of the SSM output value for S4D-Legs initialization on sequential image datasets with different resize rates (ranging from 0.5 to 4) and different dependencies between  $\Delta$  and  $L$ .

## B AUXILIARY LEMMAS

In this section, we provide the description for Itô’s isometry and a few auxiliary lemmas that we will need for the proofs of Theorem 4.1, Proposition 4.2 and Theorem 4.3.

**Lemma B.1.** *If  $\Re(z) \leq 0$ , then*

$$\left| \frac{e^z - 1}{z} \right| \leq 1.$$

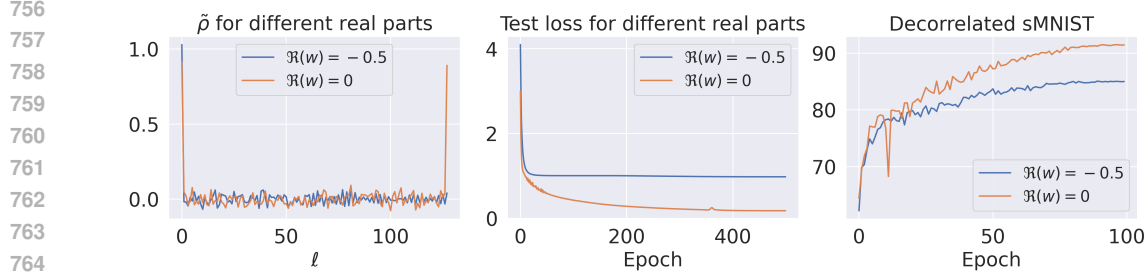


Figure 11: (Left) Training a diagonal SSM (3) with S4D-Legs initialization on a task that requires long-term memory. The learned memory function  $\tilde{\rho}$  effectively captures the spike in long-range dependencies. However, it struggles to do so when the real part is negative. (Middle) Test loss on the long-term memory task when initializing  $\Re(w) = 0$  and  $\Re(w) = -0.5$ . (Right) Test accuracy for training a diagonal SSM with S4D-Legs initialization on decorrelated sequential MNIST dataset with different real parts at initialization.

*Proof.* Notice that

$$\frac{|e^z - 1|}{|z|} = \frac{|\int_0^z e^s ds|}{|z|} \leq \frac{\int_0^z |e^s| |ds|}{|z|} = \frac{\int_0^z e^{\Re(z)} |ds|}{|z|} \leq \frac{\int_0^z |ds|}{|z|} = 1,$$

which finishes the proof.  $\square$

**Lemma B.2** (Itô's isometry). *Let  $W : [0, T] \times \Omega \rightarrow \mathbb{R}$  denote the canonical real-valued Wiener process defined up to time  $T > 0$ , and let  $X : [0, T] \times \Omega \rightarrow \mathbb{R}$  be a stochastic process that is adapted to the natural filtration of the Wiener process. Then*

$$\mathbb{E} \left[ \left( \int_0^T X_t dW_t \right)^2 \right] = \mathbb{E} \left[ \int_0^T X_t^2 dt \right],$$

where  $\mathbb{E}$  denotes expectation with respect to classical Wiener measure.

**Lemma B.3** (Gershgorin circle theorem). *Let  $A$  be a complex  $n \times n$  matrix, with entries  $a_{ij}$ . For  $i \in \{1, \dots, n\}$ , let  $R_i$  be the sum of the absolute value of the non-diagonal entries in the  $i$ -th row:  $R_i = \sum_{j \neq i} |a_{ij}|$ . Let  $D(a_{ii}, R_i) \subseteq \mathbb{C}$  be a closed disc centered at  $a_{ii}$  with radius  $R_i$ . Then every eigenvalue of  $A$  lies within at least one of the discs  $D(a_{ii}, R_i)$ .*

**Lemma B.4.** *For any  $t \in \mathbb{R}$ ,*

$$\sum_{n=1}^{\infty} \frac{1}{n^2 + t^2} = -\frac{1}{2t^2} + \frac{\pi}{2t} \coth(\pi t).$$

*Proof.* This is a side result of the Basel problem. The related proof can be found in the [Wiki page](#). We omit it here.  $\square$

**Lemma B.5.** *For any  $v_j, v_k \in \mathbb{R}$ , we have*

$$\int_0^{\infty} e^{-s} \cos(v_j s) \cos(v_k s) ds = \frac{1}{2} \left( \frac{1}{1 + (v_j - v_k)^2} + \frac{1}{1 + (v_j + v_k)^2} \right).$$

810 *Proof.* Notice that

$$\begin{aligned}
811 & \\
812 & \int_0^\infty e^{-s} \cos(v_j s) \cos(v_k s) ds \\
813 & \\
814 & = \frac{1}{2} \int_0^\infty e^{-s} \cos((v_j - v_k)s) ds + \frac{1}{2} \int_0^\infty e^{-s} \cos((v_j + v_k)s) ds \\
815 & \\
816 & = \frac{1}{2} \int_0^\infty \Re(\exp(-s + i \cdot (v_j - v_k)s)) ds + \frac{1}{2} \int_0^\infty \Re(\exp(-s + i \cdot (v_j + v_k)s)) ds \\
817 & \\
818 & = \frac{1}{2} \Re \left( \frac{1}{1 - i \cdot (v_j - v_k)} + \frac{1}{1 - i \cdot (v_j + v_k)} \right) \\
819 & \\
820 & = \frac{1}{2} \left( \frac{1}{1 + (v_j - v_k)^2} + \frac{1}{1 + (v_j + v_k)^2} \right). \\
821 & \\
822 & \\
823 & \\
824 & \square
\end{aligned}$$

## 826 C PROOF OF THEOREM 4.1

828 In this section, we prove the upper bound on the second moment of the model output value in  
829 Theorem 4.1.

830 *Proof.* First, we may express the model output  $y_L$  in a matrix form. To do so, we rewrite  $c$  as a  
831  $2m \times 1$  vector  $(\Re(c_1), \dots, \Re(c_m), \Im(c_1), \dots, \Im(c_m))^\top$  that contains the real and imaginary part of  
832  $c$ , and let  $V$  to be a  $2m \times L$  Vandermonde-like matrix

$$\begin{aligned}
834 & \\
835 & V := \begin{bmatrix} \Re\left(\frac{e^{\Delta w_1 - 1}}{\Delta w_1} e^{\Delta w_1 0}\right) & \Re\left(\frac{e^{\Delta w_1 - 1}}{\Delta w_1} e^{\Delta w_1 1}\right) & \dots & \Re\left(\frac{e^{\Delta w_1 - 1}}{\Delta w_1} e^{\Delta w_1 (L-1)}\right) \\
836 & \vdots & \vdots & \vdots \\
837 & \Re\left(\frac{e^{\Delta w_m - 1}}{\Delta w_m} e^{\Delta w_m 0}\right) & \Re\left(\frac{e^{\Delta w_m - 1}}{\Delta w_m} e^{\Delta w_m 1}\right) & \dots & \Re\left(\frac{e^{\Delta w_m - 1}}{\Delta w_m} e^{\Delta w_m (L-1)}\right) \\
838 & -\Im\left(\frac{e^{\Delta w_1 - 1}}{\Delta w_1} e^{\Delta w_1 0}\right) & -\Im\left(\frac{e^{\Delta w_1 - 1}}{\Delta w_1} e^{\Delta w_1 1}\right) & \dots & -\Im\left(\frac{e^{\Delta w_1 - 1}}{\Delta w_1} e^{\Delta w_1 (L-1)}\right) \\
839 & \vdots & \vdots & \vdots & \vdots \\
840 & -\Im\left(\frac{e^{\Delta w_m - 1}}{\Delta w_m} e^{\Delta w_m 0}\right) & -\Im\left(\frac{e^{\Delta w_m - 1}}{\Delta w_m} e^{\Delta w_m 1}\right) & \dots & -\Im\left(\frac{e^{\Delta w_m - 1}}{\Delta w_m} e^{\Delta w_m (L-1)}\right) \end{bmatrix}. \\
841 & \\
842 & \\
843 &
\end{aligned}$$

844 Then  $y_L$  can be written in a matrix form

$$845 y_L = \Delta \cdot c^\top V J x,$$

846 where  $J \in \mathbb{R}^{L \times L}$  is a row reversed identity matrix, i.e.

$$847 J = \begin{pmatrix} 0 & \dots & 0 & 1 \\ 0 & \dots & 1 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 1 & \dots & 0 & 0 \end{pmatrix}.$$

848 Furthermore, we may connect  $V$  with a standard Vandermonde matrix  $V_L$ , by noticing that

$$849 \Phi V = D V_L,$$

850 where  $V_L$  is a  $2m \times L$  complex Vandermonde matrix with  $2m$  nodes  $e^{\Delta w_1}, e^{\Delta \bar{w}_1}, \dots, e^{\Delta w_m}, e^{\Delta \bar{w}_m}$ :

$$\begin{aligned}
851 & \\
852 & \\
853 & V_L = \begin{pmatrix} 1 & e^{\Delta \bar{w}_1} & \dots & e^{\Delta \bar{w}_1 (L-1)} \\
854 & \vdots & \vdots & \vdots \\
855 & 1 & e^{\Delta \bar{w}_m} & \dots & e^{\Delta \bar{w}_m (L-1)} \\
856 & 1 & e^{\Delta w_1} & \dots & e^{\Delta w_1 (L-1)} \\
857 & \vdots & \vdots & \vdots & \vdots \\
858 & 1 & e^{\Delta w_m} & \dots & e^{\Delta w_m (L-1)} \end{pmatrix} \in \mathbb{C}^{2m \times L}, \\
859 & \\
860 & \\
861 & \\
862 & \\
863 &
\end{aligned}$$



When  $w \in \mathbb{C}^m$  with distinct imaginary parts, we can always find a scaling factor  $\gamma > 0$  such that  $e^{\gamma w_1}, \dots, e^{\gamma w_m}, e^{\gamma \bar{w}_1}, \dots, e^{\gamma \bar{w}_m}$  are distinct, where  $\bar{w}$  is the conjugate of  $w$ . Then by the argument of Vandermonde matrix, the only solution of the equation  $\sum_{j=1}^m \xi_j e^{w_j s} + \sum_{j=1}^n \hat{\xi}_j e^{\bar{w}_j s} = 0$  for  $s \geq 0$  is that  $\xi_j = \hat{\xi}_j = 0$  for  $j = 1, \dots, m$ . Since  $2\Re(e^{w_j s}) = e^{w_j s} + e^{\bar{w}_j s}$ , then  $\sum_{j=1}^m \xi_j \Re(e^{w_j s}) = 0$  only has zero solution.

Combining these two cases we finish the proof.  $\square$

## E PROOF OF THEOREM 4.3

In this section, we prove Theorem 4.3 based on the Gershgorin circle theorem (Lemma B.3).

*Proof.* First, we need to bound both the diagonal entry and the off-diagonal sum. The diagonal entry  $G_{j,j} = \frac{1}{2}(1 + \frac{1}{1+4v_j^2})$ , which can be bounded as

$$\frac{1}{2} \left( 1 + \frac{1}{1+4v_j^2} \right) \leq G_{j,j} \leq 1, \quad j = 1, \dots, m.$$

For the off-diagonal sum, we have  $\forall j = 1, \dots, m$ ,

$$\begin{aligned} 2R_j &= 2 \sum_{k \neq j} |G_{j,k}| \\ &= \sum_{k \neq j} \frac{1}{1 + (v_j - v_k)^2} + \sum_{k \neq j} \frac{1}{1 + (v_j + v_k)^2} \\ &< \sum_{k=1}^{\infty} \frac{2}{1 + \delta^2 k^2} + \sum_{k=1}^{\infty} \frac{1}{1 + (v_j + v_k)^2} - \frac{1}{1 + 4v_j^2} \\ &< \sum_{k=1}^{\infty} \frac{2}{1 + \delta^2 k^2} + \sum_{k=1}^{\infty} \frac{1}{1 + v_j^2 + v_k^2} - \frac{1}{1 + 4v_j^2} \\ &< \sum_{k=1}^{\infty} \frac{2}{1 + \delta^2 k^2} + \sum_{k=0}^{\infty} \frac{1}{1 + v_j^2 + \delta^2 k^2} - \frac{1}{1 + 4v_j^2} \\ &= \frac{2}{\delta^2} \sum_{k=1}^{\infty} \frac{1}{1/\delta^2 + k^2} + \frac{1}{\delta^2} \sum_{k=1}^{\infty} \frac{1}{(1 + v_j^2)/\delta^2 + k^2} + \left( \frac{1}{1 + v_j^2} - \frac{1}{1 + 4v_j^2} \right), \end{aligned}$$

where the first inequality is due to the fact that the minimal separation distance  $\min_{j \neq k} |v_j - v_k| \geq \delta$ , and the last inequality is because  $v_j > 0$  and reordering  $\{v_k\}_{k \geq 1}$  does not affect the result for  $\sum_{k=1}^{\infty} \frac{1}{1+v_j^2+v_k^2}$ . Then by Lemma B.4, we have

$$\begin{aligned} \frac{2}{\delta^2} \sum_{k=1}^{\infty} \frac{1}{1/\delta^2 + k^2} + \frac{1}{\delta^2} \sum_{k=1}^{\infty} \frac{1}{(1 + v_j^2)/\delta^2 + k^2} &< \frac{3}{\delta^2} \sum_{k=1}^{\infty} \frac{1}{1/\delta^2 + k^2} \\ &= \frac{3}{\delta^2} \left( -\frac{\delta^2}{2} + \frac{\pi\delta}{2} \coth\left(\frac{\pi}{\delta}\right) \right) \\ &= -\frac{3}{2} + \frac{3\pi}{2\delta} \coth\left(\frac{\pi}{\delta}\right). \end{aligned}$$

Hence we have,

$$\begin{aligned} G_{j,j} - R_j &> \frac{1}{2} \left( 1 + \frac{1}{1+4v_j^2} \right) - \frac{1}{2} \left( -\frac{3}{2} + \frac{3\pi}{2\delta} \coth\left(\frac{\pi}{\delta}\right) \right) - \frac{1}{2} \left( \frac{1}{1+v_j^2} - \frac{1}{1+4v_j^2} \right) \\ &> \frac{5}{4} - \frac{1}{2} \max\left( \frac{1}{1+x^2} - \frac{2}{1+4x^2} \right) - \frac{3\pi}{4\delta} \coth\left(\frac{\pi}{\delta}\right) \\ &> 1.19 - \frac{3\pi}{4\delta} \coth\left(\frac{\pi}{\delta}\right). \end{aligned}$$



972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

Under the same argument, we get

$$\begin{aligned}
 G_{j,j} + R_j &< 1 + \frac{1}{2} \left( -\frac{3}{2} + \frac{3\pi}{2\delta} \coth\left(\frac{\pi}{\delta}\right) \right) + \frac{1}{2} \max\left( \frac{1}{1+v_j^2} - \frac{1}{1+4v_j^2} \right) \\
 &< \frac{1}{4} + \frac{3\pi}{4\delta} \coth\left(\frac{\pi}{\delta}\right) + \frac{1}{2} \max\left( \frac{1}{1+v_j^2} - \frac{1}{1+4v_j^2} \right) \\
 &= \frac{5}{12} + \frac{3\pi}{4\delta} \coth\left(\frac{\pi}{\delta}\right).
 \end{aligned}$$

Combining the two bounds and Lemma B.3, we finish the proof. □