
On the Parameterization and Initialization of Diagonal State Space Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 State space models (SSM) have recently been shown to be very effective as a deep
2 learning layer as a promising alternative to sequence models such as RNNs, CNNs,
3 or Transformers. The first version to show this potential was the S4 model, which
4 is particularly effective on tasks involving long-range dependencies by using a
5 prescribed state matrix called the HiPPO matrix. While this has an interpretable
6 mathematical mechanism for modeling long dependencies, it also requires a custom
7 representation and algorithm that makes the model difficult to understand and
8 implement. On the other hand, a recent variant of S4 called DSS showed that
9 restricting the state matrix to be fully diagonal can still preserve the performance
10 of the original model when using a specific initialization based on approximating
11 S4’s matrix. This work seeks to systematically understand how to parameterize and
12 initialize diagonal state space models. While it follows from classical results that
13 almost all SSMs have an equivalent diagonal form, we show that the initialization
14 is critical for performance. First, we explain why DSS works mathematically, as
15 the diagonal approximation to S4 surprisingly recovers the same dynamics in the
16 limit of infinite state dimension. We then systematically describe various design
17 choices in parameterizing and computing diagonal SSMs, and perform a controlled
18 empirical study ablating the effects of these choices. Our final model S4D is a
19 simple diagonal version of S4 whose kernel computation requires just 3 lines of
20 code and performs comparably to S4 in almost all settings, with state-of-the-art
21 results in image, audio, and medical time-series domains, and 85% average on the
22 Long Range Arena benchmark.

23 1 Introduction

24 A core problem in deep learning is sequence modeling, which is typically characterized by defining a
25 parameterized sequence-to-sequence function that can be composed. Recent approaches based on
26 state space models (SSMs) have shown promise to outperform traditional deep sequence models
27 such as recurrent neural networks (RNNs), convolutional neural networks (CNNs), and Transformers,
28 in both computational efficiency and modeling ability. In particular, the S4 model displayed strong
29 results on a range of sequence modeling tasks, especially on long sequences [7]. Its ability to model
30 long-range dependencies arises from being defined with a particular state matrix called the “HiPPO
31 matrix” [4], which allows S4 to be viewed as a convolutional model that decomposes an input onto
32 an orthogonal system of smooth basis functions (Fig. 1).

33 However, beyond its theoretical interpretation, actually computing S4 as a deep learning model
34 requires a sophisticated algorithm with many linear algebraic techniques that are difficult to understand
35 and implement. These techniques were necessitated by parameterizing the state matrix as a **diagonal
36 plus low-rank** (DPLR) matrix, which is necessary to capture HiPPO matrices. A natural question is
37 whether simplifications of this parameterization and algorithm are possible. In particular, removing

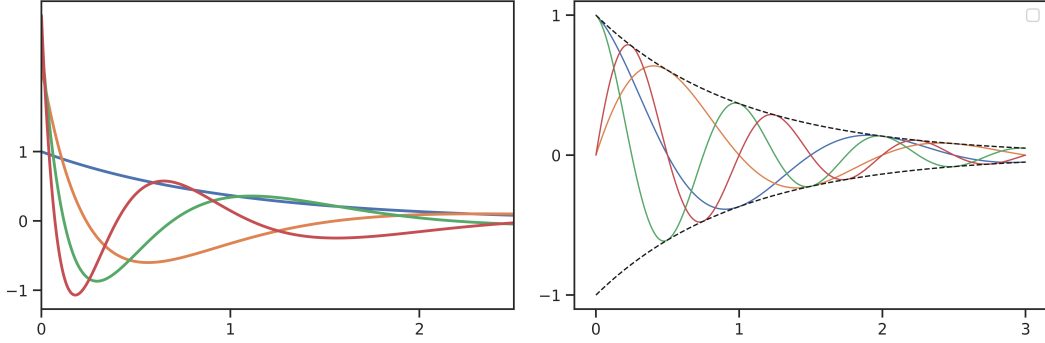


Figure 1: State space models can be interpreted as convolutional models where the convolution kernel is a linear combination of basis functions. (Left) The particular \mathbf{A} matrix chosen in S4 results in smooth basis functions with a closed form formula in terms of Legendre polynomials, resulting in a simple mathematical interpretation of the method as orthogonalizing against an exponentially-decaying measure. This basis grants long-range modeling properties, with the tradeoff of being difficult to compute the system with this matrix. (Right) We analyze much simpler state spaces with diagonal matrices. We show that it is critical to carefully initialize both the real parts which control the decay (dashed line), and the imaginary parts which control the frequencies (solid lines), allowing our diagonal SSM S4D to match S4 in nearly all settings.

38 the low-rank term would result in a **diagonal state space model** (DSSM) that is dramatically simpler
 39 to implement and understand.

40 Although it is known that almost all SSMs have an equivalent diagonal form—and therefore (complex)
 41 DSSMs are fully expressive algebraically—they may not represent all SSMs numerically, and finding
 42 a good initialization is critical. Gu et al. [7] showed that it is difficult to find a performant diagonal
 43 SSM, and that many alternative parameterizations of the state matrix – including by random diagonal
 44 matrices – are much less effective empirically, which motivated the necessity of the more complicated
 45 HiPPO matrix. However, recently Gupta [8] made the surprising empirical observation that a variant
 46 of S4 using a *particular diagonal matrix* is nearly as effective as the original S4 method. This matrix
 47 is based on the original HiPPO matrix and is defined by simply chopping off the low-rank term in
 48 the DPLR representation. The efficacy of this particular matrix compared to other diagonal matrices
 49 remains theoretically unexplained.

50 In this work, we seek to systematically understand how to train DSSMs.

- 51 • First, we categorize and consolidate various possible representations of diagonal SSMs, allow-
 52 ing the various design choices to be systematically compared – in particular the discretization,
 53 parameterization, and normalization of SSM models (Section 3).
- 54 • We theoretically provide a new mathematical analysis of DSS to understand why it works, showing
 55 that the diagonal matrix they chose based on the original HiPPO matrix is surprisingly an exact
 56 approximation of the dynamics of S4 as the state size goes to infinity. We propose even simpler
 57 variants of diagonal SSMs using different initializations of the state matrix. Our class of **S4D**
 58 methods compute the same function as S4 in the case of diagonal instead of DPLR state matrices,
 59 allowing dramatically simpler implementation (Section 4).
- 60 • We perform a controlled study of these various design choices across various settings, tasks, and
 61 sequence lengths, with both diagonal (S4D) and DPLR (S4) variants. Our best S4D methods are
 62 competitive with S4 on most settings, with near state-of-the-art results on image, audio, and medical
 63 time series benchmarks, and achieving **85%** on the Long Range Arena benchmark (Section 5).

64 2 Background

65 **Continuous State Spaces Models** S4 investigated state space models of the form that map signals
 66 $u(t) \mapsto y(t)$, which are linear time-invariant systems that can be represented either as a linear ODE
 67 (equation (1)) or convolution (equation (2)).

$$68 \quad \begin{aligned} x'(t) &= \mathbf{A}x(t) + \mathbf{B}u(t) \\ y(t) &= \mathbf{C}x(t) \end{aligned} \quad (1)$$

$$\begin{aligned} K(t) &= \mathbf{C}e^{t\mathbf{A}}\mathbf{B} \\ y(t) &= (K * u)(t) \end{aligned} \quad (2)$$

69 An intuitive way to view the convolution kernel (2) is to interpret it as a linear combination (controlled
70 by \mathbf{C}) of **basis kernels** $K_n(t)$ (controlled by \mathbf{A}, \mathbf{B})

$$K(t) = \sum_{k=0}^{N-1} \mathbf{C}_k K_k(t) \quad K_n(t) := (e^{t\mathbf{A}}\mathbf{B})_n \quad (3)$$

71 We denote this basis as $K_{\mathbf{A},\mathbf{B}}(t) = e^{t\mathbf{A}}\mathbf{B}$ if necessary to disambiguate. In the case of DSSMs, we
72 will denote $\mathbf{A} = \text{diag}(\mathbf{a})$ for $\mathbf{a} \in \mathbb{C}^N$. Note that in this simple case $K_n(t)$ is just $e^{t\mathbf{a}_n}\mathbf{B}_n$.

73 **S4: Structured State Spaces** As a deep learning model, SSMs have many elegant properties with
74 concrete empirical and computational benefits [6]. For example, the convolutional form (2) can be
75 converted into a temporal recurrence that is substantially faster for autoregressive applications [3].

76 However, making SSMs work required overcoming two key challenges: choosing appropriate values
77 for the matrices, and computing the kernel (2) efficiently.

78 First, Gu et al. [6] showed that naive instantiations of the SSM do not perform well, and instead relied
79 on a particular (real-valued) matrix \mathbf{A} called the HiPPO-LegS matrix (4). This matrix has a mathe-
80 matical interpretation of decomposing the input signal $u(t)$ onto a set of orthogonal basis functions
81 $L_n(e^{-t})$ (where $L_n(t)$ are normalized Legendre polynomials) with respect to an exponentially-
82 decaying measure, giving it long-range modeling abilities [0].¹

83 Second, S4 introduced a particular parameterization that decomposed this \mathbf{A} matrix into the sum of a
84 normal and rank-1 matrix (5), which can be unitarily conjugated into a (complex) diagonal plus rank-1
85 matrix. Leveraging this structured form, they then introduced a sophisticated algorithm for efficiently
86 computing the convolution kernel (2) for state matrices that are **diagonal plus low-rank (DPLR)**.

$$\begin{aligned} \mathbf{A}_{nk} = - \begin{cases} (2n+1)^{\frac{1}{2}}(2k+1)^{\frac{1}{2}} & n > k \\ n+1 & n = k \\ 0 & n < k \end{cases} \quad (4) \quad \mathbf{A}_{nk}^N = - \begin{cases} (n+\frac{1}{2})^{1/2}(k+\frac{1}{2})^{1/2} & n > k \\ \frac{1}{2} & n = k \\ (n+\frac{1}{2})^{1/2}(k+\frac{1}{2})^{1/2} & n < k \end{cases} \quad (5) \\ \mathbf{B}_n = \mathbf{P}_n = (2n+1)^{\frac{1}{2}} \quad \mathbf{A} = \mathbf{A}^N - \mathbf{P}\mathbf{P}^\top, \quad \mathbf{A}^D := \text{diag}(\mathbf{A}^N) \\ \text{(HiPPO-LegS matrix used in S4)} \quad \text{(Normal (diagonal) plus low-rank form)} \end{aligned}$$

88 **DSS: Diagonal State Spaces** S4 was originally motivated by searching for a *diagonal state matrix*,
89 which would be even more structured and result in very simple computation of the SSM. However,
90 the HiPPO-LegS matrix cannot be stably transformed into diagonal form [7, Lemma 3.2], and they
91 were unable to find any diagonal matrices that worked well, resulting in the DPLR formulation.

92 Gupta [8] made the surprising empirical observation that simply removing the low-rank portion of
93 the DPLR form of the HiPPO-LegS matrix results in a diagonal matrix that performs comparably
94 to the original S4 method. More precisely, their initialization is the diagonal matrix \mathbf{A}^D (5).

95 In addition to this initialization of a diagonal state matrix \mathbf{A} , they proposed a method for computing
96 the resulting state space. In Sections 3 and 4, we systematically study the components of DSS. We
97 categorize different ways to parameterize and compute the diagonal state space. Most importantly,
98 we show how to theoretically interpret this particular diagonal \mathbf{A} matrix.

99 3 Parameterizing Diagonal State Spaces

100 We describe various choices for the computation and parameterization of diagonal state spaces. Our
101 categorization of these choices leads to simple variants of the core method. Both DSS and our
102 proposed S4D can be described using a combination of these factors (Section 3.5).

103 3.1 Discretization

104 The true continuous-time SSM computes $y(t) = (K * u)(t) = \int_0^\infty \mathbf{C}e^{s\mathbf{A}}\mathbf{B}u(t-s) ds$.

105 In discrete time, we view an input sequence u_0, u_1, \dots as uniformly-spaced samples from an underly-
106 ing function $u(t)$ and must approximate this integral. Standard methods for doing so that preserve the

¹The citation [0] refers to a concurrent submission that is included in the supplemental.

107 convolutional structure of the model exist. The first step is to discretize the parameters. Two simple
 108 choices that have been used in prior work include

$$\begin{aligned} \text{(Bilinear)} \quad \overline{\mathbf{A}} &= (\mathbf{I} - \Delta/2\mathbf{A})^{-1}(\mathbf{I} + \Delta/2\mathbf{A}) & \text{(ZOH)} \quad \overline{\mathbf{A}} &= \exp(\Delta\mathbf{A}) \\ \overline{\mathbf{B}} &= (\mathbf{I} - \Delta/2\mathbf{A})^{-1} \cdot \Delta\mathbf{B} & \overline{\mathbf{B}} &= (\Delta\mathbf{A})^{-1}(\exp(\Delta \cdot \mathbf{A}) - \mathbf{I}) \cdot \Delta\mathbf{B}. \end{aligned}$$

109 With these methods, the discrete-time SSM output is just

$$y = u * \overline{\mathbf{K}} \quad \text{where } \overline{\mathbf{K}} = (\overline{\mathbf{C}\mathbf{B}}, \overline{\mathbf{C}\mathbf{A}\mathbf{B}}, \dots, \overline{\mathbf{C}\mathbf{A}^{L-1}\mathbf{B}}). \quad (6)$$

110 These integration rules have both been used in prior works (e.g. LMU and DSS use ZOH [18, 8])
 111 while S4 and its predecessors use bilinear [4, 6, 7]).

112 In Section 5, we show that there is little empirical difference between them. However, we note that
 113 there is a curious phenomenon where the bilinear transform actually perfectly smooths out the kernel
 114 used in DSS to match the S4 kernel (Section 4 Fig. 2c). We additionally note that numerical integration
 115 is a rich and well-studied topic and more stable methods of approximating the convolutional integral
 116 may exist. For example, it is well-known that simple rules like the Trapezoid rule [12] can dramatically
 117 reduce numerical integration error when the function has bounded second derivative.

118 3.2 Convolution Kernel

119 The main computational difficulty of the original S4 model is computing the convolution kernel $\overline{\mathbf{K}}$.
 120 This is extremely slow for general state matrices \mathbf{A} , and S4 introduced a complicated algorithm for
 121 DPLR state matrices. In the diagonal case, the computation is nearly trivial. Note that

$$\overline{\mathbf{K}}_\ell = \sum_{n \in [N]} \overline{\mathbf{B}}_n \mathbf{a}_n^\ell \mathbf{C}_n = (\overline{\mathbf{B}} \circ \mathbf{C}) \cdot \mathcal{V}_L(\mathbf{a}) \quad \mathcal{V}_L(\mathbf{a})_{n,\ell} = \mathbf{a}_n^\ell \quad (7)$$

122 where \circ is Hadamard product, \cdot is matrix multiplication, and \mathcal{V} is known as a **Vandermonde matrix**.

123 3.3 Parameterization

124 **Parameterization of \mathbf{A} .** Note that the kernel $K(t) = \mathbf{C}e^{t\mathbf{A}}\mathbf{B}$ blows up to ∞ as $t \rightarrow \infty$ if \mathbf{A} has
 125 any eigenvalues with positive real part. Goel et al. [3] found that this is a serious constraint that
 126 affects the stability of the model, especially when using the SSM autoregressively. They propose
 127 to force the real part of \mathbf{A} to be negative, also known as the classic left-half plane condition, by
 128 parameterizing the real part inside an exponential function $\mathbf{A} = -\exp(\mathbf{A}_r) + i \cdot \mathbf{A}_i$.

129 We note that instead of \exp , other choices can be made such as a simple clamp (i.e. ReLU).
 130 Alternatively, the original DSS uses no constraint on the real part of \mathbf{A} , which is sufficient for simple
 131 tasks involving fixed-length sequences, but would become unstable in other settings.

132 **Parameterization of \mathbf{B}, \mathbf{C} .** Note that the computation of the final discrete convolution kernel $\overline{\mathbf{K}}$
 133 depends only on the elementwise product $\mathbf{B} \circ \mathbf{C}$ (equation (7)). Therefore we can choose to store
 134 this product as the parameter, instead of \mathbf{B} and \mathbf{C} individually.

135 3.4 Normalization

136 An issue that requires consideration for SSMs is variance preservation. A simple proxy to ensure
 137 that the SSM is well behaved is [0]:

138 *For constant input $u(t) = c$, the output $y(t)$ should have mean 0, var. c^2 (w.r.t. the initialization)*

139 While this is satisfied automatically for HiPPO SSMs such as S4 [0], it is difficult to ensure for
 140 DSSMs. To address this, DSS used a **softmax** in their kernel which explicitly normalized every basis
 141 function $K_n(t)$ of the SSM (3) during every computation of the kernel. This has several drawbacks,
 142 including making the kernel more complicated, and additionally has the drawback of being calibrated
 143 only to a specific target length – in other words, changing the length L of the target kernel requires
 144 changing the underlying parameters \mathbf{A} and \mathbf{B} .

145 We would like to avoid the softmax, but address the normalization issue. While an analog to the
 146 HiPPO normalization cannot be analytically derived in the DSSM case, we simply enforce the
 147 condition at initialization by scaling \mathbf{B} by a constant. Given an arbitrary diagonal $\mathbf{A} = \text{diag}(\mathbf{a})$

148 and vector \mathbf{B} , The n -th SSM basis function is $e^{t\mathbf{a}_n} \mathbf{B}_n$. With the constant input $u(t) = 1$, the n -th
 149 component of the state (as $t \rightarrow \infty$) is given by (8), and the full state vector $x(t)$ has norm (9)

$$150 \quad x_n(t) = \int_0^\infty e^{s\mathbf{a}_n} \mathbf{B}_n ds = -\frac{\mathbf{B}_n}{\mathbf{a}_n} \quad (8) \quad \|x(t)\|^2 = \left(\sum_{n=0}^{N-1} \left| \frac{\mathbf{B}_n}{\mathbf{a}_n} \right|^2 \right)^{\frac{1}{2}}. \quad (9)$$

151 Dividing \mathbf{B} by this scalar yields a proper initialization that passes the above test. We call this the
 152 **scalar init** normalization, since it simply normalizes by a scalar at initialization.

153 3.5 S4D: the Diagonal Version of S4

154 A key component of our exposition is disentangling the various choices possible in representing
 155 and computing a diagonal state space. With this categorization, different choices can be mixed and
 156 matched to define variants of the core method. We define S4D so that the model $u \mapsto y$ *exactly*
 157 *matches the output of S4 when the low-rank component is set to 0*. This allows us to compare the
 158 performance of diagonal vs. DPLR state spaces while controlling all other factors. (We also mention
 159 how the components can be swapped out with alternatives, e.g. more sophisticated discretizations, to
 160 define promising alternate methods, and leave this as a direction for future work.) In our empirical
 161 study in Section 5, we systematically ablate the effects of each of these components.

Method	Discretization	Parameterization of \mathbf{A}	... of \mathbf{B}, \mathbf{C}	Normalization	Initialization of \mathbf{A}
162 DSS	ZOH	no restriction	combine	softmax	HiPPO-LegS-D
S4D	Bilinear	exp real part	independent	scalar init	various
(S4D-Trap)	Trapezoid	exp real part	independent	scalar init	various

163 Compared to DSS, the parameterization of S4D is guaranteed to be stable even in unbounded
 164 autoregressive settings, and its computation is much simpler due to not requiring the softmax
 165 normalization; our proposed scalar init normalization (Section 3.4) attains most of the benefits. The
 166 entire S4D method is almost trivial to implement, requiring just a few lines of code each for the
 167 initialization, kernel computation, and forward pass (Listing 1).

168 4 Initialization of Diagonal State Matrices

169 The critical question that remains is: which diagonal state matrices \mathbf{A} actually work? We comment on
 170 the limitations of DSSMs, and then provide three instantiations of S4D that are empirically effective.

171 **Expressivity and Limitations of Diagonal SSMs.** We first present a simplified view on the expres-
 172 sivity of DSSMs mentioned by [8]. First, it is well-known that almost all matrices diagonalize over
 173 the complex plane. Therefore it is critical to use complex-valued matrices in order to use DSSMs.

174 **Proposition 1.** *The set $\mathcal{D} \subset \mathbb{C}^{N \times N}$ of diagonalizable matrices is dense in $\mathbb{C}^{N \times N}$, and has full
 175 measure (i.e. its complement has measure 0).*

176 It is also well known that the state space $(\mathbf{A}, \mathbf{B}, \mathbf{C})$ is exactly equivalent to (i.e. expresses the same
 177 map $u \mapsto y$) the state space $(\mathbf{V}^{-1}\mathbf{A}\mathbf{V}, \mathbf{V}^{-1}\mathbf{B}, \mathbf{C}\mathbf{V})$, known in the SSM literature as a state space
 178 transformation. Therefore Proposition 1 says that (almost) all SSMs are equivalent to a DSSM.

179 However, we emphasize that Proposition 1 is about *expressivity* which does not guarantee strong
 180 performance of a trained model after optimization. For example, Gu et al. [7] and Gupta [8] show
 181 that randomly initialized dense real matrices or diagonal complex matrices, which are both fully
 182 expressive, perform much worse than S4.

183 Second, Proposition 1 does not take into account numerical representations of data, which was
 184 the original reason S4 required a low-rank correction term instead of a pure diagonalization. In
 185 Section 5.2, we also show that two different representations of matrices initialized with the *same*
 186 *spectrum* (i.e., are equivalent to the same diagonal \mathbf{A}) can have vastly different performance.

187 **S4D-LegS.** Gupta [8] showed that taking the DPLR form of the HiPPO-LegS matrix $\mathbf{A}^D - \mathbf{P}\mathbf{P}^\top$
 188 and simply approximating it with \mathbf{A}^D works quite well (5). Our first result is providing a clean mathe-
 189 matical interpretation of this method. Theorem 2 shows a surprising fact that does not hold in general
 190 for DPLR matrices (Appendix A.2), and arises out of the special structure of this particular matrix.

191 **Theorem 2.** *Let $\mathbf{A} = \mathbf{A}^D - \mathbf{P}\mathbf{P}^\top$ and \mathbf{B} be the HiPPO-LegS matrices, and $K_{\mathbf{A}, \mathbf{B}}$ be its basis
 192 (Fig. 1). As the state size $N \rightarrow \infty$, the SSM basis $K_{\mathbf{A}^D, \mathbf{B}}$ is unitarily equivalent to $K_{\mathbf{A}, \mathbf{B}}$ (Fig. 2).*

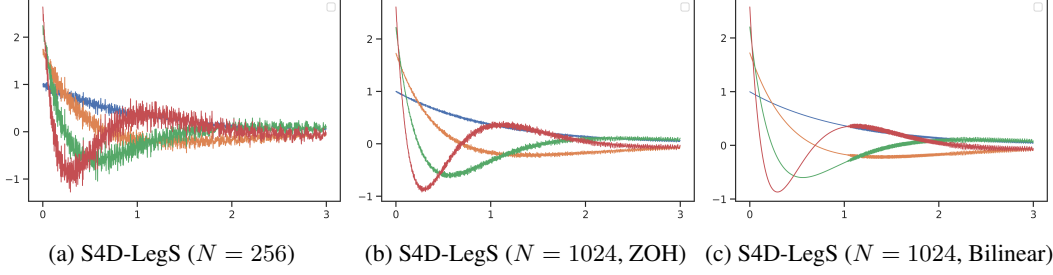


Figure 2: **(Visualization of Theorem 2)**. (a, b) By special properties of the original S4 matrix, removing its DPLR low-rank term produces the same basis functions as $N \rightarrow \infty$, explaining the empirical effectiveness of DSS. S4D-Inv is a simplified variant of S4D-LegS with approximate eigenvalues. (c) Curiously, the bilinear transform instead of ZOH smooths out the kernel to exactly match S4-LegS (Fig. 1 (Left)) as N grows.

193 (Note that *unitary* equivalence is important as it preserves the stability and timescale [0] of the system.)

194 We define **S4D-LegS** to be the S4D method for this choice of diagonal $\mathbf{A} = \mathbf{A}^D$ and \mathbf{B} . Theorem 2
 195 explains the empirical results in [8] whereby this system performed quite close to S4, but was
 196 generally consistently slightly worse. This is because DSS is a variant of S4D-LegS, which by
 197 Theorem 2 is a noisy approximation to S4-LegS. Fig. 2 illustrates this result, and also shows a curious
 198 phenomenon involving different discretization rules that is open for future work.

199 **S4D-Inv**. To further simplify S4D-LegS, we analyze the structure of $\mathbf{A}^D = \text{diag}(\mathbf{a})$ in more detail.
 200 The real part is easy to understand, which follows from the analysis in [7]: $\Re(\mathbf{a}) = -\frac{1}{2}\mathbf{1}$. Let
 201 the imaginary part be sorted, i.e. $\Im(\mathbf{\Lambda})_n$ is the n -th largest (positive) imaginary component. We
 202 empirically deduced the following conjecture for the asymptotics of the imaginary part.

203 **Conjecture 3.** As $N \rightarrow \infty$, $\Im(\mathbf{\Lambda})_0 \rightarrow \frac{1}{\pi}N^2 + c$ where $c \approx 0.5236$ is a constant. For a fixed N , the
 204 other eigenvalues satisfy an inverse scaling in n : $\Im(\mathbf{\Lambda})_n = \Theta(n^{-1})$.

205 Fig. 3 empirically supports this conjecture. Based on Conjecture 3, we propose the initialization
 206 S4D-Inv to use the following inverse-law diagonal matrix which closely approximates S4D-LegS.

207
$$\text{(S4D-Inv)} \quad \mathbf{\Lambda}_n = -\frac{1}{2} + i\frac{N}{\pi} \left(\frac{N}{2n+1} - 1 \right) \quad (10) \quad \text{(S4D-Lin)} \quad \mathbf{\Lambda}_n = -\frac{1}{2} + i\pi n \quad (11)$$

208 **S4D-Lin**. While S4D-Inv can be seen as an approximation to the original S4-LegS, we propose an
 209 even simpler scaling law for the imaginary parts that can be seen as an approximation of S4-FouT
 210 ([0]), where the imaginary parts are simply the Fourier series frequencies (.e. matches the diagonal
 211 part of the DPLR form of S4-FouT). Fig. 1 (Right) illustrates an idealized basis that S4D-Lin can
 212 produce, which are simply damped Fourier basis functions.

213 **General DSSM Basis Functions.** The empirical study in Section 5 performs many ablations of
 214 different diagonal initializations, showing that many natural variants of the proposed methods do not
 215 perform as well. The overall guiding principles for the diagonal state matrix $\mathbf{\Lambda}$ are twofold, which
 216 can be seen from the closed form of the basis functions $K_n(t) = e^{t\mathbf{a}_n} \mathbf{B}_n$ (Eq. (3)).

217 First, the real part of \mathbf{a}_n controls the decay rate of the function. $\mathbf{a}_n = -\frac{1}{2}$ is a good default that
 218 bounds the basis functions by the envelope $e^{-\frac{t}{2}}$, giving a constant timescale (Fig. 1 (Right)).

219 Second, the imaginary part of \mathbf{a}_n controls the oscillating frequencies of the basis function. Critically,
 220 these need to be spread out, which explains why random initializations of $\mathbf{\Lambda}$ do not perform well.
 221 The S4D-Inv and S4D-Lin use simple asymptotics for these imaginary components that provide
 222 interpretable methods. We believe that alternative initializations that have different mathematical
 223 interpretations may exist, which is an interesting question for future work.

224 5 Experiments

225 Our experimental study shows that S4D has strong performance in a wide variety of domains and
 226 tasks, including the well-studied Long Range Arena (LRA) benchmark where the best S4D variant is
 227 competitive with S4 on all tasks and significantly performs all non-SSM baselines.

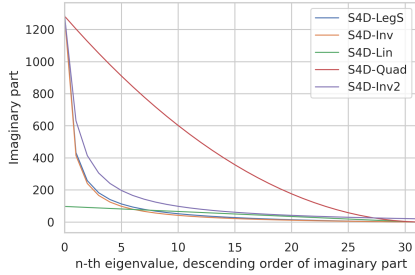


Figure 3: **(S4D eigenvalues.)** All S4D methods have eigenvalues $-\frac{1}{2} + \lambda_n i$. S4D-LegS theoretically approximates dynamics of the original (non-diagonal) S4 (Blue), and has eigenvalues following an inverse law $\lambda_n \propto n^{-1}$ (Orange). The precise law is important: other scaling laws with the same range, including an inverse law with different constant (Purple) and a quadratic law (Red), perform empirically worse (Section 5.2). A very different linear law based on Fourier frequencies also performs well (Green).

228 We begin with controlled ablations of the various representations of diagonal state space models. Sec-
 229 tions 5.1 and 5.2 ablate the proposed methods for parameterizing, computing, and initializing DSSMs
 230 from Sections 3 and 4. Section 5.3 show full results of larger models on standard benchmarks,

231 **Methodology and Datasets.** In order to study the effects of different S4 and S4D variants in a
 232 controlled setting, we propose the following protocol. We focus on three datasets covering a varied
 233 range of data modalities (image pixels, biosignal time series, audio waveforms), sequence lengths
 234 (1K, 4K, 16K), and tasks (classification and regression with bidirectional and causal models).

- 235 • **Sequential CIFAR (sCIFAR).** CIFAR-10 images are flattened into a sequence of length 1024, and
 236 a bidirectional sequence model is used to perform 10-way classification.
- 237 • **BIDMC Vital Signs.** EKG and PPG signals of length 4000 are used to predict respiratory rate
 238 (RR), heart rate (HR), and blood oxygen saturation (SpO2). We focus on SpO2 in this study.
- 239 • **Speech Commands (SC).**² A 1-second raw audio waveform comprising 16000 samples is used
 240 for 35-way spoken word classification. We use an autoregressive (AR) or *causal* model to examine
 241 different settings; the causal setting more closely imitates autoregressive speech generation, where
 242 SSMs have shown recent promise [3].

243 We fix a simple architecture and training protocol that works generically. The architecture has 4
 244 layers and hidden dimension $H = 128$, resulting in $\sim 100K$ parameters. All results are averaged
 245 over multiple seeds (full protocol and results including std. reported in Appendix B).

246 5.1 Parameterization, Normalization, Discretization

247 Given the same diagonal SSM matrices A, B , there are many variants of how to parameterize the
 248 matrices and compute the SSM kernel described in Section 3. We ablate several aspects of the
 249 parametrization, normalization, and discretization. Results are in Table 1, and show that:

- 250 (i) Different discretizations do not make a noticeable difference (Section 3.1).
- 251 (ii) Unrestricting the real part Section 3.3 may be slightly better at the expense of generation
 252 stability [3], but overall makes little difference.
- 253 (iii) Normalizing the model with either softmax or the simple scalar rescaling Section 3.4 provides
 254 a consistent improvement, validating the importance of proper normalization.

Discretization	Normalization	sCIFAR	SC (uni)	BIDMC (SpO2)
Bilinear	none	84.53	89.52	0.1204
ZOH	none	84.50	89.53	0.1341
Bilinear	scalar init	85.14	89.90	0.1090
ZOH	scalar init	84.84	90.19	0.1299
ZOH	softmax	85.14	90.36	0.1299

Discretization	Real part	sCIFAR	SC (uni)	BIDMC (SpO2)
Bilinear	Exp	85.20	89.52	0.1193
	None	85.35	90.58	0.1102
	Relu	85.06	90.22	0.1172
ZOH	Exp	85.02	89.93	0.1303
	None	85.15	90.19	0.1289
	Relu	84.98	90.03	0.1232

Table 1: (Left) Ablations on Discretization and Normalization of DSSMs. (Right) Ablations on Discretization and Parameterization of DSSMs.

255 These ablations justify the simpler parameterization and algorithm that we use for S4D (Section 3.5)
 256 compared to DSS. For the remaining of the experiments in Section 5.2 and Section 5.3, we fix the
 257 S4D parameterization and algorithm described in Section 3. Note that this computes exactly the

²We note that a line of prior work including S4 [10, 13, 7] all used a smaller 10-class subset of SC, so our numbers are not directly comparable.

258 same kernel as the original S4 algorithm when the low-rank portion is set to 0, allowing controlled
259 comparisons of the critical state matrix \mathbf{A} for the remainder of this section.

260 5.2 S4D Initialization Ablations

261 The original S4 model proposed a specific formula for the \mathbf{A} matrix, and the first diagonal version
262 [8] used a specific matrix based on it. Our new proposed variants S4D-Inv and S4D-Lin also define
263 precise formulas for the initialization of the \mathbf{A} matrix (10). This raises the question of whether the
264 initialization of the \mathbf{A} still needs to be so precise, despite the drastic simplifications from the original
265 version. We perform several natural ablations on these initializations, showing that even simple
266 variations of the precise formula will degrade performance.

267 **Imaginary part scaling factor.** The scaling rules for S4D-Inv and S4D-Lin are simple polynomial
268 laws, but how is the constant factor chosen and how important is it? These constants for the imaginary
269 parts of S4D initializations were chosen based on connections to S4 methods based on HiPPO. Note
270 that the range of imaginary components for S4D-Inv and S4D-Lin are quite different (Fig. 3); the
271 largest imaginary part is $\frac{N^2}{\pi}$ for S4D-Inv and πN for S4D-Lin.

272 We consider scaling all imaginary parts by a constant factor of 0.01 or 100.0 to investigate whether
273 the constant matters. Note that this preserves the overall shape of the basis functions (Fig. 1) and
274 simply changes the frequencies, and it is not obvious that this should degrade performance. However,
275 both changes substantially reduce the performance of both S4D methods in all settings.

276 **Randomly initialized imaginary part.** Next, we consider choosing the imaginary parts randomly.
277 For S4D-Inv, we keep the real parts equal to $-\frac{1}{2}$ and set each imaginary component to $\mathbf{a}_n =$
278 $-\frac{1}{2} + i\frac{N}{\pi}(\frac{N}{2u+1} - 1)$ for $u \sim N \cdot \text{Unif}(0, 1)$. Note that when u is equally spaced in $[0, 1]$ instead of
279 uniformly random, this exactly recovers S4D-Inv (10), so this is a sensible random approximation to it.
280 Similarly, we consider a variant of S4D-Lin that replaces the n in (11) with $N \cdot \text{Unif}(0, 1)$.

281 Table 2a (*Random Imag*) shows that this small change causes minor degradation in performance.
282 We additionally note that the randomly initialized imaginary ablation can be interpreted as follows.
283 Fig. 3 shows the asymptotics of the imaginary parts of SSM matrices, where the imaginary parts of
284 the eigenvalues correspond to y-values corresponding to uniformly spaced nodes on the x-axis. This
285 ablation then replaces the uniform spacing on the x-axis with uniformly random x values.

286 **Randomly initialized real part.** We consider initializing the real part of each eigenvalue as
287 $-\text{Unif}(0, 1)$ instead of fixing them to $-\frac{1}{2}$. Table 2a(Left, *Random Real*) shows that this also causes
288 minor but consistent degradation in performance on the diagnostic datasets. Finally, we also consider
289 randomizing both real and imaginary parts, which degrades performance even further.

290 **Ablation: Other S4D matrices.** Other simple variants of initializations (S4D-Inv2, Quad) (Fig. 3)
291 show that it is not just the range of the eigenvalues but the actual distribution that is important.

292 We include two additional methods here that are not based on the proposed S4D-Inv or S4D-Lin meth-
293 ods. First, S4D-Rand uses a randomly initialized diagonal \mathbf{A} , and validates that it performs poorly, in
294 line with earlier findings [7, 8]. Second, S4D-Real uses a particular real diagonal initialization with
295 $\Lambda_n = -(n + 1)$. This is the exact same spectrum as the original S4(-LegS) method, which validates
296 that it is not just the diagonalization that matters, highlighting the limitation of Proposition 1.

297 5.3 Full Comparisons of S4D and S4 Methods

298 **Trainable \mathbf{A}, \mathbf{B} matrices.** Table 2b shows the performance of all S4D and S4 variants [0] on
299 the diagnostics datasets. We observe several interesting phenomena: (i) Freezing the matrices
300 performs comparably to training them on sCIFAR and BIDMC, but is substantially worse on SC. We
301 hypothesize that this results from Δ being poorly initialized for SC, so that at initialization models
302 do not have context over the entire sequence, and training \mathbf{A} and \mathbf{B} helps adjust for this. As further
303 evidence, the *finite window methods* S4-LegT and S4-FouT (defined in HTTYH I) have the most
304 limited context and suffer the most when \mathbf{A} is frozen. (ii) The full DPLR versions are often slightly
305 better throughout training. We report the validation accuracy after 1 epoch of training on sCIFAR
306 and SC to illustrate this phenomenon.

Table 2: (Initialization and Trainability ablations)

Ablation	SCIFAR	SC (AR)	BIDMC	Frozen (A, B)	SCIFAR		SC (AR)		BIDMC
					Acc (first)	Acc (best)	Acc (first)	Acc (best)	RMSE (best)
S4D-Lin	85.12	90.66	0.128						
Scale 0.01	-7.27	-1.92	+0.040	S4-LegS	53.63	86.19	33.87	85.33	0.1049
Scale 100	-7.91	-4.04	+0.077	S4-LegT	54.76	86.30	8.77	57.35	0.1106
Random Imag	-0.42	-3.08	-0.001	S4-FouT	55.28	86.05	9.27	69.57	0.1072
Random Real	-0.73	-0.87	+0.011	S4-LegS+FouT	54.38	86.53	34.06	83.37	0.0887
Random Both	-1.28	-5.88	+0.007	S4D-LegS	50.87	84.81	22.76	77.18	0.0960
				S4D-Inv	53.19	84.40	18.49	76.53	0.0995
				S4D-Lin	51.75	84.96	19.09	75.58	0.0935
				Trainable (A, B)					
S4D-Inv	84.79	90.27	0.114	S4-LegS	54.23	86.29	62.19	90.68	0.1033
Scale 0.01	-5.03	-0.08	+0.028	S4-LegT	55.16	86.12	55.86	90.42	0.1146
Scale 100	-7.77	-52.31	+0.034	S4-FouT	55.89	85.93	60.56	90.83	0.1136
Random Imag	-0.29	-0.52	+0.010	S4-LegS+FouT	55.00	86.18	61.76	91.01	0.0970
Random Real	0.12	-2.18	+0.032	S4D-LegS	50.41	85.64	47.54	88.47	0.1148
Random Both	-1.55	-0.55	+0.024	S4D-Inv	53.42	84.59	45.73	89.69	0.1132
S4D-Inv2	-2.62	-39.84	+0.005	S4D-Lin	52.23	85.75	47.68	89.56	0.1032
S4D-Quad	-1.83	-0.62	+0.024						
S4D-Random	-6.32	-1.95	+0.034						
S4D-Real	-5.45	-10.17	+0.066						

(a) Ablations of the initialization of the diagonal \mathbf{A} matrix in S4D. Very simple changes that largely preserve the structure of the diagonal eigenvalues all degrade performance.

(b) Results for all S4 and S4D methods on the diagnostic datasets, when the \mathbf{A} and \mathbf{B} matrices are either frozen (*Top*) or trained (*Bottom*). Diagonal state matrices are highly competitive with full DPLR versions, achieving strong results on all datasets.

Table 3: (Diagnostic datasets: Full results with larger models.) For Speech Commands, we show both an autoregressive model as in the ablations, and an unconstrained bidirectional model.

MODEL	SCIFAR		SC		BIDMC		
	TEST	AR	BI.	HR	RR	SP02	
S4-LegS	91.80 (0.43)	93.60 (0.13)	96.08 (0.15)	0.332 (0.013)	0.247 (0.062)	0.090 (0.006)	
S4-FouT	91.22 (0.25)	91.78 (0.10)	95.42 (0.20)	<u>0.339</u> (0.020)	0.301 (0.030)	0.068 (0.003)	
S4D-LegS	89.92 (1.69)	93.57 (0.09)	95.97 (0.14)	0.367 (0.001)	0.248 (0.036)	0.102 (0.001)	
S4D-Inv	90.69 (0.06)	93.40 (0.67)	96.18 (0.27)	0.373 (0.024)	0.254 (0.022)	0.110 (0.001)	
S4D-Lin	90.42 (0.03)	93.37 (0.05)	96.25 (0.03)	0.379 (0.006)	0.226 (0.008)	0.114 (0.003)	

307 **Large models on diagnostic datasets.** Finally, we relax the strict requirements on model size and
308 regularization for the diagnostic datasets, and show the performance of S4 and S4D variants on the
309 test sets with a larger model (architecture and training details in Appendix B) when the model size
310 and regularization is simply increased (Table 3). We note that results for each dataset are better than
311 the original S4 model, which was already state-of-the-art on these datasets.

312 **Long Range Arena** We use the same hyperparameter setting for the state-of-the-art S4 model in [0]
313 on the Long Range Arena benchmark for testing long dependencies in sequence models. S4D variants
314 are highly competitive on all datasets except Path-X, and outperform the S4 variants on several of
315 them. On Path-X using this hyperparameter setting with bidirectional models, only S4D-Inv, our
316 simpler approximation to the original S4-LegS model, achieves above random chance, and has an
317 average of 85% on the full LRA suite, more than 30 points better than the original Transformer [16].

Table 4: (Long Range Arena) Accuracy on full suite of LRA tasks. Hyperparameters in Appendix B.

MODEL	LISTOPS	TEXT	RETRIEVAL	IMAGE	PATHFINDER	PATH-X	AVG
S4-LegS	59.60 (0.07)	86.82 (0.13)	90.90 (0.15)	<u>88.65</u> (0.23)	<u>94.20</u> (0.25)	96.76 (-)	86.16
S4-FouT	57.88 (1.90)	86.34 (0.31)	89.66 (0.88)	89.07 (0.19)	94.46 (0.24)	\times	77.90
S4D-LegS	<u>60.47</u> (0.34)	86.18 (0.43)	89.46 (0.14)	88.19 (0.26)	93.06 (1.24)	\times	77.89
S4D-Inv	60.18 (0.35)	87.34 (0.20)	91.09 (0.01)	87.83 (0.37)	93.78 (0.25)	<u>91.22</u> (-)	<u>85.24</u>
S4D-Lin	60.52 (0.51)	<u>86.97</u> (0.23)	<u>90.96</u> (0.09)	87.93 (0.34)	93.96 (0.60)	\times	78.39
S4 (original)	58.35	76.02	87.09	87.26	86.05	88.10	80.48
Transformer	36.37	64.27	57.46	42.44	71.40	\times	53.66

318 **References**

- 319 [1] Shaojie Bai, J Zico Kolter, and Vladlen Koltun. Trellis networks for sequence modeling. In *The*
320 *International Conference on Learning Representations (ICLR)*, 2019.
- 321 [2] N Benjamin Erichson, Omri Azencot, Alejandro Queiruga, Liam Hodgkinson, and Michael W
322 Mahoney. Lipschitz recurrent neural networks. In *International Conference on Learning*
323 *Representations*, 2021.
- 324 [3] Karan Goel, Albert Gu, Chris Donahue, and Christopher Ré. It’s raw! audio generation with
325 state-space models. *arXiv preprint arXiv:2202.09729*, 2022.
- 326 [4] Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent
327 memory with optimal polynomial projections. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia
328 Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information*
329 *Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020,*
330 *NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper/2020/hash/102f0bb6efb3a6128a3c750dd16729be-Abstract.html)
331 [cc/paper/2020/hash/102f0bb6efb3a6128a3c750dd16729be-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/102f0bb6efb3a6128a3c750dd16729be-Abstract.html).
- 332 [5] Albert Gu, Caglar Gulcehre, Tom Le Paine, Matt Hoffman, and Razvan Pascanu. Improving the
333 gating mechanism of recurrent neural networks. In *The International Conference on Machine*
334 *Learning (ICML)*, 2020.
- 335 [6] Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré.
336 Combining recurrent, convolutional, and continuous-time models with the structured learnable
337 linear state space layer. In *Advances in Neural Information Processing Systems (NeurIPS)*,
338 2021.
- 339 [7] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured
340 state spaces. In *The International Conference on Learning Representations (ICLR)*, 2022.
- 341 [8] Ankit Gupta. Diagonal state spaces are as effective as structured state spaces. *arXiv preprint*
342 *arXiv:2203.14343*, 2022.
- 343 [9] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):
344 1735–1780, 1997.
- 345 [10] Patrick Kidger, James Morrill, James Foster, and Terry Lyons. Neural controlled differential
346 equations for irregular time series. *arXiv preprint arXiv:2005.08926*, 2020.
- 347 [11] James Morrill, Cristopher Salvi, Patrick Kidger, James Foster, and Terry Lyons. Neural rough
348 differential equations for long time series. *The International Conference on Machine Learning*
349 *(ICML)*, 2021.
- 350 [12] Anthony Ralston and Philip Rabinowitz. *A first course in numerical analysis*. Courier Corpora-
351 tion, 2001.
- 352 [13] David W Romero, Anna Kuzina, Erik J Bekkers, Jakub M Tomczak, and Mark Hoogendoorn.
353 Ckconv: Continuous kernel convolution for sequential data. *arXiv preprint arXiv:2102.02611*,
354 2021.
- 355 [14] T Konstantin Rusch and Siddhartha Mishra. Unicornn: A recurrent model for learning very
356 long time dependencies. *The International Conference on Machine Learning (ICML)*, 2021.
- 357 [15] Chang Wei Tan, Christoph Bergmeir, Francois Petitjean, and Geoffrey I Webb. Time series
358 extrinsic regression. *Data Mining and Knowledge Discovery*, pages 1–29, 2021. doi: <https://doi.org/10.1007/s10618-021-00745-9>.
- 359
- 360 [16] Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng
361 Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena : A benchmark for
362 efficient transformers. In *International Conference on Learning Representations*, 2021. URL
363 <https://openreview.net/forum?id=qVyeW-grC2k>.
- 364 [17] Trieu H Trinh, Andrew M Dai, Minh-Thang Luong, and Quoc V Le. Learning longer-term
365 dependencies in RNNs with auxiliary losses. In *The International Conference on Machine*
366 *Learning (ICML)*, 2018.

367 [18] Aaron Voelker, Ivana Kajić, and Chris Eliasmith. Legendre memory units: Continuous-time
368 representation in recurrent neural networks. In *Advances in Neural Information Processing*
369 *Systems*, pages 15544–15553, 2019.

370 Checklist

- 371 1. For all authors...
- 372 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
373 contributions and scope? [Yes]
- 374 (b) Did you describe the limitations of your work? [Yes] See Section 5, e.g. not matching
375 the baseline S4 on Path-X.
- 376 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 377 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
378 them? [Yes]
- 379 2. If you are including theoretical results...
- 380 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 381 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix A.
- 382 3. If you ran experiments...
- 383 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
384 imental results (either in the supplemental material or as a URL)? [Yes] The code is
385 a simple modification from the original S4 [7] repository, which is publicly available.
386 All experimental details to reproduce results are in the Appendix.
- 387 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
388 were chosen)? [Yes] See Appendix B.
- 389 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
390 ments multiple times)? [Yes] See the Appendix.
- 391 (d) Did you include the total amount of compute and the type of resources used (e.g., type
392 of GPUs, internal cluster, or cloud provider)? [No] Our datasets were small scale, and
393 the Long Range Arena details match those of S4 and are publicly available.
- 394 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 395 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 396 (b) Did you mention the license of the assets? [Yes]
- 397 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 398 (d) Did you discuss whether and how consent was obtained from people whose data you’re
399 using/curating? [N/A]
- 400 (e) Did you discuss whether the data you are using/curating contains personally identifiable
401 information or offensive content? [N/A] We used standard benchmarks and synthetic
402 data.
- 403 5. If you used crowdsourcing or conducted research with human subjects...
- 404 (a) Did you include the full text of instructions given to participants and screenshots, if
405 applicable? [N/A]
- 406 (b) Did you describe any potential participant risks, with links to Institutional Review
407 Board (IRB) approvals, if applicable? [N/A]
- 408 (c) Did you include the estimated hourly wage paid to participants and the total amount
409 spent on participant compensation? [N/A]