

# STATE-SPACE MODELS CAN LEARN IN-CONTEXT BY GRADIENT DESCENT

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Deep state-space models (Deep SSMs) have shown capabilities for in-context learning on autoregressive tasks, similar to transformers. However, the architectural requirements and mechanisms enabling this in recurrent networks remain unclear. This study demonstrates that state-space model architectures can perform gradient-based learning and use it for in-context learning. We prove that a single structured state-space model layer, augmented with local self-attention, can reproduce the outputs of an implicit linear model with least squares loss after one step of gradient descent. Our key insight is that the diagonal linear recurrent layer can act as a gradient accumulator, which can be ‘applied’ to the parameters of the implicit regression model. We validate our construction by training randomly initialized augmented SSMs on simple linear regression tasks. The empirically optimized parameters match the theoretical ones, obtained analytically from the implicit model construction. Extensions to multi-step linear and non-linear regression yield consistent results. The constructed SSM encompasses features of modern deep state-space models, with the potential for scalable training and effectiveness even in general tasks. The theoretical construction elucidates the role of local self-attention and multiplicative interactions in recurrent architectures as the key ingredients for enabling the expressive power typical of foundation models.

## 1 INTRODUCTION

The current generation of Large Language Models (LLMs) and foundation models are extremely capable and have started proliferating in several real-world applications. These models are based on the transformer architecture (Vaswani et al., 2017), and a big part of their capability has been attributed to in-context learning (Wei et al., 2023; Lu et al., 2023). In-context learning in transformers is relatively well studied (Wies et al., 2023; Pan et al., 2023; Guo et al., 2023; Garg et al., 2023). A prominent explanation for the mechanism used by transformers to do in-context learning is that the model performs in-context learning by gradient descent (von Oswald et al., 2023; Akyürek et al., 2024). But, the quadratic dependence of transformers on the input length makes them computationally expensive. To mitigate this, there has been much work on alternatives based on recurrent networks, such as state-space models (Gu et al., 2021; Gu & Dao, 2023) and linear recurrent networks (Orvieto et al., 2023). These recurrent models can perform inference efficiently since the computational complexity of recurrent networks is linear in the sequence length. At the same time, linear recurrent networks allow parallelization across the sequence during training using associative scan. The latest versions of these models are competitive with transformers at scale (Dao & Gu, 2024; De et al., 2024), and also capable of in-context learning (Grazzi et al., 2024). However, the mechanism they use for in-context learning remains unclear.

This work focuses on the in-context learning mechanism of State-Space Models (SSMs). Multiple variations of state-space models have recently shown competitive performance at scale (Gu & Dao, 2023; De et al., 2024), while earlier generations struggled with scaling (Smith et al., 2022; Poli et al., 2023). All SSMs, linear attention models, and other linear recurrent networks share a common formalism of being linear recurrent networks interleaved with non-linear layers. On the other hand, the ability to do in-context learning seems to be a hallmark of most recent scalable variants (Grazzi et al., 2024). Which features of these successful models contribute to in-context learning, as opposed to earlier variants? Using a constructive approach, we pinpoint input-dependent input and output processing, as the key features required for in-context learning.

We show that SSMs with local self-attention, a form of input-dependent input processing, can perform in-context learning analogously to transformers, i.e. through gradient descent steps on an implicit linear regression problem. The key insight we use is that the state of the recurrent network can be used to aggregate gradients of the parameters of the implicit linear model, which can later be ‘applied’ to the initial parameters of the implicit model. Using the general SSM formalism, we show how to design the recurrent and output equations that enable them to do in-context learning. Our construction, which we call *GD-SSM*, is not restricted to in-context learning tasks and performs well on general-purpose prediction problems.

In summary, our contributions are to show that:

- one-layer SSMs with diagonal recurrence, two-dimensional state, and input-dependent input and output processing can perform one step of minibatch gradient descent on an implicit least-squares loss function;
- multi-step (minibatch) gradient descent can be achieved by stacking the 1-layer model;
- gradient descent for an implicit non-linear regression problem can be achieved by augmenting the SSMs with non-linearities;
- a randomly initialized model trained on regression tasks learns parameters that match our construction for in-context learning tasks based on regression.

## 2 BACKGROUND

A sequence model operates on an input sequence  $\mathcal{S} = \{\mathbf{s}_t\}_{t=1}^T \in \mathbb{R}^{T \times f}$ , where  $T$  is also referred to as the context-length and  $f$  is the feature dimension.

Contemporary sequence models based on transformers interleave self-attention with MLP layers to perform sequence processing. The self-attention performs sequence-mixing while the MLPs perform channel-mixing. The most common form of self-attention has been the scaled dot-product self-attention, which embeds the sequence into a query  $\mathbf{Q} = \mathbf{S}\mathbf{W}_Q$ , key  $\mathbf{K} = \mathbf{S}\mathbf{W}_K$  and value  $\mathbf{V} = \mathbf{S}\mathbf{W}_V$ , where  $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{f \times m}$ ,  $\mathbf{W}_V \in \mathbb{R}^{f \times d}$ , and then calculates the output of attention as

$$\text{SA}(\mathcal{S}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{m}}\right)\mathbf{V}. \quad (1)$$

Local self-attention (Beltagy et al., 2020) uses the same form as Eq. 1, but on an input sequence that is a subset of the full sequence, with a sliding window.

By discarding the softmax (and scaling for simplicity), self-attention can be written in a linear form

$$\text{LSA}(\mathcal{S}) = \mathbf{Q}\mathbf{K}^T\mathbf{V}. \quad (2)$$

Deep SSMs are sequence models that replace the self-attention with a linear recurrent network, to perform the sequence-mixing. In the most general form, the recurrent SSM block consists of a recurrent state  $\mathbf{Z}_t \in \mathbb{R}^{d \times m}$  updated iteratively as

$$\mathbf{Z}_t = \mathbf{A}(\mathbf{s}_t) * \mathbf{Z}_{t-1} + \mathbf{B}(\mathbf{s}_t), \quad (3)$$

where  $\mathbf{A}, \mathbf{B} : \mathbb{R}^f \rightarrow \mathbb{R}^{d \times m}$ , and  $*$  is some multiplication operator (e.g. matrix or element-wise multiplication).

The output of the SSM is often calculated by an input-dependent linear transformation of the current state

$$\mathbf{o}_t = \mathbf{Z}_t\mathbf{U}(\mathbf{s}_t), \quad (4)$$

combined with a non-linearity as

$$\mathbf{o}'_t = \mathbf{U}_2(\mathbf{o}_t \odot \sigma(\mathbf{U}_1\mathbf{o}_t)),$$

where  $\odot$  is elementwise multiplication.

Transformers with linear attention can also be written in this recurrent form Katharopoulos et al. (2020a)

$$\mathbf{Z}_t = \mathbf{Z}_{t-1} + \mathbf{v}(\mathbf{s}_t) \otimes \mathbf{k}(\mathbf{s}_t), \quad (5)$$

where  $\mathbf{v}(\mathbf{s}_t) = \mathbf{W}_V^T \mathbf{s}_t \in \mathbb{R}^d$ ,  $\mathbf{k}(\mathbf{s}_t) = \mathbf{W}_K^T \mathbf{s}_t \in \mathbb{R}^m$  and  $\otimes$  is the outer product. This is a recurrent reformulation of the part of Eq. 2 containing  $\mathbf{V}$  and  $\mathbf{K}$  in recurrent form.

S4 (Gu et al., 2021) and multi-headed S5 (Smith et al., 2022) can also be written in this form

$$\mathbf{Z}_t = \mathbf{A}\mathbf{Z}_{t-1} + \mathbf{B}\mathbf{s}_t, \quad (6)$$

where  $\mathbf{A}$  and  $\mathbf{B}$  are learnable parameters of appropriate dimension, and  $\mathbf{Z}_t$  consists of  $m$  heads of dimension  $d$  each (although multiple heads are not used in the original paper).

If local self-attention is used to process the input before being fed into the SSM, all instances of  $\mathbf{s}_t$  would be replaced by a context vector

$$\mathbf{c}_t = \text{SA}(\mathbf{S}') \quad \text{or} \quad \mathbf{c}_t = \text{LSA}(\mathbf{S}'),$$

where  $\mathbf{S}' \subset \mathbf{S}$  is a subsequence of the whole sequence, and (L)SA denotes the (linear) self-attention operation.

### 3 SSMS CAN EMULATE GRADIENT DESCENT ON LINEAR REGRESSION TASKS

We will now show that an SSM as described in Section 2 can perform gradient descent on an implicit linear model to minimize a least squares loss (for particular choices of parameters). Extensions to non-linear regression models are considered in Section 3.3.

Consider a linear regression problem. The goal is to minimize the corresponding least squares loss using gradient descent. The linear model will be the implicit model on which the SSM performs gradient descent. Performing mini-batch (batch size  $> 1$ ) gradient descent on the parameters of this implicit model involves two steps: (i) to accumulate gradients of the loss with respect to the parameters, and (ii) apply the accumulated gradient to the initial value of the parameters of the linear model to calculate the updated parameters. Predictions can be made with the updated parameters by combining them linearly with the input.

Assume the training samples for the linear regression problem are provided as a sequence of inputs and targets. A large enough SSM can then accumulate the gradients of the loss function in its state *if a local self-attention-like layer processes the sequence inputs before the recurrence*<sup>1</sup>.

Given the accumulated gradients, the next-step emission of the SSM is equivalent to i) updating the parameters of the implicit model gradient and ii) computing the model output with the updated parameters. Multiple steps of gradient descent can be achieved by stacking multiple layers, while nonlinearity in the implicit model can be handled by adding nonlinear input-output embedding layers. We argue that the architecture that allows a single layer to perform gradient descent provides the inductive bias for the model to do in-context learning.

#### 3.1 SINGLE STEP 1-DIMENSIONAL LINEAR REGRESSION

Consider a linear regression model with 1-d output for simplicity

$$y = \mathbf{w}^T \mathbf{x}, \quad (7)$$

for parameter  $\mathbf{w} \in \mathbb{R}^f$ . This is the implicit linear model we aim to reproduce for in the 1-dimensional target case. Given dataset of  $N$  samples  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=0}^N$ ,  $\mathbf{x} \in \mathbb{R}^f$ ,  $y \in \mathbb{R}$ , the associated least squares loss is

$$\mathcal{L}(\mathcal{D}; \mathbf{w}) = \frac{1}{2N} \sum_{i=1}^N \|\hat{y}_i - y_i\|_2^2 = \frac{1}{2N} \sum_i (\mathbf{w}^T \mathbf{x}_i - y_i)^2. \quad (8)$$

The best fit for  $\mathbf{w}$  is the minimum of  $\mathcal{L}$  over  $\mathbf{w} \in \mathbb{R}^f$ . The gradient of the loss calculated on the first  $t$  samples of the dataset is

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathcal{D}_{1:t}; \mathbf{w}_0) = \frac{1}{t} \sum_{i=1}^t (\mathbf{w}_0^T \mathbf{x}_i - y_i) \mathbf{x}_i,$$

<sup>1</sup>A SSM layer with input-dependent recurrence would be able to simulate a local self-attention layer, but with significantly increased computational and conceptual complexity.

where  $\mathcal{D}_{1:t}$  denotes the first  $t$  samples in  $\mathcal{D}$ . The unscaled gradient,

$$\mathbf{g}_{\mathbf{w}_0}(\mathcal{D}_{1:t}) = \sum_{i=1}^t (\mathbf{w}_0^T \mathbf{x}_i - y_i) \mathbf{x}_i,$$

can be recursively calculated as

$$\mathbf{g}_{\mathbf{w}_0}(\mathcal{D}_{1:t}) = \mathbf{g}_{\mathbf{w}_0}(\mathcal{D}_{1:t-1}) + (\mathbf{w}_0^T \mathbf{x}_t - y_t) \mathbf{x}_t. \quad (9)$$

Scaling the  $\mathbf{g}_{\mathbf{w}_0}(\mathcal{D}_{1:t})$  gives the mini-batch gradient  $\nabla_{\mathbf{w}} \mathcal{L}(\mathcal{D}_{1:t}; \mathbf{w}_0) = \frac{1}{t} \mathbf{g}_{\mathbf{w}_0}(\mathcal{D}_{1:t})$  for minibatch size  $t$ .

To make a prediction,  $\hat{y}$ , we apply the gradient to the parameters of the linear model in Eq. 7 and compute the corresponding output, i.e.

$$\begin{aligned} \hat{y}_{t+1} &= (\mathbf{w}_0 - \eta \nabla_{\mathbf{w}_0} \mathcal{L}(\mathcal{D}_{1:t}; \mathbf{w}_0))^T \mathbf{x}_{t+1}, \\ &= (\mathbf{w}_0 - \frac{\eta}{t} \mathbf{g}_{\mathbf{w}_0}(\mathcal{D}_{1:t}))^T \mathbf{x}_{t+1}. \end{aligned}$$

When  $\mathbf{w}_0 = \mathbf{0}$ , this reduces to

$$\hat{y}_{t+1} = -\frac{\eta}{t} \mathbf{g}_{\mathbf{w}_0}(\mathcal{D}_{1:t})^T \mathbf{x}_{t+1}. \quad (10)$$

**Implementation as an SSM:** Equation 9, which is a linear recurrence equation, can be implemented by an appropriately constructed SSM.

**Proposition 1** *Given a diagonal linear recurrent layer, and tokens  $\mathbf{s}_j = \mathbf{c}_j = [\mathbf{x}_j y_j, \mathbf{x}_{j+1}]$ , for  $j = 1, \dots, N$ , and  $[\dots]$  concatenation,  $\mathbf{x}_j, y_j$  drawn from a linear model, one can construct recurrent matrix  $\mathbf{A}(\mathbf{s}_j)$ , input  $\mathbf{B}(\mathbf{s}_j)$  and output matrix  $\mathbf{U}(\mathbf{s}_j)$  such that each recurrent step for every token  $\mathbf{s}_j$  produces  $\hat{y}_{j+1} = -(\Delta \mathbf{w})^T \mathbf{x}_{j+1}$  as output, where  $\Delta \mathbf{w}$  is one step of gradient descent, i.e.  $\Delta \mathbf{w} = \eta \nabla_{\mathbf{w}} \mathcal{L}$ . The test input  $\mathbf{x}_{N+1}$  is contained in token  $\mathbf{c}_N$ , and produces the test prediction  $\hat{y}_{N+1}$ .*

Specifically, if we use  $\mathbf{z}_t \in \mathbb{R}^f$  to denote the state of the recurrent network and let it directly correspond to the vector  $\mathbf{g}_{\mathbf{w}_0}(\mathcal{D}_{1:t})$ , the equivalent SSM layer is a linear recurrence equation,

$$\mathbf{z}_t = \mathbf{I} \mathbf{z}_{t-1} + (\mathbf{w}_0^T \mathbf{x}_t - y_t) \mathbf{x}_t. \quad (11)$$

The state of the SSM,  $\mathbf{z}_t$ , represents the implicit linear regression problem through the *unscaled accumulated gradients* of the least squares loss with respect to the parameters  $\mathbf{w}_0$  of the implicit linear model.

As linear regression is performed on the training dataset  $\mathcal{D} = \{\langle \mathbf{x}_i, y_i \rangle\}_{i=0}^N$ , the SSM receives the training data in the form of a sequence as input. In the most general case, this is a sequence  $\mathbf{s}_1 = \mathbf{x}_1, \mathbf{s}_2 = [0, \dots, y_1], \dots$  where  $[\dots]$  denoting concatenation and  $y_i$  is padded with  $f - 1$  zeros for its dimensions to match that of  $\mathbf{x}_i$ . This more general case is discussed in the next section. But here, we will consider a case which simplifies our construction.

Let  $\mathbf{s}_1, \mathbf{s}_2, \dots$  be a sequence of constructed context vectors  $\mathbf{s}_t = \mathbf{c}_t$ , where each  $\mathbf{c}_t = [\mathbf{x}_t y_t, \mathbf{x}_{t+1}] \in \mathbb{R}^{2f}$ , and let us assume  $\mathbf{w}_0 = \mathbf{0}$  for simplicity<sup>2</sup>. If the sequence input weights  $\Psi \in \mathbb{R}^{2f \times f}$  are such that  $\Psi^T \mathbf{c}_t = \mathbf{x}_t y_t$ , Eq. 11 can be written as an SSM (Eq. 3), i.e.

$$\mathbf{z}_t = \mathbf{I} \mathbf{z}_{t-1} + \Psi \mathbf{c}_t. \quad (12)$$

A parameter matrix,  $\Psi$ , satisfies equation 12 if all but the first  $f$  diagonal entries are zero. The state of this network is the unscaled gradient  $t \nabla_{\mathbf{w}_0} \mathcal{L}(\mathcal{D}_{1:t}; \mathbf{w}_0)$  and the state recursion accumulates the gradients as in step (i) above.

The accumulated gradient is then ‘applied’ to the implicit model’s initial parameters,  $\mathbf{w}_0$ , before computing the  $(N + 1)$ -th output. With  $\mathbf{w}_0 = \mathbf{0}$ , the output is

$$\mathbf{o}_t = \beta \mathbf{z}_t^T \Theta \mathbf{c}_t, \quad (13)$$

<sup>2</sup>The more general case is treated for the case of multi-step GD in Appendix A.2.

where  $\beta = -\frac{\eta}{N}$ ,  $\eta$  is the learning rate, and  $N$  is the number of training points or, equivalently, the total length of the context. The SSM final output,  $o_t$  above, corresponds to a prediction of the trained linear model,  $\hat{y}_{t+1}$  in Eq. 10, if  $\Theta$  obeys  $\Theta \mathbf{c}_t = \mathbf{x}_{t+1}$ . It is easy to check that  $\Theta$  satisfies this condition if all but its last  $f$  diagonal entries are zero (see Figure 3 for a concrete example). Note that the output in Eq. 13 matches the general form of the SSM output in Eq. 4 (without the non-linearity).

The above shows that the following SSM

$$\mathbf{c}_t = [\mathbf{x}_t \mathbf{y}_t, \mathbf{x}_{t+1}], \quad (14)$$

$$\mathbf{z}_t = \mathbf{I} \mathbf{z}_{t-1} + \Psi \mathbf{c}_t, \quad (15)$$

$$o_t = \beta \mathbf{z}_t^T \Theta \mathbf{c}_t = \hat{y}_{t+1}. \quad (16)$$

can perform gradient descent on the parameters  $\mathbf{w}_0$  of the implicit linear model and use this mechanism for in-context learning. The specific structure of the SSM in equation 14 demonstrates the importance of multiplicative processing, for both the inputs and outputs.

### 3.2 SINGLE STEP N-DIMENSIONAL LINEAR REGRESSION

In this section, we generalize the construction above to the N-dimensional case. Without loss of generality, we assume the input and the target,  $\mathbf{x}$  and  $\mathbf{y}$ , have both dimensions  $f$  i.e.  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^f$ . If this is not the case, the input and output dimensionality can be matched by defining appropriate embeddings<sup>3</sup>. We can then treat the N-dimensional system as  $f$  1-D linear regression problems, one for each element of  $\mathbf{y}$ .

**Proposition 2** *Given a diagonal linear recurrent layer augmented with local self-attention with sliding window of size 3, and tokens  $\mathbf{s}_{2j} = \mathbf{x}_j$  and  $\mathbf{s}_{2j+1} = \mathbf{y}_j$ , for  $j = 1, \dots, N$ ,  $\mathbf{x}_j, \mathbf{y}_j$  drawn from a linear model, one can construct recurrent matrix  $\mathbf{A}(\mathbf{s}_j)$ , input  $\mathbf{B}(\mathbf{s}_j)$  and output matrix  $\mathbf{U}(\mathbf{s}_j)$  such that each recurrent step for every token  $\mathbf{s}_j$  produces  $\hat{\mathbf{y}}_{j+1} = -(\Delta \mathbf{W})^T \mathbf{x}_{j+1}$  as output, where  $\Delta \mathbf{W}$  is one step of gradient descent, i.e.  $\Delta \mathbf{W} = \eta \nabla_{\mathbf{W}} \mathcal{L}$ . The test input  $\mathbf{x}_{N+1}$  is contained in token  $\mathbf{s}_{2N+2}$ , and produces the test prediction  $\hat{\mathbf{y}}_{N+1}$ .*

Similar to Eq. 11, we show the above by writing the SSM as

$$\mathbf{Z}_t = \mathbf{Z}_{t-1} + \mathbf{y}_t \mathbf{x}_t^T, \quad (17)$$

where  $\mathbf{Z}_t$  corresponds to the parameters of the implicit linear model,  $\mathbf{W} \in \mathbb{R}^{f \times f}$ , and we assume, for simplicity, that  $\mathbf{W} = \mathbf{0}^4$ . The output is

$$\mathbf{o}_t = \beta \mathbf{Z}_t \mathbf{x}_{t+1}. \quad (18)$$

See Appendix A.1 for the full derivation.

To see how this can be written in the form of Eq. 3, let the input sequence consist of the training dataset of the implicit linear regression problem (as before). This time, we cast the training dataset into a standard sequence  $\mathbf{s}_1, \mathbf{s}_2, \dots$ , where

$$\mathbf{s}_{2j} = \mathbf{x}_j, \quad (19)$$

$$\mathbf{s}_{2j+1} = \mathbf{y}_j, \quad (20)$$

and  $\mathbf{s}_i$  denotes the  $i$ -th sequence element of the input, such that  $\mathbf{s}_{2j+2} = \mathbf{x}_{j+1}$ .

At each step, the state update, in Eq. 17, and the output, in Eq. 18, include  $\mathbf{x}_t, \mathbf{y}_t, \mathbf{x}_{t+1}$ . This necessitates the introduction of a quadratic (in  $\mathbf{s}_i$ ) local self-attention mechanism. Let  $\mathbf{C}_t$  be the context matrix for the local self-attention, i.e. a sliding window of length three running through the sequence,

$$\mathbf{C}_t = \begin{bmatrix} \vdots & \vdots & \vdots \\ \mathbf{x}_t & \mathbf{y}_t & \mathbf{x}_{t+1} \\ \vdots & \vdots & \vdots \end{bmatrix}. \quad (21)$$

<sup>3</sup>E.g. with dimensions  $k, l$  into the same higher dimension  $k + l = f$  by concatenation with appropriately sized zero vectors, or through a linear transformation to dimension  $f$ .

<sup>4</sup>The more general case is discussed in Appendix A.2.

The local self-attention operation,  $CQC^T$ , is a truncated form of Eq. 2. When  $Q = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 0 & 1 & 0 \end{pmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ , plugging  $CQC^T$  into the second term in Eq. 17 produces an SSM with input  $\mathbf{y}_t \mathbf{x}_t^T$ .

The SSM can now be written as in Eq. 3, i.e.

$$\mathbf{Z}_t = \mathbf{Z}_{t-1} + \mathbf{C}_t \mathbf{Q} \mathbf{C}_t^T,$$

with output

$$\mathbf{o}_t = \beta \mathbf{Z}_t \mathbf{C}_t \mathbf{q},$$

where  $\mathbf{q} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$  makes the output  $\hat{\mathbf{y}}_{t+1}$  as in Eq. 10. The output also matches the form of the general SSM output in Eq. 4. The construction allows us to perform gradient descent on the parameters of an arbitrary dimensional implicit model. We call this type of SSM a GD-SSM.

Eq. 3 shows that the recurrent state of the SSM is 2-dimensional, as in most recently proposed SSMs (Dao & Gu, 2024). The local self-attention reproduces the local linear self-attention of Eqs. 2 and 5. As for the 1-dimensional case, the self-attention higher-order dependencies,  $CQC^T$ , is key for making the SSM learn (in-context) an implicit (linear) regression model.

### 3.3 GENERALISING TO ANY REGRESSION PROBLEM

**Multi-step gradient descent:** The proposed construction can be extended to multi-step gradient descent. Since each layer of a GD-SSM produces the parameters of the implicit linear model updated by one step of GD, this is equivalent to stacking together multiple layers. In our derivations above, we assumed the initial parameter of the implicit linear model is  $\mathbf{0}$ . In the multi-step GD, all layers other than the first will correspond to a *non-zero initialised implicit linear model*. Technically, extra gradient steps in the implicit model introduce one additional term in the gradient accumulation equation. Each of the two terms requires a separate (parallel) recurrence and is performed by a dedicated layer. At the end, the states are combined to obtain the multiple-step GD update, with minor extra computational burdens. See Appendix A.2 for a detailed construction of multi-step GD.

**Non-linear regression:** Non-linear regression can be handled by adding MLP layers to the GD-SSM. In the previous sections, we let GD-SSM accumulate the gradients of a linear regressor. Additional MLP layers can learn to transform the state of these linear layers into quantities corresponding to the gradient of the implicit non-linear model. See Appendix A.3 for a detailed explanation.

**Regularisation terms in the loss:** As the recurrent layers only accumulate the gradients, we can separate the gradient calculation and accumulation from its application. This has a practical advantage. Any input-independent regularisation term, e.g. L2 norms, can be added to the model without changing the recurrence structure. See Appendix A.4 for details.

## 4 TRAINED LINEAR RECURRENT NETWORKS DO EMULATE GRADIENT DESCENT ON LINEAR REGRESSION TASKS

We investigated if the GD-SSM variant of the general SSM architecture does do gradient descent in-context learning. To do this, we trained a randomly initialised model on various in-context learning tasks for linear and non-linear regression. In each case, the sequence token input to the model consisted of the inputs  $\mathbf{x}$  and target values  $\mathbf{y}$  from the training dataset  $\mathcal{D}$ , and the model was expected to output the prediction for the query (test)  $\mathbf{x}$  given in the last timestep. The models were trained to minimize the mean squared error between the test prediction and target.

### 4.1 SINGLE STEP 1-DIMENSIONAL LINEAR REGRESSION

We first tested the simplest case of one-step gradient descent on a linear regression problem with scalar predictions/targets. This corresponds directly to the construction in Section 3.1. To do this, similar to Garg et al. (2023); von Oswald et al. (2023), we randomly generated linear regression

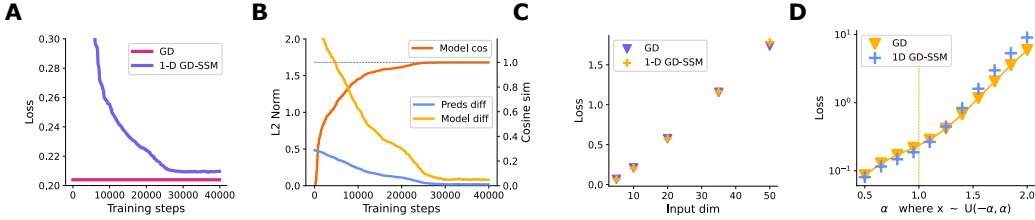


Figure 1: **Comparing one step of GD with a trained single layer GD-SSM for 1-dimensional regression:** **A:** Trained single layer SSM loss and Gradient descent based weight-construction loss are identical. **B:** Cosine similarity and the L2 distance between models as well as their predictions. **C:** Comparison of loss between Gradient Descent (GD) and the SSM layer model for different input sizes  $N = N_x$ . **D:** The trained single 1-D SSM layer, and gradient descent show identically loss (in log-scale) when provided input data different than during training i.e. with scale of 1. We display the mean/std. or the single runs of 5 seeds.

tasks consisting of training and test points, and trained the model to make a prediction for the test input using the training input as context.

We generated randomly sampled linear regression tasks  $\tau$  in the following way: Each task (context)  $\tau$  consisted of a sequence of in-context training data  $\mathcal{D}_\tau = \{\langle \mathbf{x}_{\tau,i}, y_{\tau,i} \rangle\}_{i=1}^N$  and test point  $\langle \mathbf{x}_{\tau,N+1}, y_{\tau,N+1} \rangle$ . To generate this, the  $\mathbf{x}_{\tau,i}$ s are sampled from a uniform distribution  $\mathbf{x}_{\tau,i} \sim U(-1, 1)^f$ . Then, for each task  $\tau$ , the parameters of its implicit linear model  $\mathbf{w}_\tau$  is sampled from a normal distribution, so that each element  $[\mathbf{w}_\tau]_i \sim \mathcal{N}(0, 1)$ . This is used to calculate the  $y_{\tau,i}$ s for each corresponding  $\mathbf{x}_{\tau,i}$  using  $y_{\tau,i} = \mathbf{w}_\tau^T \mathbf{x}_{\tau,i}$ .

The sequence  $\mathcal{S} = \{\mathbf{s}_{\tau,1}, \dots, \mathbf{s}_{\tau,N}\}$  is constructed so that  $\mathbf{s}_{\tau,t} = \mathbf{c}_{\tau,t} = [\mathbf{x}_{\tau,t} y_{\tau,t}, \mathbf{x}_{\tau,t+1}]$ , with  $[\dots]$  denoting the vector concatenation operation, and  $\mathbf{c}_t$  is the constructed context vector. Note that this includes the query  $\mathbf{x}_{\tau,N+1}$  in  $\mathbf{c}_N$ . We will use a more general construction in the next section. The outputs of the GD-SSM at time  $T = N$  is the prediction for  $\mathbf{x}_{N+1}$  i.e.  $SSM(\mathcal{S}_\tau)_N = \hat{y}_{\tau,N+1}$ , with target  $y_{\tau,N+1} = \mathbf{w}_\tau^T \mathbf{x}_{\tau,N+1}$ . The model was trained to minimize the expected squared prediction error, averaged over linear regression tasks  $\tau$ :

$$\min_{\theta} \mathbb{E}_\tau \left[ \|\hat{y}_\theta(\mathbf{c}_{\tau,1}, \dots, \mathbf{c}_{\tau,N}) - y_{\tau,\text{test}}\|^2 \right],$$

where  $\theta$  are the randomly initialised parameters of the GD-SSM. We evaluate our model on multiple metrics:

1. L2 norm between the difference in predictions  $\|\hat{y}_\theta(x_{\tau,\text{test}}) - \hat{y}_{\theta_{\text{GD}}}(x_{\tau,\text{test}})\|_2$  where  $\hat{y}_{\theta_{\text{GD}}}$  is the prediction from the GD based construction.
2. Cosine similarity between the sensitivities  $\frac{\partial \hat{y}_{\theta_{\text{GD}}}(x_{\tau,\text{test}})}{\partial x_{\text{test}}}$  and  $\frac{\partial \hat{y}_\theta(x_{\tau,\text{test}})}{\partial x_{\text{test}}}$ .
3. L2 norm between the sensitivities.
4. Loss of the two models where validation data is sampled from  $U(-\alpha, \alpha)^{n_i}$ , with a different  $\alpha$  than the one used in training.
5. The loss of two models when trained for different number of feature dimensions  $f$ .

Figure 1 shows the results of this comparisons. We find that these metrics show excellent agreement between the trained model and GD, over different hyperparameters. To test if the trained network has learned a general purpose learning rule as opposed to fitting to the dataset, that is to test if the trained network generalises to linear regression tasks outside the training distribution, we draw values of inputs from  $U(-\alpha, \alpha)$ . We see that our model generalises to both cases.

## 4.2 SINGLE STEP N-DIMENSIONAL LINEAR REGRESSION

In the more general case where we allow the targets  $\mathbf{y}_{\tau,i}$  to be of arbitrary dimension, we also relax other assumptions in the form of the input. Notably, instead of constructing the inputs in a particular way, we let the sequence be similar to a more natural sequence of  $\mathbf{x}_1, \mathbf{y}_1, \mathbf{x}_2, \mathbf{y}_2, \dots$ , i.e.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

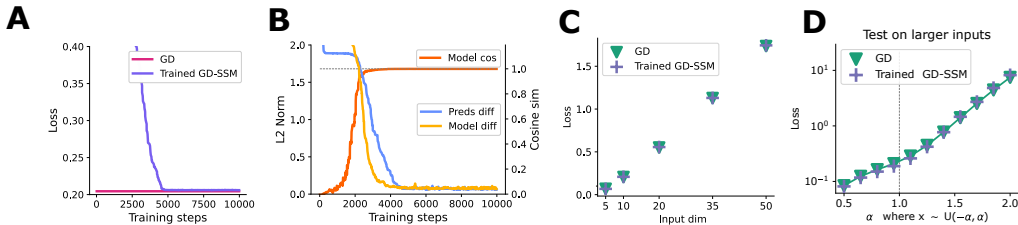


Figure 2: **Comparing one step of GD with a trained single layer GD-SSM for N-dimensional regression:** **A:** Trained single layer GD-SSM loss and Gradient descent based weight-construction loss are identical. **B:** Cosine similarity and the L2 distance between models as well as their predictions. **C:** Comparison of loss between Gradient Descent (GD) and the SSM layer model for different input sizes  $f$ . **D:** The trained GD-SSM layer, and gradient descent show identically loss (in log-scale) when provided input data different than during training i.e. with scale of 1. We display the mean/std. or the single runs of 5 seeds.

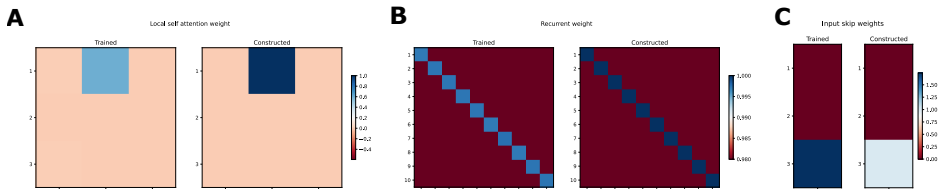


Figure 3: Comparison of learned weights a GD-SSM with the Gradient descent based construction. **A:** Comparison of local self attention weights of trained GD-SSM with the Gradient descent based construction. **B:** Comparison of recurrent parameters of trained GD-SSM with the Gradient descent based construction. Since the recurrence parameters are tensors, for the ease of visualisation, each diagonal entry is the mean of the corresponding diagonal recurrence matrix. **C:** Comparison of skip connection weights of trained GD-SSM with the Gradient descent based construction.

$s_{2j} = x_j; s_{2j+1} = y_j$ . This requires the introduction of a local self-attention mechanism with attention window of 3 as discussed in Section 3.2. The local self-attention mechanism allows the model to access inputs from multiple adjacent timesteps, and not just the current one. We perform the same comparisons as in Section 4.1, and find that the model trained from a random initialisation has an excellent match with the construction, both in the output/loss metrics (Figure 2) and in the parameters it ends up with (Figure 3).

We compared loss metrics for GD-SSM with other standard models including one- and two-layer transformers, and S5 (Smith et al., 2022). We found that, on exactly the same sequence token inputs, only the GD-SSM and a 2-layer transformer reached the same loss as actual GD, while other SSMs such as S5 struggle to perform ICL with 1 layer (Figure 4C).

### 4.3 MULTI-STEP AND NON-LINEAR REGRESSION

To perform multi-step GD, we tested a GD-SSM with multiple-layers, where each layer of the GD-SSM does 1-step of GD as shown in Appendix A.2. The token construction is identical to that used for N-dimensional linear regression, and each layer includes the local self-attention mechanism. In Figure 4A, we show that that trained network is able to reach the same loss as the multi-step GD.



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

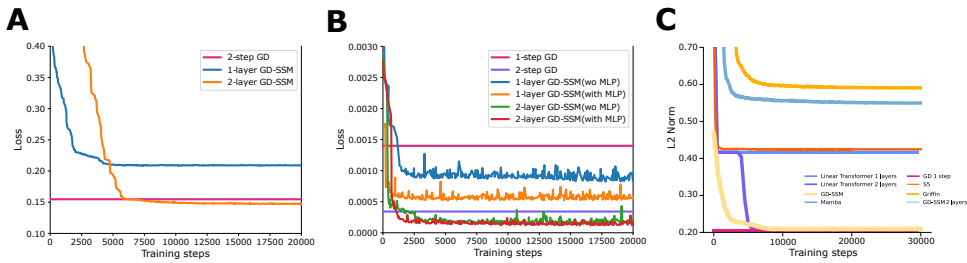


Figure 4: **A:** Comparison of GD-SSM performance with single layer and two layers on regression task. **B:** Comparison of GD-SSM performance with and without MLP layers in single layer and multi-layer setup on non-linear regression task. **C:** Comparison with other models on 1-D regression. The GD-SSM model was evaluated with both 1-layer and 2-layer configurations, and the S5, Mamba, and Griffin models were included for comparison. TF refers to linear Transformer models, with both 1-layer and 2-layer variants tested to evaluate their performance. A similar plot for N-D regression is shown in Figure 5.

Finally, to be able to handle any regression task, we tested GD-SSM with an MLP layer on non-linear regression tasks. Again, we see that the trained model is able to reach the same level of performance as performing GD directly on the dataset (Figure 4B).

## 5 RELATED WORK

In-context learning in transformers has been studied extensively, and various mechanisms have been proposed to explain it (Hendel et al., 2023; von Oswald et al., 2023; Akyürek et al., 2022). Of those, the most prominent is that the self-attention mechanism performs gradient descent on a linear loss. A construction with linear self-attention was demonstrated by von Oswald et al. (2023). While linear-self attention can be written as an RNN (Katharopoulos et al., 2020b), the results of von Oswald et al. (2023) depends on the tokens being constructed in a specific way, which limits the generality of their construction. Moreover, the most common type of self-attention used is softmax scaled dot-product self-attention rather than linear self-attention, and the basic linear self-attention mechanism used in von Oswald et al. (2023) is not competitive with softmax self-attention. Whereas, we show a construction that uses local-self attention with state-space models, which has been shown to be competitive with softmax self-attention transformers (De et al., 2024).

Liu et al. (2024) formulate the SSM layer as an online optimization objective with exact solution. Similarly, Sun et al. (2024) use the state of the SSM to perform GD. But both these papers consider online updates, that is minibatch size 1 updates. One layer of their model is not capable of performing larger minibatch updates. Moreover, without local self-attention, their online optimization objective is limited to using single sequence tokens with a single-layer. And finally their goal is not to explain how in-context learning happens with existing architectures, but rather to develop entirely new SSM variants using the idea that an optimization process can be used for compressing information as well, which is the key property required of a recurrent network.

The combination of local-self attention with Mamba has been used in for improving performance in multiple instances, including most recently (De et al., 2024; Ren et al., 2024). But these papers do not specifically construct their model for ICL, nor do they attempt to distill inductive biases required for ICL in state-space models.

## 6 DISCUSSION

This work establishes a clear connection between state-space models (SSMs) and gradient-based in-context learning. We have demonstrated, both theoretically and empirically, that SSMs can emulate gradient descent on implicit regression models, providing a mechanistic explanation for their in-

486 context learning capabilities. Our construction, GD-SSM, reveals the crucial role of architectural  
487 features such as local self-attention and multiplicative output interactions in enabling this behavior.  
488 These findings not only explain the success of recent SSM variants in in-context learning tasks but  
489 also provide valuable insights for the design of future sequence models.

490 The alignment between our theoretical construction and the behavior of trained models suggests  
491 that the gradient descent mechanism may be a natural inductive bias in these architectures. This  
492 understanding opens new avenues for analyzing and improving sequence models, potentially leading  
493 to more efficient and effective architectures for a wide range of tasks.

494 Future research could explore the implications of this mechanism in larger-scale models, more com-  
495 plex tasks, and real-world applications. Additionally, investigating how this understanding can be  
496 leveraged to enhance the design and training of state-space models could yield significant advance-  
497 ments in the field of sequence modeling.

## 500 REFERENCES

- 501 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning  
502 algorithm is in-context learning? Investigations with linear models, November 2022. URL <http://arxiv.org/abs/2211.15661>. arXiv:2211.15661 [cs].
- 503  
504 Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-Context Language Learning: Ar-  
505 chitectures and Algorithms, January 2024. URL <http://arxiv.org/abs/2401.12973>.  
506 arXiv:2401.12973 [cs] version: 2.
- 507  
508 Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The Long-Document Transformer,  
509 December 2020. URL <http://arxiv.org/abs/2004.05150>. arXiv:2004.05150 [cs].
- 510 Tri Dao and Albert Gu. Transformers are SSMs: Generalized Models and Efficient Algorithms  
511 Through Structured State Space Duality, May 2024. URL [http://arxiv.org/abs/2405.](http://arxiv.org/abs/2405.21060)  
512 [21060](http://arxiv.org/abs/2405.21060). arXiv:2405.21060 [cs].
- 513 Soham De, Samuel L. Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Al-  
514 bert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, Guillaume Des-  
515 jardins, Arnaud Doucet, David Budden, Yee Whye Teh, Razvan Pascanu, Nando De Freitas,  
516 and Caglar Gulcehre. Griffin: Mixing Gated Linear Recurrences with Local Attention for Ef-  
517 ficient Language Models, February 2024. URL <http://arxiv.org/abs/2402.19427>.  
518 arXiv:2402.19427 [cs].
- 519 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of  
520 deep networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International*  
521 *Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp.  
522 1126–1135. PMLR, 06–11 Aug 2017. URL [https://proceedings.mlr.press/v70/](https://proceedings.mlr.press/v70/finn17a.html)  
523 [finn17a.html](https://proceedings.mlr.press/v70/finn17a.html).
- 524 Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What Can Transformers Learn  
525 In-Context? A Case Study of Simple Function Classes, August 2023. URL [http://arxiv.](http://arxiv.org/abs/2208.01066)  
526 [org/abs/2208.01066](http://arxiv.org/abs/2208.01066). arXiv:2208.01066 [cs].
- 527  
528 Riccardo Grazi, Julien Siems, Simon Schrod, Thomas Brox, and Frank Hutter. Is Mamba Capable  
529 of In-Context Learning?, February 2024. URL <http://arxiv.org/abs/2402.03170>.  
530 arXiv:2402.03170 [cs].
- 531 Albert Gu and Tri Dao. Mamba: Linear-Time Sequence Modeling with Selective State Spaces,  
532 December 2023. URL <http://arxiv.org/abs/2312.00752>. arXiv:2312.00752 [cs].
- 533 Albert Gu, Karan Goel, and Christopher Re. Efficiently Modeling Long Sequences with Structured  
534 State Spaces. In *International Conference on Learning Representations*, October 2021. URL  
535 <https://openreview.net/forum?id=uYLFoz1vlAC>.
- 536  
537 Tianyu Guo, Wei Hu, Song Mei, Huan Wang, Caiming Xiong, Silvio Savarese, and Yu Bai.  
538 How Do Transformers Learn In-Context Beyond Simple Functions? A Case Study on Learn-  
539 ing with Representations. October 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=ikwEDvalJZ)  
[ikwEDvalJZ](https://openreview.net/forum?id=ikwEDvalJZ).

- 540 Roei Hendel, Mor Geva, and Amir Globerson. In-Context Learning Creates Task Vectors, October  
541 2023. URL <http://arxiv.org/abs/2310.15916>. arXiv:2310.15916 [cs].
- 542  
543 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are  
544 RNNs: Fast Autoregressive Transformers with Linear Attention, June 2020a. URL <https://arxiv.org/abs/2006.16236v3>.
- 545  
546 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers  
547 are RNNs: Fast Autoregressive Transformers with Linear Attention. In *Proceedings of the*  
548 *37th International Conference on Machine Learning*, pp. 5156–5165. PMLR, November 2020b.  
549 URL <https://proceedings.mlr.press/v119/katharopoulos20a.html>. ISSN:  
550 2640-3498.
- 551  
552 Bo Liu, Rui Wang, Lemeng Wu, Yihao Feng, Peter Stone, and Qiang Liu. Longhorn: State Space  
553 Models are Amortized Online Learners, July 2024. URL [http://arxiv.org/abs/2407.](http://arxiv.org/abs/2407.14207)  
554 [14207](http://arxiv.org/abs/2407.14207). arXiv:2407.14207 [cs].
- 555  
556 Sheng Lu, Irina Bigoulaeva, Rachneet Sachdeva, Harish Tayyar Madabushi, and Iryna Gurevych.  
557 Are Emergent Abilities in Large Language Models just In-Context Learning?, September 2023.  
URL <http://arxiv.org/abs/2309.01809>. arXiv:2309.01809 [cs].
- 558  
559 Antonio Orvieto, Samuel L. Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pas-  
560 canu, and Soham De. Resurrecting Recurrent Neural Networks for Long Sequences, March 2023.  
URL <http://arxiv.org/abs/2303.06349>. arXiv:2303.06349 [cs].
- 561  
562 Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. What In-Context Learning “Learns”  
563 In-Context: Disentangling Task Recognition and Task Learning, May 2023. URL <http://arxiv.org/abs/2305.09731>. arXiv:2305.09731 [cs].
- 564  
565 Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y. Fu, Tri Dao, Stephen Baccus, Yoshua  
566 Bengio, Stefano Ermon, and Christopher Ré. Hyena Hierarchy: Towards Larger Convo-  
567 lutional Language Models, April 2023. URL <http://arxiv.org/abs/2302.10866>.  
568 arXiv:2302.10866 [cs].
- 569  
570 Liliang Ren, Yang Liu, Yadong Lu, Yelong Shen, Chen Liang, and Weizhu Chen. Samba: Simple  
571 Hybrid State Space Models for Efficient Unlimited Context Language Modeling, June 2024. URL  
572 <http://arxiv.org/abs/2406.07522>. arXiv:2406.07522 [cs].
- 573  
574 Jimmy T. H. Smith, Andrew Warrington, and Scott Linderman. Simplified State Space Layers  
575 for Sequence Modeling. September 2022. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Ai8Hw3AXqks)  
[Ai8Hw3AXqks](https://openreview.net/forum?id=Ai8Hw3AXqks).
- 576  
577 Yu Sun, Xinhao Li, Karan Dalal, Jiarui Xu, Arjun Vikram, Genghan Zhang, Yann Dubois, Xinlei  
578 Chen, Xiaolong Wang, Sanmi Koyejo, Tatsunori Hashimoto, and Carlos Guestrin. Learning to  
579 (Learn at Test Time): RNNs with Expressive Hidden States, August 2024. URL <http://arxiv.org/abs/2407.04620>. arXiv:2407.04620 [cs].
- 580  
581 Masayuki Tanaka. Weighted sigmoid gate unit for an activation function of deep neural network.  
582 *Pattern Recognition Letters*, 135:354–359, 2020. ISSN 0167-8655. doi: [https://doi.org/10.1016/](https://doi.org/10.1016/j.patrec.2020.05.017)  
583 [j.patrec.2020.05.017](https://doi.org/10.1016/j.patrec.2020.05.017). URL [https://www.sciencedirect.com/science/article/](https://www.sciencedirect.com/science/article/pii/S0167865518307773)  
584 [pii/S0167865518307773](https://www.sciencedirect.com/science/article/pii/S0167865518307773).
- 585  
586 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
587 Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In I. Guyon, U. V. Luxburg,  
588 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neu-*  
589 *ral Information Processing Systems 30*, pp. 5998–6008. Curran Associates, Inc., 2017. URL  
<http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>.
- 590  
591 Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, Joao Sacramento, Alexander Mord-  
592 vintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers Learn In-Context by Gra-  
593 dient Descent. In *Proceedings of the 40th International Conference on Machine Learning*,  
pp. 35151–35174. PMLR, July 2023. URL [https://proceedings.mlr.press/v202/](https://proceedings.mlr.press/v202/von-oswald23a.html)  
[von-oswald23a.html](https://proceedings.mlr.press/v202/von-oswald23a.html). ISSN: 2640-3498.

594 Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao  
595 Liu, Da Huang, Denny Zhou, and Tengyu Ma. Larger language models do in-context learning  
596 differently, March 2023. URL <http://arxiv.org/abs/2303.03846>. arXiv:2303.03846  
597 [cs].

598  
599 Noam Wies, Yoav Levine, and Amnon Shashua. The Learnability of In-Context Learning, March  
600 2023. URL <http://arxiv.org/abs/2303.07895>. arXiv:2303.07895 [cs].  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

## 648 A APPENDIX

### 649 A.1 N-D LINEAR REGRESSION

650 The implicit linear model of N-dimensional ( $N = f$ ) linear regression is

$$651 \hat{\mathbf{y}} = \mathbf{W}^T \mathbf{x},$$

652 for  $\mathbf{W} \in \mathbb{R}^{f \times f}$  and  $\hat{\mathbf{y}}, \mathbf{x} \in \mathbb{R}^f$ . The loss is

$$653 \mathcal{L}(\mathcal{D}; \mathbf{W}_0) = \frac{1}{2N} \sum_i (\mathbf{W}_0^T \mathbf{x}_i - \mathbf{y}_i)^T (\mathbf{W}_0^T \mathbf{x}_i - \mathbf{y}_i). \quad (22)$$

654 The gradient of this loss calculated using the first  $t$  samples  $\mathcal{D}_{1:t}$  is

$$655 \nabla_{\mathbf{W}} \mathcal{L}(\mathcal{D}_{1:t}; \mathbf{W}_0) = \frac{1}{t} \sum_{i=1}^t \frac{\partial}{\partial \mathbf{W}} (\mathbf{W}^T \mathbf{x}_i - \mathbf{y}_i)^T (\mathbf{W}^T \mathbf{x}_i - \mathbf{y}_i) \Big|_{\mathbf{W}=\mathbf{W}_0}.$$

656 For notational simplicity, we will write this as  $f$  1-D regression problems, one for each element of  $\mathbf{y}$ . Using  $[\mathbf{W}]_{:,i}$  to denote the  $i$ -th column of matrix  $\mathbf{W}$ , using  $[\mathbf{W}_0]_{:,i}$  to denote the  $i$ -th column of matrix  $\mathbf{W}_0$  and using  $[\mathbf{y}_t]_i = [\mathbf{W}]_{:,i}^T \mathbf{x}_t$  to denote the  $i$ -th element of vector  $\mathbf{y}_t$ ,

$$657 \mathbf{g}(\mathcal{D}_{1:t}; [\mathbf{W}_0]_{:,i}) = \mathbf{g}(\mathcal{D}_{1:t-1}; [\mathbf{W}_0]_{:,i}) + ([\mathbf{W}_0]_{:,i}^T \mathbf{x}_t - [\mathbf{y}_t]_i) \mathbf{x}_t, \quad (23)$$

658 where  $\frac{1}{t} \mathbf{g}(\mathcal{D}_{1:t}; [\mathbf{W}_0]_{:,i}) = [\nabla_{\mathbf{W}} \mathcal{L}(\mathcal{D}_{1:t}; \mathbf{W}_0)]_{:,i}$ .

659 **Implementation as a linear recurrent network:** The input is given as in Eqs. 19 & 20, with context matrix constructed as in Eq. 21.

660 The local self attention calculates  $\mathbf{C}\mathbf{Q}\mathbf{C}^T$ , which is then provided as input to the SSM layer. If  $\mathbf{Q} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}^T = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ , this corresponds to the input to the SSM being  $[\mathbf{y}_t]_i \mathbf{x}_t$  (assuming  $\mathbf{W} = \mathbf{0}$ ).

661 Defining  $\mathbf{Z}_t \in \mathbb{R}^{f \times f}$  as the state of the SSM, where  $[\mathbf{Z}_t]_{:,i} = \mathbf{g}(\mathcal{D}_{1:t}; [\mathbf{W}_0]_{:,i})^T$  (note the transpose), and assuming  $\mathbf{W}_0 = \mathbf{0}$ , the equations Eq. 23 for all  $i$ s (all columns) can be written as a single equation

$$662 \mathbf{Z}_t = \mathbf{Z}_{t-1} + \mathbf{C}_t \mathbf{Q} \mathbf{C}_t^T,$$

663 with output

$$664 \mathbf{o}_t = -\frac{\eta}{N} \mathbf{Z}_t \mathbf{C}_t \mathbf{q}.$$

665 If  $\mathbf{q} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ , this corresponds exactly to the output  $\hat{\mathbf{y}}_{t+1} = -\eta \nabla_{\mathbf{W}} \mathcal{L}(\mathcal{D}_{1:t}; \mathbf{W}_0)^T \mathbf{x}_{t+1}$ . This provides a construction that allows us to perform gradient descent on the parameters of an arbitrary dimensional input model.

### 693 A.2 MULTI-STEP GD

694 **Proposition 3** Given  $1, \dots, L$  diagonal linear recurrent layers, each augmented with local self-attention with sliding window of size 3, and sequence tokens  $\mathbf{s}_{2j} = \mathbf{x}_j$  and  $\mathbf{s}_{2j+1} = \mathbf{y}_j$ , for  $j = 1, \dots, N$ , drawn from a linear model, for each layer  $l$ , one can construct pairs of recurrent matrices  $\mathbf{A}^{(l,1)}(\mathbf{s}_j)$ ,  $\mathbf{A}^{(l,2)}(\mathbf{s}_j)$ , inputs  $\mathbf{B}^{(l,1)}(\mathbf{s}_j)$ ,  $\mathbf{B}^{(l,2)}(\mathbf{s}_j)$  and output matrices  $\mathbf{U}^{(l,1)}(\mathbf{s}_j)$ ,  $\mathbf{U}^{(l,2)}(\mathbf{s}_j)$  such that each recurrent step for every token  $\mathbf{s}_j$  produces  $\hat{\mathbf{y}}_{j+1} = (\Delta_l \mathbf{W})^T \mathbf{x}_{j+1}$  as output, where  $\Delta_l \mathbf{W}$  is the update for  $l$  steps of gradient descent, i.e.  $\Delta_l \mathbf{W} = ((\mathbf{W}_0 - \eta \nabla_{\mathbf{W}} \mathcal{L}(\mathcal{D}; \mathbf{W}_0)) - \eta \nabla_{\mathbf{W}} \mathcal{L}(\mathcal{D}; (\mathbf{W}_0 - \eta \nabla_{\mathbf{W}} \mathcal{L}(\mathcal{D}; \mathbf{W}_0)))) \dots l$  times). The test input  $\mathbf{x}_{N+1}$  is contained in token  $\mathbf{s}_{2N+2}$ , and produces the test prediction  $\hat{\mathbf{y}}_{N+1}$ .

We have so far assumed that  $\mathbf{W}_0 = \mathbf{0}$  thus far, which for the first layer corresponds to initialising the parameters of the equivalent linear model to all zeros. For multi-step GD, the second layer onwards have a non-zero initial value of parameters, so let's derive a general form for a layer that performs GD with non-zero initialisation.

Starting from Eq. 23, repeated below for convenience,

$$\mathbf{g}(\mathcal{D}_{1:t}; [\mathbf{W}_0]_{:,i}) = \mathbf{g}(\mathcal{D}_{1:t-1}; [\mathbf{W}_0]_{:,i}) + ([\mathbf{W}_0]_{:,i}^T \mathbf{x}_t - [\mathbf{y}_t]_i) \mathbf{x}_t,$$

we note that this accumulates two different components –  $([\mathbf{W}_0]_{:,i}^T \mathbf{x}_t) \mathbf{x}_t$  and  $[\mathbf{y}_t]_i \mathbf{x}_t$ .

We propose having two different heads (per layer) to accumulate these quantities separately (i)  $[\mathbf{y}_t]_i \mathbf{x}_t$  and (ii)  $\mathbf{x}_t \mathbf{x}_t$ .

Defining  $\mathbf{Z}_t \in \mathbb{R}^{m \times m}$  as the state of the SSM, where  $[\mathbf{Z}_t]_{:,i} = \mathbf{g}(\mathcal{D}_{1:t}; [\mathbf{W}_0]_{:,i})^T$  (note the transpose), the recurrent network corresponding to the accumulation of (i) is the same as before:

$$\mathbf{Z}_t = \mathbf{Z}_{t-1} + \mathbf{C}_t \mathbf{Q} \mathbf{C}_t^T.$$

For (ii), we can write an equivalent recurrent network

$$\tilde{\mathbf{Z}}_t = \tilde{\mathbf{Z}}_{t-1} + \mathbf{C}_t \tilde{\mathbf{Q}} \mathbf{C}_t^T,$$

If  $\tilde{\mathbf{Q}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix}^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ , this corresponds to the input to the SSM being  $\mathbf{x}_t \mathbf{x}_t^T$ .

The output at layer  $l$  is

$$\mathbf{o}_t^{(l)} = (\mathbf{W}_{l-1} + \Delta_l \mathbf{W}) \mathbf{C}_t \mathbf{q},$$

where

$$\Delta_l \mathbf{W} = -\frac{\eta}{N} \left( \mathbf{W}_{l-1}^T \tilde{\mathbf{Z}}_t^{(l)} - \mathbf{Z}_t^{(l)} \right).$$

In summary, at each layer, there are two linear recurrent layers, and the output includes a multiplicative combination across layers and with the external input.

Since each recurrent layer is performing the same operation, one could also loop the output of each layer back to do multiple steps of GD.

### A.3 NON-LINEAR GD

Let us consider a non-linear regression problem with a least squares loss, where the output labels are again 1-dimensional for simplicity. For a given dataset of  $N$  samples  $\mathcal{D} = \{ \langle \mathbf{x}_i, y_i \rangle \}_{i=0}^N$ ,  $\mathbf{x} \in \mathbb{R}^f$ ,  $y \in \mathbb{R}$ , predictions from a non-linear model are generated using

$$\hat{y} = g(\mathbf{w}^T \mathbf{x}),$$

for  $\mathbf{w} \in \mathbb{R}^f$  and where  $g$  is some non-linear function such as sigmoid or an MLP. This will be our implicit non-linear model for this 1-dimensional target case.

The best fit for  $\mathbf{w}$  is found by minimizing the loss

$$\mathcal{L}(\mathcal{D}; \mathbf{w}_0) = \frac{1}{2N} \sum_{i=1}^N \|\hat{y}_i - y_i\|_2^2 = \frac{1}{2N} \sum_i (g(\mathbf{w}_0^T \mathbf{x}_i) - y_i)^2.$$

The gradient of the loss calculated on the first  $t$  samples of the dataset is

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathcal{D}_{1:t}; \mathbf{w}_0) = \frac{1}{t} \sum_{i=1}^t (g(\mathbf{w}_0^T \mathbf{x}_i) - y_i) g'(\mathbf{w}_0^T \mathbf{x}_i) \mathbf{x}_i,$$

where  $\mathcal{D}_{1:t}$  denotes the first  $t$  samples in  $\mathcal{D}$ , and  $g'$  denotes the first derivative of  $g$ .

We define the term  $\mathbf{g}$  but without containing  $g$ , as

$$\mathbf{g}_{\mathbf{w}_0}(\mathcal{D}_{1:t}) = \sum_{i=1}^t (\mathbf{w}_0^T \mathbf{x}_i - y_i) \mathbf{x}_i,$$

which can be recursively calculated as before. Note that this quantity is not the unscaled gradient anymore. In this case we have the following Proposition:

**Proposition 4** *Given a diagonal linear recurrent layer augmented with local self-attention with sliding window of size 3, followed by a MLP layer, and tokens  $\mathbf{s}_{2j} = \mathbf{x}_j$  and  $\mathbf{s}_{2j+1} = \mathbf{y}_j$ , for  $j = 1, \dots, N$ , drawn from a non-linear model, one can construct recurrent matrix  $\mathbf{A}(\mathbf{s}_j)$ , input  $\mathbf{B}(\mathbf{s}_j)$  and output matrix  $\mathbf{U}(\mathbf{s}_j)$  such that each recurrent step for every token  $\mathbf{s}_j$  produces  $\hat{\mathbf{y}}_{j+1} = -(\Delta\mathbf{W})^T \mathbf{x}_{j+1}$  as output, where  $\Delta\mathbf{W}$  is one step of gradient descent, i.e.  $\Delta\mathbf{W} = \eta \nabla_{\mathbf{W}} \mathcal{L}$ . The test input  $\mathbf{x}_{N+1}$  is contained in token  $\mathbf{s}_{2N+2}$ , and produces the test prediction  $\hat{\mathbf{y}}_{N+1}$ .*

Writing this as an SSM exactly as in 12 (with the same  $\Psi$ )

$$\mathbf{z}_t = \mathbf{I} \mathbf{z}_{t-1} + \Psi \mathbf{c}_t,$$

we now need to calculate the output as

$$\hat{\mathbf{y}}_{t+1} = -\eta \nabla_{\mathbf{w}_0} \mathcal{L}(\mathcal{D}_{1:t}; \mathbf{w}_0)^T \mathbf{x}_{t+1}.$$

The interleaved MLP layer  $f$  would need to learn a mapping such that

$$f\left(\sum_{i=1}^t (\mathbf{w}_0^T \mathbf{x}_i - y_i) \mathbf{x}_i\right) = \sum_{i=1}^t (g(\mathbf{w}_0^T \mathbf{x}_i) - y_i) g'(\mathbf{w}_0^T \mathbf{x}_i) \mathbf{x}_i,$$

which it will be able to since it is an universal approximator. This can then be used to calculate the output as

$$\mathbf{o}_t = f(\mathbf{z}_t)^T \Theta \mathbf{c}_t.$$

#### A.4 REGULARISATION TERMS IN THE LOSS

For example, if we wished to change the loss function in Eq. 22 to include L2 regression,

$$\mathcal{L}(\mathcal{D}; \mathbf{W}) = \frac{1}{2N} \sum_{i=1}^N \|\hat{\mathbf{y}}_i - \mathbf{y}_i\|_2^2 + \|\mathbf{W}\|_2^2,$$

this will change the gradient to be

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathcal{D}_{1:t}; \mathbf{W}_0) = \frac{1}{t} \sum_{i=1}^t \frac{\partial}{\partial \mathbf{W}} \|\hat{\mathbf{y}}_i - \mathbf{y}_i\|_2^2 + 2\mathbf{W}_0.$$

To construct an SSM that does GD on this loss instead of the one in Eq. 8, all we'd have to do is change the  $\Delta\mathbf{W}$  in the output to be

$$\Delta_t \mathbf{W} = -\frac{\eta}{N} \left( \mathbf{W}_{t-1}^T \tilde{\mathbf{Z}}_t^{(l)} - \mathbf{Z}_t^{(l)} + \mathbf{W}_{t-1} \right).$$

#### A.5 EXPERIMENTAL DETAILS OF 1-D REGRESSION

The dataset construction and the model evaluation methods are same as that used in von Oswald et al. (2023). Each task (context)  $\tau$  consists of in-context training data  $D_\tau = \{(x_{\tau,i}, y_{\tau,i})\}_{i=1}^N$  and test point  $(x_{\tau,N+1}, y_{\tau,N+1})$ . At every optimization step of the model, we sample the regression parameters  $W_\tau \sim \mathcal{N}(0, I)$ . We then sample  $x_{\tau,i} \sim U(-1, 1)^f$  and construct a scalar target  $y_{\tau,i} = W_\tau x_{\tau,i}$ , where  $f$  is the feature size of the input. To evaluate the model, a set of  $10^4$  tasks are sampled and mean squared error is calculated.

To evaluate the proposed model, we conducted experiments using a constructed token dataset with an input feature size of 10. The model architecture is a single-layer, state space model (SSM) with a hidden dimension of 20. Initialization was performed with 2 blocks and an SSM latent size of 10. No activation function was used during training.

The training process spanned a maximum of 300,000 epochs with a batch size of 64. A cosine annealing schedule and linear warmup were utilized for the learning rate, beginning with an initial SSM learning rate of  $1 \times 10^{-4}$  for optimizing the SSM parameters using the AdamW optimizer. The global learning rate for the remaining parameters was set to  $2 \times 10^{-4}$ , and these parameters were also optimized using AdamW. A weight decay of 0.05 was applied to regularize the model and prevent overfitting.

## 810 A.6 EXPERIMENTAL DETAILS OF N-D REGRESSION

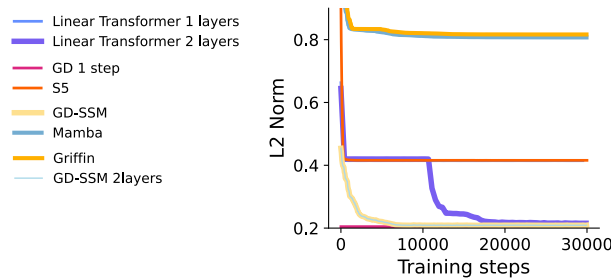
811  
812 The setup is similar to 1-d regression, but we sample  $n_o$  different regression parameters  $W_\tau^k \sim$   
813  $\mathcal{N}(0, I)$ , where  $n_o$  is the dimension of vector  $y_{\tau,i}$  and  $1 \leq k \leq n_o$ . We then construct target  
814  $y_{\tau,i}^k = W_\tau^k x_{\tau,i}$  for all  $k$ . For all our experiments, we choose  $n_o = f$ . We have evaluated the model  
815 based on all the experiments that are used for the 1-D regression evaluation.

816 To have the model prediction equivalent to one step gradient descent, we train the single layer GD-  
817 SSM using Adam optimizer, with an initial learning rate of 0.0001 for recurrent parameters with  
818 cosine annealing. For all the other parameters we double the learning rate that is used for recurrent  
819 parameters. In all our experiments, each optimization step contains 64 tasks.  
820

## 821 A.7 EXPERIMENTAL DETAILS OF MULTI-STEP AND NON-LINEAR REGRESSION

822 To emulate the multi-step regression task, we train the multi-layer GD-SSM architecture, without  
823 applying any non-linearity between the layers. We use the same hyper-parameters that are used for  
824 the N-d regression tasks. We compare the validation loss of the regression task between the learned  
825 model and the model based on construction.  
826

827 For non-linear regression tasks, we use the experiments from Finn et al. (2017) and follow the data  
828 construction of von Oswald et al. (2023). To the output of the linear GD-SSM layer, we add a  
829 non-linear function, weighted sigmoid gated unit (Tanaka, 2020), as used in (Smith et al., 2022).  
830 For all the non-linear regression tasks, an input embedding layer is used in addition to the model  
831 architecture used for linear regression tasks. The hyper-parameters are the same that is used for all  
832 the previous tasks, and the trained model loss is compared with a model with GD-SSM layer/layers  
833 based on gradient descent construction and trained non-linearity layer/layers.  
834



835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846 Figure 5: Comparison with other models on N-D regression. The GD-SSM model was evaluated  
847 with both 1-layer and 2-layer configurations, and the S5, Mamba, and Griffin models were included  
848 for comparison. TF refers to linear Transformer models, with both 1-layer and 2-layer variants tested  
849 to evaluate their performance.  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863