# Understanding Input Selectivity in Mamba:
# Impact on Approximation Power, Memorization, and Associative Recall Capacity

**Ningyuan Huang** [1 2] **Miguel Sarabia** [3] **Abhinav Moudgil** [1 4] **Pau Rodríguez** [3] **Luca Zappella** [3] **Federico Danieli** [3]

## Abstract

State-Space Models (SSMs), and particularly Mamba, have recently emerged as a promising alternative to Transformers. Mamba introduces input selectivity to its SSM layer (S6) and incorporates convolution and gating into its block definition. While these modifications do improve Mamba's performance over its SSM predecessors, it remains largely unclear how Mamba leverages the additional functionalities provided by input selectivity, and how these interact with the other operations in the Mamba architecture. In this work, we demystify the role of input selectivity in Mamba, investigating its impact on function approximation power, long-term memorization, and associative recall capabilities. In particular: (i) we prove that the S6 layer of Mamba can represent projections onto *Haar wavelets*, providing an edge over its Diagonal SSM (S4D) predecessor in approximating discontinuous functions commonly arising in practice; (ii) we show how the S6 layer can dynamically counteract memory decay; (iii) we provide analytical solutions to the MQAR associative recall task using the Mamba architecture with different mixers — Mamba, Mamba-2, and S4D. We demonstrate the tightness of our theoretical constructions with empirical results on concrete tasks. Our findings offer a mechanistic understanding of Mamba and reveal opportunities for improvement.

## 1. Introduction

State Space Models (SSMs) have recently emerged as a promising approach for long-range sequence modeling, due to their computational efficiency compared to Transformers, and parallelizability compared to (nonlinear) Recurrent Neural Networks (RNNs). In particular, Mamba (Gu & Dao, 2023; Dao & Gu, 2024) demonstrated state-of-the-art performance on various language modeling tasks, with smaller model size and faster inference than Transformers. The success of Mamba has largely been attributed to the fact that the parameters in its SSM layer (S6) are input-dependent, leading to improved expressivity compared to its SSM predecessors (Cirone et al., 2024). The goal of this work is to provide a more structural explanation for Mamba's superior performance. We do so by answering two questions: (i) how does the S6 layer's expressivity translate into its practical performance? (ii) how can the S6 layer interact with the rest of the Mamba block to solve concrete tasks?

We answer question (i) by providing a fine-grained analysis of the S6 layer via the lens of function approximation and long-term memory. We prove that the S6 layer can represent projections onto *Haar wavelets* and thus efficiently model discontinuous signals, which is relevant for solving practical tasks. Moreover, we show that the S6 layer still suffers from exponential memory decay, but highlight a mechanism which allows it to dynamically counteract such decay.

Building on the understanding of the S6 layer, we then investigate how S6 interacts with the other components of the Mamba block (which also encompasses a short-convolution and a gate branch, see Tab. 2) to tackle the Multiple-Query Associative Recall (MQAR) task (Arora et al., 2024a). We prove that 1-layer Mamba and Mamba-2 models can both solve MQAR, even without gating; we describe how S6 and the convolution interact to achieve this, and how Mamba-2 leverages its independent convolutions to get more parameter-efficient solutions. We also show that a 1-layer Mamba can solve MQAR exactly *even without input-dependence* in its SSM: this (perhaps unexpected) result helps cementing the importance of convolution and gating in the Mamba architecture. This analysis further informs us on how a variation to the functional form of Mamba, and particularly to how the SSM state matrix is affected by the input, can improve its performance on the INDUCTION HEADS task (Bietti et al., 2024; Sanford et al., 2024a).

We complement our theoretical findings with numerical experiments on synthetic sequence modeling tasks. For the S6

layer, we demonstrate its approximation power on discontinuous functions, and its counteraction of memory decay, via the KEEP $n$-TH task — a generalization of KEEP FIRST from (Chiang & Cholak, 2022) requiring to memorize the $n$-th token in a sequence. Finally, for the full Mamba model, we confirm empirically that the model sizes prescribed theoretically by our analytical solutions to the MQAR and INDUCTION HEADS tasks are tight in practice. Overall, our contributions can be summarized as follows:

- We prove that the S6 layer of Mamba can represent projections onto *Haar wavelets*, providing an edge over the S4D layer in approximating discontinuous functions commonly arising in practice (Sec. 4.1).

- We use sensitivity analysis to show the S6 layer generally suffers from exponential decay of memory, and describe how it can dynamically counteract it (Sec. 4.2).

- We show how the Mamba architecture can exactly solve the MQAR task using different SSM mixers, with an explicit characterization of the required model size, which helps explaining their performance difference (Sec. 5.1).

- Our findings reveal opportunities to further improve Mamba, such as by changing the way input dependence is incorporated into the SSM state matrix (Sec. 5.2).

## 2. Related Work

**State-Space Models (SSMs)** SSMs (i.e., linear RNNs) have recently emerged as a promising sequence-modeling approach, with faster training than (nonlinear) RNNs, and faster inference time than Transformers. To enable long-term memorization capability, Gu et al. (2020) designed SSMs as polynomial approximations of signals, by prescribing non-normal HiPPO matrices as state matrices. To improve computational efficiency, Gu et al. (2022b) proposed the S4 model: first by reparameterizing HiPPO as a sum of normal and low-rank matrices, then by further considering a *diagonal* simplification in the S4D model (Gu et al., 2022a). All these SSMs are Linear Time-Invariant and thus computationally efficient, but consequently lack the ability to process information in a time-varying, input-dependent manner. Mamba (Gu & Dao, 2023) overcomes this by using input-dependent SSM parameters, without sacrificing computational efficiency, showing performance competitive with Transformers on long-range language modeling tasks. Mamba-2 (Dao & Gu, 2024) simplifies the state matrix to a scalar multiple of identity and modifies the Mamba model architecture, showing further empirical improvements.

**Expressivity of SSMs** The analyses on the expressivity of SSMs can be divided into two main approaches: formal language theory and approximation theory. Studies following the first approach investigate what formal languages can SSM recognize (Merrill et al., 2024; Sarrof et al., 2024; Grazzi et al., 2025). Studies ascribing to the second approach — including this work — characterize what function classes can SSMs approximate: Li et al. (2022) showed that SSMs can approximate linear functionals with exponential memory decay; Wang et al. (2024) extended this analysis to nonlinear RNNs, showing however that adding nonlinearity (in the hidden state recurrence) does not fix the memory decay issues. Orvieto et al. (2024) proved that SSMs augmented with MLPs are universal approximators of regular functionals, but this improvement over Li et al. (2022) and Wang et al. (2024) requires the hidden state size to grow linearly with sequence length. Cirone et al. (2024) extended the results from Li et al. (2022) to *input-dependent* SSMs such as Mamba, showing their universal approximation on the class of nonlinear functionals arising from controlled differential equations. While Cirone et al. (2024) highlighted how Mamba can approximate a *larger* class of functionals than S4D, we explore the consequences of this observation by characterizing specific function classes arising in practice.

**Associative Recall Capability** Associative Recall (AR) describes the ability of a model to retrieve information from its memory, based on the input context. Tasks for evaluating associative recall capabilities include INDUCTION HEADS (Olsson et al., 2022; Sanford et al., 2024a), $k$-HOP INDUCTION HEADS (Sanford et al., 2024b), MULTIPLE-QUERY ASSOCIATIVE RECALL (MQAR) (Arora et al., 2024a), and NEEDLE-IN-THE-HAYSTACK (Kamradt, 2023). Empirically, Mamba (Gu & Dao, 2023) and Mamba-2 (Dao & Gu, 2024) demonstrated performance competitive with Transformers on (a simple version of) INDUCTION HEADS and MQAR, respectively. Theoretical understanding of how language models perform associative recall begins to emerge: Bietti et al. (2024) constructed a 2-layer Transformer with positional encoding that solves the INDUCTION HEADS task; Sanford et al. (2024b) extended such construction to a log-depth Transformer that solves the $k$-HOP INDUCTION HEADS task; Arora et al. (2024a) showed that a gated convolution model can solve MQAR. However, these constructions all require the model size to scale with the input sequence length. In this work, we show that 1-layer Mamba models can solve INDUCTION HEADS and MQAR with model size *independent of* the sequence length, highlighting the role of the convolution operation and the gate branch, that has been previously overlooked in the literature.

## 3. Preliminaries: Linear RNNs as SSMs

The application of a Linear RNN can be interpreted as the discrete solution of a *linear* dynamical system, or *SSMs*,

$$\dot{\boldsymbol{h}}(t) = \boldsymbol{A}(t)\boldsymbol{h}(t) + \boldsymbol{B}(t)\,x(t), \qquad \boldsymbol{y}(t) = \boldsymbol{C}(t)\boldsymbol{h}(t). \quad (1)$$

Here, the input $x(t) \in \mathbb{R}$ acts as forcing term for the *hidden state* $\boldsymbol{h} \in \mathbb{R}^N$ through the application of the *input matrix* $\boldsymbol{B}(t) \in \mathbb{R}^{N \times 1}$. The natural evolution of the hidden state is dictated by the *state matrix* $\boldsymbol{A}(t) \in \mathbb{R}^{N \times N}$. Finally, the output $\boldsymbol{y} \in \mathbb{R}^{d_y}$ is obtained by linearly transforming the hidden state via the *output matrix* $\boldsymbol{C}(t) \in \mathbb{R}^{d_y \times N}$. Collectively, we call $\boldsymbol{A}(t), \boldsymbol{B}(t), \boldsymbol{C}(t)$ the *SSM parameters*. For computational efficiency, modern Linear RNNs consider *diagonal* state matrices, with negative eigenvalues $\boldsymbol{A}(t) = \boldsymbol{\Lambda}(t) = -\mathrm{diag}([\lambda_1(t), \dots, \lambda_N(t)])$. This simplifies the solution of (1), and ensures stability in the evolution of the hidden state. Under this assumption, and considering an initial state $\boldsymbol{h}(0) \equiv \boldsymbol{0}$, we can explicitly write the hidden state solution as an integral function of the input (Dahleh et al., 2011):

$$\boldsymbol{h}(t) = \int_0^t e^{\int_s^t \boldsymbol{\Lambda}(r)\, dr} \boldsymbol{B}(s)\, x(s)\, ds. \tag{2}$$

Generally, inputs to Linear RNNs are provided as (discrete) sequences of values $[\boldsymbol{x}_t]_{t=1}^T$, and thus we consider a *discretization* of system (1). This amounts to substituting the *differential* equation with an (approximating) *recurrent* one:

$$(1) \approx \boldsymbol{h}_t = \overline{\boldsymbol{\Lambda}}_t \boldsymbol{h}_{t-1} + \overline{\boldsymbol{B}}_t x_t, \qquad \boldsymbol{y}_t = \overline{\boldsymbol{C}}_t \boldsymbol{h}(t), \quad (3)$$

where $\overline{\boldsymbol{\Lambda}}_t, \overline{\boldsymbol{B}}_t, \overline{\boldsymbol{C}}_t$ are the discrete counterparts of $\boldsymbol{\Lambda}(t)$, $\boldsymbol{B}(t), \boldsymbol{C}(t)$, respectively. Notice we can recover an explicit solution to (3) by unrolling the recurrence relation starting from $\boldsymbol{h}_0 \equiv \boldsymbol{0}$, to obtain

$$\boldsymbol{h}_t = \sum_{s=1}^t \left( \prod_{r=s+1}^t \overline{\boldsymbol{\Lambda}}_r \right) \overline{\boldsymbol{B}}_s x_s. \tag{4}$$

The discretized SSM parameters are usually obtained following a Zeroth-Order Hold (ZOH) scheme (Tóth et al., 2008). Given a discrete time-step $\Delta(t) \in \mathbb{R}^+$, ZOH prescribes

$$
\begin{aligned}
\overline{\boldsymbol{\Lambda}}_t &= e^{\boldsymbol{\Lambda}(t)\Delta(t)}, & \overline{\boldsymbol{C}}_t &= \boldsymbol{C}(t), \\
\overline{\boldsymbol{B}}_t &= (\boldsymbol{\Lambda}(t)\Delta(t))^{-1}(e^{\boldsymbol{\Lambda}(t)\Delta(t)} - I)(\boldsymbol{B}(t)\Delta(t)).
\end{aligned}
\tag{5}
$$

Note that in the original Mamba formulation the authors use Forward Euler for $\boldsymbol{B}(t)$ for simplicity, $\overline{\boldsymbol{B}}_t = \boldsymbol{B}(t)\Delta(t)$.

In general, there is some flexibility in the choice of the functional form that the system parameters can take. The main requirements are that: (i) $\overline{\boldsymbol{\Lambda}}_t$ is diagonal(-izable), so that the product $\prod_{r=s+1}^t \overline{\boldsymbol{\Lambda}}_r$ in (4) can be computed efficiently; (ii) the eigenvalues of $\overline{\boldsymbol{\Lambda}}_t$ are bounded in $[-1, 1]$, to ensure stability; and that (iii) the SSM parameters do *not* depend on the state $\boldsymbol{h}(t)$, so to keep the system linear in $\boldsymbol{h}(t)$, and allow to compute its solution in parallel along $t$. The most relevant choices analyzed in this paper are: S4D (Gu et al., 2022a), which models linear time-invariant dynamics,

$$\Delta(t) := 1, \ \ \boldsymbol{\Lambda}(t) := \boldsymbol{\Lambda}, \ \ \boldsymbol{B}(t) := \boldsymbol{B}, \ \ \boldsymbol{C}(t) := \boldsymbol{C}; \ \ (6)$$

and Mamba (Gu & Dao, 2023), which models time-varying (input-dependent) dynamics,

$$
\begin{aligned}
\Delta(t) &:= \mathrm{SoftPlus}(\mathrm{Linear}(x(t))), \quad \boldsymbol{\Lambda}(t) := \boldsymbol{\Lambda}, \\
\boldsymbol{B}(t) &\equiv \boldsymbol{B}(x(t)) := \mathrm{Linear}(x(t)), \\
\boldsymbol{C}(t) &\equiv \boldsymbol{C}(x(t)) := \mathrm{Linear}(x(t)).
\end{aligned}
\tag{7}
$$

Note that both the Mamba-2 mixer and *Linear Attention* (Katharopoulos et al., 2020) can be interpreted as specializations of Mamba (Dao & Gu, 2024), where $\overline{\boldsymbol{\Lambda}}_t \equiv \lambda_t \boldsymbol{I}$ and $\overline{\boldsymbol{\Lambda}}_t \equiv \boldsymbol{I}$, respectively.

## 4. Mamba SSM Mixer Layer Analysis

In this section, we focus on the core of the Mamba architecture: its SSM mixer layer, also referred to as S6. The goal is to understand how input-selectivity impacts its expressivity and long-range memorization capacity. We show that an S6 layer can represent projections onto wavelets, thus efficiently modeling discontinuous signals. This is useful in practice, e.g., for isolating specific tokens in a sequence. For long-range memorization, we use sensitivity analysis to show that, while S6 suffers from exponential decay of memory (akin to S4D), input-selectivity allows to decrease the rate of such memory decay by "freezing time". We validate our theoretical insights by experimenting on the KEEP $n$-TH task.

### 4.1. Function Approximation Power: Expressing Wavelets via the S6 Layer

We begin by analyzing the expressivity of the S6 layer and the power of its input-dependent discretization in terms of function approximation capabilities. We prove that an S6 layer can approximate projections onto wavelets arbitrarily well (Thm. 1), while an S4D layer can at best project onto Fourier bases. Consequently, in approximating target functions with discontinuities, S6 achieves a faster approximation rate than S4D (Cor. 1). The advantage of S6 over S4D in approximating discontinuous functions translates into their performance differences in memorization tasks.

**Linear RNNs as Time-Projection Onto Basis Functions**
The idea that Linear RNNs could perform projections onto specific sets of basis functions is not new, and indeed served as theoretical grounding for the HiPPO work (Gu et al., 2020). However, to our knowledge, this interpretation has not yet been leveraged to explain the capabilities of modern Linear RNNs. In what follows, we take such interpretation to analyze their expressivity. To streamline the analysis and slim notation, we focus on 1D inputs and view them as continuous signals, $\boldsymbol{x}_s \equiv x(s) \in \mathbb{R}$, $s \in [0, T]$, without loss of generality. For direct comparison with S4D, we consider a simplified version of S6 where the input-dependence only affects the state matrix via $\Delta(x)$, and does *not* af-

3

fect $\boldsymbol{B}(x_t)$ (i.e., $\boldsymbol{B}(x_t) = [B_1, \ldots, B_N]^\top$ independent of the input $x_t$). Substituting this into (2), and recalling that $\boldsymbol{\Lambda} = -\operatorname{diag}([\lambda_1, \ldots, \lambda_N])$, each component $n = 1, \ldots, N$ of the hidden state can be evaluated separately as an inner product between time-dependent functions:

$$h_n^{\text{M}}(t) = \int_0^t \underbrace{e^{-\lambda_n \int_s^t \Delta(x_r)\,dr} B_n}_{=:g_n^{\text{M}}(s;t,x)} x(s)\,ds = \langle g_n^{\text{M}}, x\rangle. \quad (8)$$

We refer to $g_n^{\text{M}}(s; t, x)$ as the *Mamba basis function*, with the notation emphasizing its general dependency on the input signal $x$ up to time $t$. For ease of comparison, we can recover an analogous formula to (8) also for S4D, by letting $\Delta(x_t) \equiv 1 \,\forall t$ (see also (6)). This gives

$$h_n^{\text{S4D}}(t) = \int_0^t \underbrace{e^{-\lambda_n(t-s)} B_n}_{=:g_n^{\text{S4D}}(s;t)} x(s)\,ds = \langle g_n^{\text{S4D}}, x\rangle. \quad (9)$$

As we can see, S4D can provide only exponentials as basis functions. The approximation properties of these functions are limited: this is established in the literature, and ties back to the theory of Vandermonde matrices (Gautschi & Inglese, 1987), as recently pointed out by Orvieto et al. (2024). Their poor performance are mainly due to: (i) the stability constraint, $\operatorname{Re}(-\lambda_n) \leq 0$, which causes an exponentially fast decay to 0, *de-facto* limiting the effective support of said bases; (ii) the large degree of overlap between different bases (obtained by varying the only free parameter $\lambda_n$). The only way to curb these negative effects is by pushing the eigenvalues to be equispaced onto the complex unit disk, namely $e^{-\lambda_n} \to e^{i\frac{2\pi n}{N}}$. This is precisely the strategy recommended by Orvieto et al. (2024), however it reduces the application of the S4D layer to simply performing a Fourier transform. In contrast, the additional flexibility provided by the input-selectivity in Mamba allows for a much richer variety of basis functions (8) to be employed in the projection, an example of which is shown next.

**Mamba Bases Can Represent Haar Wavelets**  Here we provide the main theoretical result in this section, namely that the Mamba S6 layer can perform projections onto Haar wavelets. Due to their ability to capture local aspects of a function such as spikes and discontinuities, wavelets are generally better suited than Fourier bases in solving certain signal processing tasks, (e.g., needles-in-the-haystack (Kamradt, 2023), transient signals (Mallat, 2012)). Recall the Haar wavelets are defined by dilation and translation,

$$\begin{aligned}
\psi_{0,0}(s) = \psi(s) &:= \mathbb{1}_{[0,\frac{1}{2})}(s) - \mathbb{1}_{[\frac{1}{2},1]}(s) \\
\psi_{j,k}(s) &:= 2^{j/2}\psi(2^j s - k),
\end{aligned} \quad (10)$$

with $j \in \mathbb{N}$ denoting the dilation scale, and $k = 0, \ldots, 2^j - 1$ the translation. Higher-order wavelets correspond to localized and "spiky" bases; see Fig. A.1 (left) for an illustration.

**Theorem 1.** *Consider a Haar wavelet $\psi_{j,k} : [0,1] \to \mathbb{R}$, and the Mamba basis function (8) at $t = 1$, $g_{j,k}^{\text{M}}(s; 1, x) = e^{-\lambda_{j,k} \int_s^1 \Delta_{j,k}(x_r)\,dr} B_{j,k}$. Let $\tilde{x}_s := \operatorname{concat}[x_s; s]$ be the input signal augmented with time positional encoding. For any $\epsilon > 0$, there exist 3 Mamba basis functions $g_{j,k}^{\text{M}_1}, g_{j,k}^{\text{M}_2}, g_{j,k}^{\text{M}_3}$ such that the approximation error*

$$\left| \psi_{j,k}(s) - \left( g_{j,k}^{\text{M}_1}(s; 1, \tilde{x}) + g_{j,k}^{\text{M}_3}(s; 1, \tilde{x}) - 2g_{j,k}^{\text{M}_2}(s; 1, \tilde{x}) \right) \right|$$

*is smaller than $\epsilon$, $\forall s \in [0, 1]$.*

The proof relies on tweaking the input-dependent discretization $\Delta(s)$ to output $\Delta(s) \to \infty$ or $\Delta(s) \to 0$ (note that $\Delta$ can directly depend on the time variable $s$ instead of the input signal $x_s$, due to time Positional Encoding (PE)). This effectively allows Mamba to represent Heaviside functions as bases: by linearly combining shifted Heaviside bases, one can immediately recover the required Haar wavelets (see Fig. A.1, middle-right subplots). The details of the proof are reported in App. A.2, where we also proceed to relax the inclusion of PE as an assumption for Thm. 1.

Theorem 1 translates into practical advantages of Mamba over S4D, as Haar wavelets are much better than Fourier bases for approximating *discontinuous* functions common in practice. This is formalized in the following corollary.

**Corollary 1.** *For a piecewise-constant function $\rho(t)$ with $m \geq 1$ discontinuities, there exist $N$ Mamba basis functions (8) such that the $L^2$ approximation error $\|\rho - \sum_{n=1}^N g_n^{\text{M}}\|_{L^2}$ is of order $\mathcal{O}(2^{-\frac{N}{3m}})$. On the other hand, S4D basis functions can achieve an approximation error of $\mathcal{O}(N^{-1})$.*

Corollary 1 stems from Thm. 1, and from approximation results using Haar wavelets and Fourier bases available in the literature (Vetterli, 2001; Eckhoff, 1993); see proof in App. A.2. In the following, we illustrate how approximating wavelets translates into advantages on concrete tasks.

**Task KEEP $n$-TH**  The goal of the KEEP $n$-TH task is to recover the $n$-th element in a randomly-generated sequence, $y_t = x_n$. This generalizes the KEEP FIRST task in (Chiang & Cholak, 2022) where $n = 1$. The solution of KEEP $n$-TH can be directly represented by combining the projections of a piecewise-constant signal $x(s) = \sum_{i=1}^t x_i \mathbb{1}_{[i-1,i)}(s)$ onto two Heaviside functions $H(s - n), H(s - (n - 1))$. As we have shown in Thm. 1, one S6 layer in Mamba can reproduce precisely this type of projections, provided the input is augmented with time-positional information. Thus, we arrive at the following Corollary:

**Corollary 2.** *There exists an S6 layer that solves KEEP $n$-TH on input augmented with time Positional Encoding.*

The proof of Cor. 2 is reported in App. A.3; the results in Tab. 1 verify empirically that Mamba with PE can perfectly

Table 1. KEEP FIFTH experimental results. Average accuracy across 3 seeds with $T = 50$ and $|V| = 128$. Standard error across 3 seeds is 0.00 for all models. Mamba and S4D models consist of *embedding*, *SSM* ($\boldsymbol{h} \in \mathbb{R}^{8 \times 32}$), and *linear* layers (without convolution and gating, see Tab. 2). Positional Encoding (PE) encodes the position in the last element of the embedding, resulting in $|V|$ fewer parameters.

| | MAMBA+PE | MAMBA | S4D | S4D+PE | TRANSFORMER |
|---|---|---|---|---|---|
| Accuracy↑ | **1.00** | 0.08 | 0.09 | 0.08 | **1.0** |
| Parameters | 9.2k | 9.3k | 8.8k | 8.7k | 6.4k |

solve the task (same as Transformers), whereas Mamba without PE and S4D both fail, highlighting the advantage of Mamba over S4D in approximating discontinuous functions in practice. Additional ablations on model size, sequence length, and the role of PE are reported in App. D.2.

## 4.2. Long-Range Modelling: Sensitivity Analysis

In this section, we examine the long-range memorization capacity of SSM layers, by performing a sensitivity analysis of the layer output with respect to changes to the input, as the sequence length increases. To this end, we analyze the derivative of the SSM hidden state at time $t$ with respect to the past input at time $j$, $|\frac{\partial h_t}{\partial x_j}|$. We argue that preserving sensitivity (i.e., a non-zero derivative) is *necessary* for memorization: if the past input has no impact on the current state, one cannot hope for any information about it to be retained. With Lem. 1, we show how generally the sensitivity of an S6 layer decays exponentially fast, similarly to S4D. Thanks to input selectivity, however, the S6 layer can adjust the rate of this decay, thus dynamically tweaking the amount of information to retain. In Lem. 2 we illustrate this mechanism, which proves to be useful for solving the task in Sec. 5.2.

For simplicity, we consider 1D inputs $x_t \in \mathbb{R}$. Given a generic, input-dependent recurrence relationship as in (4), we show in App. B that the sensitivity of the state with respect to its inputs at the $j$-th instant $x_j \in \mathbb{R}$ is given by

$$
\begin{aligned}
\frac{\partial \boldsymbol{h}_t}{\partial x_j} &= \frac{\partial}{\partial x_j} \left( \sum_{s=1}^t \left( \prod_{r=s+1}^t \overline{\boldsymbol{\Lambda}}_r \right) \overline{\boldsymbol{B}}_s x_s \right) \\
&= \left( \prod_{r=j+1}^t \overline{\boldsymbol{\Lambda}}_r \right) \left( \frac{\partial}{\partial x_j} (\overline{\boldsymbol{B}}_j x_j) \right. \\
&\quad \left. + \frac{\partial \overline{\boldsymbol{\Lambda}}_j}{\partial x_j} \sum_{s=1}^{j-1} \left( \prod_{r=s+1}^{j-1} \overline{\boldsymbol{\Lambda}}_r \right) \overline{\boldsymbol{B}}_s x_s \right).
\end{aligned}
\tag{11}
$$

**Lemma 1.** *Consider the hidden states arising from the S4D and S6 SSMs defined in (6) and (7). The sensitivity of the $n$-th component of their states at time $t$ with respect to the input at time $j \ll t$ is given by, respectively,*

$$
\left| \frac{\partial h_t^{S4D}}{\partial x_j} \right| = \tilde{c}(\lambda_n, B_n, x_{\leq j}) \, e^{-\lambda_n (t-(j+1))},
$$

$$
\left| \frac{\partial h_t^M}{\partial x_j} \right| = \tilde{c}(\Delta, \lambda_n, B_n, x_{\leq j}) \, e^{-\lambda_n \sum_{r=j+1}^t \Delta(x_r)},
$$

*where $\tilde{c}(\Delta, \lambda_n, B_n, x_{\leq j})$ depends on the input subsequence $x_{\leq j}$, independent of the sequence length $t$.*

Lemma 1 shows that both S4D and S6 have exponential decay of sensitivity when the sequence length $t$ increases. For S4D, the only mitigation strategy is to set $\lambda_n \to 0$ (and thus $e^{\lambda_n} \to 1$). While S6 can implement the same strategy, it can also counteract the decay by adapting the input-dependent discretization, as formalized in Lem. 2.

**Lemma 2.** *Consider the discrete-time S6 in (7) where $\overline{\boldsymbol{B}}_t = [B_1(x_t), \ldots, B_n(x_t)]^\top \in \mathbb{R}^N$. Suppose there exists a constant $c \geq 0$ such that*

$$
\lim_{t \to \infty} \lambda_n \sum_{r=1}^t \Delta(x_r) \leq c.
\tag{12}
$$

*Then the sensitivity of the $n$-th component of the state at time $t$ with respect to any input $x_j$ is lower bounded by*

$$
\lim_{t \to \infty} \left| \frac{\partial h_t^M}{\partial x_j} \right| \geq e^{-c} \left| \frac{\partial}{\partial x_j} B_n(x_s) \, \Delta(x_s) \, x_s \right|.
\tag{13}
$$

This implies that, to retain sensitivity for longer sequences, Mamba must necessarily push $\lambda \Delta(x_t) \to 0$ (or equivalently, $e^{-\lambda \Delta(x_t)} \to 1$). We illustrate this with the KEEP $n$-TH task shown in Fig. 1, where we report the distribution of the learned parameters $e^{-\lambda \Delta(\boldsymbol{x}_t)}$ as we increase the sequence length. The shift towards 1 appears clear, validating the condition discussed in Lem. 2.

**Remark 1.** *While in this section we mainly focus on the properties of the* original *Mamba mixer layer (Gu & Dao, 2023), we note that the results proven in Thm. 1, Lem. 1 and Lem. 2 hold analogously for the Mamba-2 mixer (Dao & Gu, 2024). We remind that the Mamba-2 mixer layer can be seen as a simplification of Mamba's, whereby the state matrix is parameterized by a single scalar, $\boldsymbol{\Lambda} = \lambda \boldsymbol{I}$, rather than its full diagonal. Nonetheless, we do not rely on Mamba's additional flexibility for our derivations; choosing a suitable scalar $\lambda$ suffices (see details in App. A.2).*

## 5. Full Mamba Architecture Analysis

In this section, we describe the full Mamba architecture and study how its SSM mixer coordinates with the other components in the model to efficiently solve associative-recall tasks. A full Mamba architecture includes an embedding layer, a number of Mamba mixer blocks, and an output layer. The mixer block is further composed of a short convolution, an SSM, and a gate — we refer to Tab. 2 for details in the differences between Mamba and Mamba-2, but point out that Mamba-2 leverages *three independent* short convolutions (rather than a *single common* one) to compute its SSM parameters, as outlined in (14b).
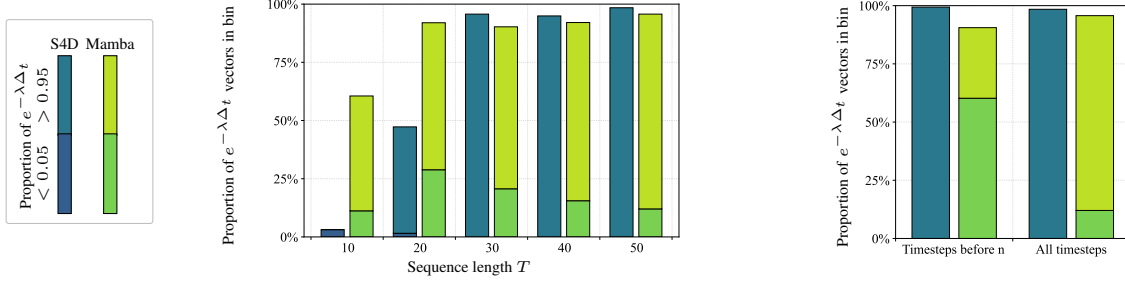
*Figure 1.* Distribution of $e^{-\lambda\Delta_t}$ computed on test inputs by models trained to successfully solve KEEP $n$-TH tasks for various sequence lengths $T$ (same setup as Tab. 1). The left histogram confirms that the Mamba model must push $e^{-\lambda\Delta_t} \to 1$ as $T$ increases, as implied by Lem. 2, to retain information throughout the sequence. Also the S4D must behave similarly. Meanwhile, the right histogram shows that, compared to S4D, Mamba has additional flexibility in forgetting irrelevant information ($e^{-\lambda\Delta_t}$ is mostly 0 before timestep $n$) and memorizing information selectively ($e^{-\lambda\Delta_t}$ is mostly 1 from timestep $n$ onwards), while S4D is forced to memorize indiscriminately.

*Table 2.* Comparison of Mamba (including S4D as a special case) and Mamba-2 (single-head) mixers. We denote with $\odot$ the Hadamard (elementwise) product, $\otimes$ the Kronecker (outer) product, $\mathrm{conv}(\boldsymbol{x})_{[t]}$ the $t$-th output of the convolution, and $\sigma$ the pointwise nonlinearity.

| MAMBA | | MAMBA-2 | |
|---|---|---|---|
| $\boldsymbol{\Lambda} \in \mathbb{R}^{d \times N}$, $\boldsymbol{x}_t \in \mathbb{R}^d$, $\boldsymbol{h}_t \in \mathbb{R}^{d \times N}$ | | $\boldsymbol{\Lambda} = \lambda \in \mathbb{R}$, $\boldsymbol{x}_t \in \mathbb{R}^d$, $\boldsymbol{h}_t \in \mathbb{R}^{d \times N}$ | (14a) |
| $\hat{\boldsymbol{x}}_t = \sigma(\mathrm{conv}(\boldsymbol{x})_{[t]}) \in \mathbb{R}^d$ | | $\hat{\boldsymbol{x}}_t = \sigma(\mathrm{conv}_u(\mathrm{Linear}(\boldsymbol{x}))_{[t]}) \in \mathbb{R}^d$ | (14b) |
| $\Delta_t = \mathrm{SoftPlus}(\mathrm{Linear}(\hat{\boldsymbol{x}}_t)) \in \mathbb{R}^d \quad (\Delta_t^{\mathrm{S4D}} = \mathbf{1})$ | | $\Delta_t = \mathrm{SoftPlus}(\mathrm{Linear}(\boldsymbol{x}_t)) \in \mathbb{R}$ | (14c) |
| $\boldsymbol{B}_t = \mathrm{Linear}(\hat{\boldsymbol{x}}_t) \in \mathbb{R}^N \quad (B_t^{\mathrm{S4D}} = B)$ | | $\boldsymbol{B}_t = \sigma(\mathrm{conv}_B(\mathrm{Linear}(\boldsymbol{x}))_{[t]}) \in \mathbb{R}^N$ | (14d) |
| $\boldsymbol{C}_t = \mathrm{Linear}(\hat{\boldsymbol{x}}_t) \in \mathbb{R}^N \quad (C_t^{\mathrm{S4D}} = C)$ | | $\boldsymbol{C}_t = \sigma(\mathrm{conv}_C(\mathrm{Linear}(\boldsymbol{x}))_{[t]}) \in \mathbb{R}^N$ | (14e) |
| $\boldsymbol{h}_t = e^{\boldsymbol{\Lambda} \odot (\Delta_t \otimes \mathbf{1}_N)} \odot \boldsymbol{h}_{t-1} + (\Delta_t \odot \hat{\boldsymbol{x}}_t) \otimes \boldsymbol{B}_t$ | | $\boldsymbol{h}_t = e^{\lambda \Delta_t} \boldsymbol{h}_{t-1} + (\Delta_t \hat{\boldsymbol{x}}_t) \otimes \boldsymbol{B}_t$ | (14f) |
| $\boldsymbol{y}_t = \boldsymbol{h}_t \boldsymbol{C}_t \in \mathbb{R}^d, \quad \tilde{\boldsymbol{y}}_t = g(\boldsymbol{x}_t) \odot \boldsymbol{y}_t := \sigma(\mathrm{Linear}(\boldsymbol{x}_t)) \odot \boldsymbol{y}_t \in \mathbb{R}^d$ | | | (14g) |

While Mamba and Mamba-2 achieve performance competitive with Transformers and outperform their SSM predecessors in solving MQAR and INDUCTION HEADS, the details of how this solution can be assembled by the architecture remain elusive, with only lower bounds on the SSM mixer size available in the literature (Arora et al., 2024b; Sanford et al., 2024b), lacking the consideration of other components of Mamba such as convolution and gating. Here we close this gap by providing analytical constructions for a 1-layer Mamba model that can exactly solve these tasks for any input. Perhaps counterintuitively, as we prove in Thm. 4, the components of a single Mamba mixer block are already powerful enough to solve MQAR exactly, even just using an S4D mixer layer (replacing the S6). With Thm. 2 and Thm. 3 we further show how, thanks to the input selectivity of S6, Mamba and Mamba-2 can leverage leaner mechanisms to solve MQAR. Particularly, the S6 layer can use $\boldsymbol{B}_t$ and $\boldsymbol{C}_t$ to efficiently structure information within its hidden state, and retrieve it when required. Finally, we show how the ability to structure the hidden state (used for MQAR) can be combined with the capacity to dynamically adjust the rate of memory decay (investigated for KEEP $n$-TH) to exactly solve INDUCTION HEADS with a variant of S6. We name this variant Mamba-$\Delta^\top$, and discuss it in Sec. 5.2.

We remark that the constructions described in this section are just *possible* solutions that the Mamba architectures can implement, and we do not exclude the existence of alternative ones. Nonetheless, in Fig. 2 we verify empirically that our solutions are tight in terms of model size.

### 5.1. 1-Layer Mamba Can Solve MQAR

The MQAR task (Arora et al., 2024a) prescribes input and output sequences as follows

$$\boldsymbol{x} = [\underbrace{k_1, v_1, \ldots, k_\kappa, v_\kappa,}_{\kappa \text{ key-value pairs}} \mid \underbrace{\ldots, k_{i_1}, \ldots, k_{i_\kappa}, \ldots}_{\text{shuffled keys, interwoven with noise}} ],$$

$$\boldsymbol{y} = [\times, \times, \ldots, \times, \times, \mid \ldots, v_{i_1}, \ldots, v_{i_j}, \ldots ].$$

The keys $k_i$ are randomly chosen from a key set of size $\kappa$, whereas the values $v_i$ and the noise are randomly taken from a vocabulary of size $|V|$. The goal is to correctly predict the value associated with the corresponding key at the query positions, while other non-query positions (denoted with $\times$) are ignored.

In this section we prove that the MQAR task can be exactly solved by three architectures: vanilla Mamba (Thm. 2), Mamba-2 (Thm. 3), and Mamba with an S4D mixer (Thm. 4), which we refer to as Mamba-S4D. Next we pro-

*Table 3.* Overview of exact solutions to the MQAR task that can be implemented by the Mamba model with S4D-mixer in Thm. 4 (top), Mamba-mixer in Thm. 2 (middle), and Mamba-2-mixer in Thm. 3 (bottom). While S4D-mixer lacks input selectivity in its SSM layer, it can solve MQAR via the gated-convolution mechanism on a larger embedding space (top). By contrast, both Mamba and Mamba-2 can solve MQAR with the same selective SSM layer construction *without gating*, and differ on the choice of convolutions (middle, bottom).

| S4D | ■ Embedding | ■ Convolution | ■ SSM | ■ Gate |
|---|---|---|---|---|
| | Encodings for any (key, value) combination can be uniquely identified. | Combines pairs with non-linear size 2 kernels. Filters out all except for (key, value)-pairs. | Hidden state is a $\kappa \times |V|$ vector, summing over all the (key, value)-pairs. | Retrieves chunk of hidden state corresponding to key. |

| MAMBA | ■ Embedding | ■ Convolution | ■ SSM | |
|---|---|---|---|---|
| | Orthogonal encodings. Keys have higher weights than values. | Combines pairs with non-linear size 2 kernels. Implements: $\hat{x}_t := \mathrm{ReLU}\left((x_{t-1} + 2x_t) - 1\right)$ | The hidden state columns are indexed by the keys, storing associated values. Retrieval selects appropriate column. $\Omega$ denotes a scalar $\geq 9$. | |

| MAMBA-2 | ■ Embedding | ■ Convolution | ■ SSM | |
|---|---|---|---|---|
| | One-hot orthogonal encodings. | There are 3 non-linear convolutions with kernel size 2. $\hat{x} := \mathrm{Identity}(x_t)$  $\hat{x}_t^B := \mathrm{Shift}(x_t)$ | The hidden state columns are indexed by the keys, storing associated values *only* (due to shift-1 convolution). Retrieval selects appropriate column. | |

vide a sketch of the proofs, which we further illustrate in Tab. 3, and defer the full details to App. C.

**Theorem 2.** *There exists a 1-layer Mamba model without gating that solves MQAR with $\kappa$ pairs using embedding size $d = O(\kappa + \log|V|)$, and state size $N = \kappa$.*

*Proof sketch.* We first prove a construction using standard basis embeddings with $d = \kappa + |V|$, and then apply Johnson-Linderstrass (JL) Lemma (see Lem. 4) to reduce dimensionality. We specify the Mamba model as follows. A size-2 (nonlinear) convolution kernel combines token *pairs*, extracting key-value pairs and removing non-informative ones (e.g., value-key, value-value). The SSM layer (S6) organizes the hidden state *matrix* (14f) via $\boldsymbol{B}_t$ (14d) such that each column corresponds to a specific key, and holds the value embedding associated with said key (hence the state size $N = \kappa$). At query time, $\boldsymbol{C}_t$ (14e) uses the query(=key) embedding to retrieve the desired value from the correct column in the hidden state matrix. □

**Theorem 3.** *There exists a 1-layer Mamba-2 model without gating that solves MQAR with $\kappa$ pairs using embedding size $d = O(\log\kappa + \log|V|)$, and state size $N = \log\kappa$.*

*Proof sketch.* The construction is similar to Thm. 2 for Mamba, but with some simplifications due to the additional flexibility provided by *three* independent convolutions in Mamba-2. Specifically, we set $\mathrm{conv}_B = (1, 0)$ as a shift-1 convolution kernel (i.e., shifting the input sequence to the right by one position), while $\mathrm{conv}_u \equiv \mathrm{conv}_C$ as the identity maps. Once again, JL Lemma allows us to reduce dimensionalities for both keys and values embeddings. □

**Remark 2.** *Our construction of Mamba-2 is similar to the construction in (Li et al., 2025) that shows how convolution-augmented Transformers can solve MQAR. Indeed, in light of the convolution layer in Tab. 2, we can interpret Mamba-2 as a convolution-augmented subquadratic Transformer.*

**Remark 3.** *With the results from Thm. 2 and 3, we can infer that the extra convolutions included in the Mamba-2 layer over vanilla Mamba in general allow for a more parameter-efficient solution of MQAR, by noting that $d = \log\kappa + \log|V| < \kappa + \log|V|$, and $N = \log\kappa < \kappa$.*

**Theorem 4.** *There exists a 1-layer Mamba model with an S4D mixer that solves MQAR with $\kappa$ pairs, using embedding size $d = O(\kappa \log|V|)$, and state size $N = 1$.*
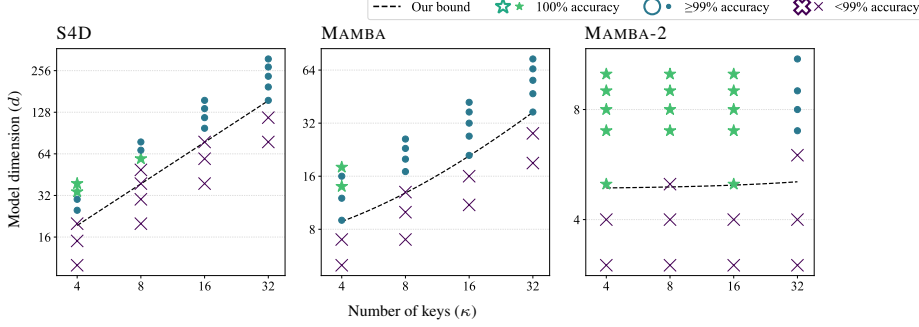
*Figure 2.* Trained models accuracy on MQAR task (best of 7 seeds), varying $\kappa$ and $d$. For S4D $N = 4$, for Mamba $N = 2\kappa$, and for Mamba-2 $N = 8\ln\kappa$. We use $T = 100$ and $|V| = 128$ for all runs. The theoretical bounds on model size for assembling the solutions proposed in Thm. 2 to 4 (black lines) separate reasonably well models that can achieve 100% accuracy (above black lines) from those that do not (below). In terms of model size efficiency, Mamba-2 is better than Mamba, which in turn is better than S4D.
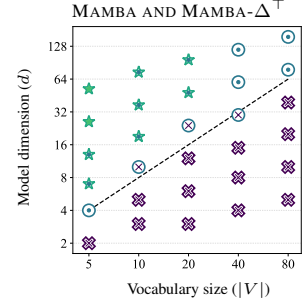
*Figure 3.* Trained models accuracy on INDUC-TION HEADS task (best of 5 seeds), varying $|V|$ and $d$, $N = 4|V|$. Mamba-$\Delta^\top$'s performance (outlined) is equal or better than Mamba's (filled) and only hits 100% above the theoretical bound from Lem. 3 (black).

*Proof sketch.* We first present the construction using $d = O(\kappa|V|)$, and again apply JL Lemma to reduce to $d = O(\kappa \log |V|)$. To overcome the limitation given by the lack of input selectivity in the SSM layer, we organize the hidden state along the embedding dimension only, and partition it so that each chunk holds the value associated with a specific key. At query time, the retrieval of the desired key-value chunk is done by leveraging the gating layer. □

**Remark 4.** *Our construction of the S4D-mixer relies both on gated convolution and the SSM recurrence. Arora et al. (2024a) provided a construction to solve* MQAR *based on gated convolution only. In contrast, our construction relies on size-2 kernels only, thanks to the SSM recurrence.*

To validate our theoretical constructions, we train Mamba, Mamba-2, and Mamba-S4D models with varying embedding size $d$, on MQAR tasks with different number of key-value pairs $\kappa$. The goal is to check how tight in practice are the theoretical bounds on model dimension derived above, and whether indeed trained models must respect them to solve the tasks. Results are reported in Fig. 2, where dashed curves denote our theoretical bounds and markers indicate empirical results. Notice our bounds in Thm. 2 to 4 are close to the empirical threshold between models sizes that can recover exact solutions or not, illustrating the tightness of our theoretical results. See additional details in App. D.3.

## 5.2. 1-Layer Mamba-$\Delta^\top$ Can Solve INDUCTION HEADS

The INDUCTION HEADS task was first introduced by Olsson et al. (2022) to study Associative Recall capabilities in Transformers. Here we use the formulation from Sanford et al. (2024a): given an input sequence of tokens $[x_1, \dots, x_t]$ from a finite vocabulary $x_t \in V$, the goal is to report, for each $x_t$, the token coming immediately after the latest previous occurrence in the input of token $x_i$. That is, the output $y_t$ must be $y_t = x_{j(t)+1}$ where $j(t) = \max\{j : j < t, x_j = x_t\}$ (or a "blank" token, $y_i = \times$, if $x_t$ appears for the first time).

Note the INDUCTION HEADS task is similar to MQAR, but with two significant differences. (i) There is no logical distinction between keys and values, so the model needs to identify the role of each token, but this can be handled by the short convolution, as we will show. More importantly, (ii) we need to retain information about only the *latest* previous occurrence of a token, so the model should dynamically forget and remember information pertaining different tokens. Since the latter memorization ability was investigated in Sec. 4.2, it is natural to leverage those insights in our solution. To do so efficiently, we introduce a slight variation to the S6 layer: the Mamba-$\Delta^\top$ SSM mixer, prescribing the following hidden state evolution

$$\boldsymbol{h}_t = e^{\boldsymbol{\Lambda} \odot (\mathbf{1}_d \otimes \Delta(\hat{\boldsymbol{x}}_t))} \odot \boldsymbol{h}_{t-1} + \hat{\boldsymbol{x}}_t \otimes \boldsymbol{B}_t. \quad (15)$$

Comparing this to (14f), the only difference lies in the action of $\Delta(\hat{\boldsymbol{x}}_t)$, which now varies along the *state* dimension $N$, rather than the *embedding* dimension $d$. While Gu & Dao (2023) hypothesized a similar performance for both versions, our findings reveal that the dependence along state dimension is better suited to solving the INDUCTION HEADS task:

**Lemma 3.** *There exists a 1-layer Mamba model with the Mamba-$\Delta^\top$ SSM mixer (15) that solves* INDUCTION HEADS *with vocabulary $V$ using embedding size $d = 2|V|$ and state size $N = |V|$.*

*Proof sketch.* The proof follows closely the MQAR construction, in that we leverage the matrix structure in the hidden state such that its columns are indexed by the keys and store the associated values. In the INDUCTION HEADS task, though, each token $x_i$ acts as key in the $(x_i, x_{i+1})$ pair, and as value in the $(x_{i-1}, x_i)$ pair. To handle this distinction,

we simply duplicate the embedding, and let the convolution layer correctly combine tokens information pairwise, so that after convolution each token encapsulates information both regarding a value and its preceding key. Moreover, we use $\Delta(x_t)$ to selectively erase outdated information, or to retain currently valid information, depending on the input observed. Pushing $\Delta(x_t) \to \infty$ flushes a previously memorized value, while $\Delta(x_t) \to 0$ preserves it. We refer to App. C.2 for the detailed proof. □

**Remark 5.** *To solve* INDUCTION HEADS*, Bietti et al. (2024); Sanford et al. (2024b) constructed* 2*-layer Transformers relying on PE and with size scaling as sequence length. On the other hand, we propose a* 1*-layer Mamba composing a convolution and a variant SSM layer. Notably, this allows us to drop the PE and thus have the model size depend only on* $|V|$*, and* not *the sequence length, improving upon the constructions for Transformers.*

To demonstrate the efficiency of our Mamba-$\Delta^\top$ variant on the INDUCTION HEADS task, we compare it against the Mamba baseline, and report results in Fig. 3. For all model sizes considered, Mamba-$\Delta^\top$ performs equally or better, demonstrating that selectivity along the state dimension (in the state matrix) improves Mamba's ability to solve the INDUCTION HEADS task.

## 6. Conclusion and Future Work

In this work, we demystify the role of input selectivity in Mamba, showing its impact on approximation power, long-term memory, and associative recall capabilities. We prove that the S6 layer can efficiently represent discontinuous signals and adaptively mitigate sensitivity decay. We also uncover the role of other architectural components in Mamba, particularly convolution and gating. We present a mechanistic explanation of how Mamba solves memorization and associative recall tasks, with tight theoretical model size bounds matching empirical results. Our findings reveal opportunities to further improve Mamba, such as an alternative way to inject input dependence within the SSM state matrix.

Our current theory does not consider the aspects of optimization and generalization, both of which are interesting future directions to explore. Moreover, our analysis focuses on simple associative recall tasks; extending it to more complicated tasks such as $k$-HOP INDUCTION HEADS (Sanford et al., 2024b), SEQUENTIAL FUNCTION COMPOSITION (Chen et al., 2024), and POINTER VALUE RETRIEVAL (Zhang et al., 2021) would be a natural next step. Overall, our work and proposed improvements add to the growing understanding of SSMs and could accelerate their development.

## Impact Statement

The goal of this paper is to improve the understanding of the Mamba architecture, specifically through analyzing the role of input selectivity. Our findings contribute to the advancement of State Space Models, which may in turn further democratize access to Large Language Models, sharpening both the existing positive and negative aspects of LLMs. No additional societal impact is expected from this work.

## References

Arora, S., Eyuboglu, S., Timalsina, A., Johnson, I., Poli, M., Zou, J., Rudra, A., and Re, C. Zoology: Measuring and Improving Recall in Efficient Language Models. In *International Conference on Learning Representations*, 2024a.

Arora, S., Eyuboglu, S., Zhang, M., Timalsina, A., Alberti, S., Zou, J., Rudra, A., and Re, C. Simple Linear Attention Language Models Balance the Recall-Throughput Trade-off. In *International Conference on Machine Learning*, volume 235, pp. 1763–1840, 2024b.

Bietti, A., Cabannes, V., Bouchacourt, D., Jegou, H., and Bottou, L. Birth of a Transformer: A Memory Viewpoint. *Advances in Neural Information Processing Systems*, 36, 2024.

Chen, L., Peng, B., and Wu, H. Theoretical Limitations of Multi-Layer Transformer. arXiv preprint, 2024. URL https://arxiv.org/abs/2412.02975.

Chiang, D. and Cholak, P. Overcoming a Theoretical Limitation of Self-Attention. In *Annual Meeting of the Association for Computational Linguistics*, volume 1, pp. 7654–7664, 2022. doi: 10.18653/v1/2022.acl-long.527.

Cirone, N. M., Orvieto, A., Walker, B., Salvi, C., and Lyons, T. Theoretical Foundations of Deep Selective State-Space Models. In *Advances in Neural Information Processing Systems*, volume 37, 2024.

Dahleh, M., Dahleh, M. A., and Verghese, G. Lectures on Dynamic Systems and Control. MIT OpenCourseWare, 2011. URL https://ocw.mit.edu/courses/6-241j-dynamic-systems-and-control-s

pring-2011/996025f6db0d90b00f11c44fc4
9b85f9_MIT6_241JS11_textbook.pdf.

Dao, T. and Gu, A. Transformers are SSMs: Generalized Models and Efficient Algorithms Through Structured State Space Duality. In *International Conference on Machine Learning*, volume 235, pp. 10041–10071, 2024.

Eckhoff, K. S. Accurate and Efficient Reconstruction of Discontinuous Functions from Truncated Series Expansions. *Mathematics of Computation*, 61(204):745–763, 1993. doi: 10.1090/S0025-5718-1993-1195430-1.

Gautschi, W. and Inglese, G. Lower Bounds for the Condition Number of Vandermonde Matrices. *Numerische Mathematik*, 52:241–250, 1987. doi: 10.1007/BF013988 78.

Grazzi, R., Siems, J., Franke, J. K., Zela, A., Hutter, F., and Pontil, M. Unlocking State-Tracking in Linear RNNs Through Negative Eigenvalues. In *International Conference on Learning Representations*, 2025.

Gu, A. and Dao, T. Mamba: Linear-Time Sequence Modeling with Selective State Spaces. arXiv preprint, 2023. URL https://arxiv.org/abs/2312.00752.

Gu, A., Dao, T., Ermon, S., Rudra, A., and Ré, C. HiPPO: Recurrent Memory with Optimal Polynomial Projections. *Advances in Neural Information Processing Systems*, 33: 1474–1487, 2020.

Gu, A., Goel, K., Gupta, A., and Ré, C. On the Parameterization and Initialization of Diagonal State Space Models. *Advances in Neural Information Processing Systems*, 35: 35971–35983, 2022a.

Gu, A., Goel, K., and Ré, C. Efficiently Modeling Long Sequences with Structured State Spaces. In *International Conference on Learning Representations*, 2022b.

Kamradt, G. Needle In A Haystack - Pressure Testing LLMs. Github, 2023. URL https://github.com/gkamr adt/LLMTest_NeedleInAHaystack. [Accessed 2025-05-29].

Katharopoulos, A., Vyas, A., Pappas, N., and Fleuret, F. Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention. In *International Conference on Machine Learning*, pp. 5156–5165, 2020.

Kingma, D. P. and Ba, J. Adam: A Method for Stochastic Optimization. arXiv preprint, 2014. URL https://ar xiv.org/abs/1412.6980.

Li, M., Zhang, X., Huang, Y., and Oymak, S. On the Power of Convolution Augmented Transformer. In *Conference on Artifical Intelligence*, number 17, pp. 18393–18402, 2025. doi: 10.1609/aaai.v39i17.34024.

Li, Z., Han, J., Weinan, E., and Li, Q. Approximation and Optimization Theory for Linear Continuous-Time Recurrent Neural Networks. *Journal of Machine Learning Research*, 23(42):1–85, 2022.

Mallat, S. Group Invariant Scattering. *Communications on Pure and Applied Mathematics*, 65(10):1331–1398, 2012. doi: 10.1002/cpa.21413.

Merrill, W., Petty, J., and Sabharwal, A. The Illusion of State in State-Space Models. In *International Conference on Machine Learning*, 2024.

Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen, A., Conerly, T., Drain, D., Ganguli, D., Hatfield-Dodds, Z., Hernandez, D., Johnston, S., Jones, A., Kernion, J., Lovitt, L., Ndousse, K., Amodei, D., Brown, T., Clark, J., Kaplan, J., McCandlish, S., and Olah, C. In-context Learning and Induction Heads. *Transformer Circuits Thread*, 2022. URL https://transformer-cir cuits.pub/2022/in-context-learning-a nd-induction-heads/index.html.

Orvieto, A., De, S., Gulcehre, C., Pascanu, R., and Smith, S. L. Universality of Linear Recurrences Followed by Non-linear Projections: Finite-Width Guarantees and Benefits of Complex Eigenvalues. In *International Conference on Machine Learning*, 2024.

Sanford, C., Hsu, D., and Telgarsky, M. One-Layer Transformers Fail to Solve the Induction Heads Task. arXiv preprint, 2024a. URL https://arxiv.org/abs/ 2408.14332.

Sanford, C., Hsu, D., and Telgarsky, M. Transformers, Parallel Computation, and Logarithmic Depth. In *International Conference on Machine Learning*, 2024b.

Sarrof, Y., Veitsman, Y., and Hahn, M. The Expressive Capacity of State Space Models: A Formal Language Perspective. *Advances in Neural Information Processing Systems*, 37:41202–41241, 2024.

Tóth, R., Felici, F., Heuberger, P., and Van den Hof, P. Crucial Aspects of Zero-Order Hold LPV State-Space System Discretization. *IFAC Proceedings Volumes*, 41(2): 4952–4957, 2008.

Vershynin, R. *High-Dimensional Probability: An Introduction with Applications in Data Science*. 2018. doi: 10.1017/9781108231596.

Vetterli, M. Wavelets, Approximation, and Compression. *IEEE Signal Processing Magazine*, 18(5):59–73, 2001. doi: 10.1109/79.952805.

Wang, S., Li, Z., and Li, Q. Inverse Approximation Theory for Nonlinear Recurrent Neural Networks. In *International Conference on Learning Representations*, 2024.

Zhang, C., Raghu, M., Kleinberg, J., and Bengio, S. Pointer Value Retrieval: A New Benchmark for Understanding the Limits of Neural Network Generalization. arXiv preprint, 2021. URL https://arxiv.org/abs/2107.12580.

## A. Approximation Power of Mamba

### A.1. Notation

We typically use bold upper case $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}$ to denote matrices and bold lower case $\boldsymbol{x}, \boldsymbol{y}$ to denote vectors or sequence. In Sec. 3 and Sec. 4, the hidden state at time $t$ is denoted as $\boldsymbol{h}(t)$ (in the continuous setting) or $\boldsymbol{h}_t$ (in the discrete setting). In Sec. 5, with a slight abuse of notation, the hidden state $\boldsymbol{h}_t \in \mathbb{R}^{d \times N}$ denotes a matrix. We use $\overline{\boldsymbol{A}}, \overline{\lambda}$ to denote the discretized versions of $\boldsymbol{A}, \lambda$. We use $\mathbb{R}, \mathbb{N}$ to denote the reals and the natural number. The identity matrix is denoted as $\boldsymbol{I}$, where the all-ones vector is denoted as $\boldsymbol{1}$. We let $\mathrm{diag}(\boldsymbol{v})$ be the diagonal matrix with diagonal filled with the vector $\boldsymbol{v}$. We denote $\mathrm{SoftPlus}, \mathrm{ReLU}, \mathrm{SiLU}$ as the corresponding pointwise nonlinearity $\sigma$, and $\mathrm{Linear}$ as the linear layer. We let $\mathbb{1}$ be the indicator function, $H(s)$ be the heaviside function. We denote with $\odot$ the Hadamard (elementwise) product and $\otimes$ the Kronecker (outer) product. We use $d$ for embedding size, $N$ for state size, and $t$ or $T$ for sequence length.

### A.2. Mamba Approximates Haar Wavelets

In the following, we recall and outline the complete proof for Thm. 1.

**Theorem 1.** *Consider a Haar wavelet $\psi_{j,k} : [0,1] \to \mathbb{R}$, and the Mamba basis function (8) at $t = 1$, $g_{j,k}^M(s; 1, x) = e^{-\lambda_{j,k} \int_s^1 \Delta_{j,k}(x_r) dr} B_{j,k}$. Let $\tilde{x}_s := \mathrm{concat}[x_s; s]$ be the input signal augmented with time positional encoding. For any $\epsilon > 0$, there exist 3 Mamba basis functions $g_{j,k}^{M_1}, g_{j,k}^{M_2}, g_{j,k}^{M_3}$ such that the approximation error*

$$\left| \psi_{j,k}(s) - \left( g_{j,k}^{M_1}(s; 1, \tilde{x}) + g_{j,k}^{M_3}(s; 1, \tilde{x}) - 2g_{j,k}^{M_2}(s; 1, \tilde{x}) \right) \right|$$

*is smaller than $\epsilon$, $\forall s \in [0,1]$.*

*Proof.* The Haar wavelet at scale $j$ with translation $k$ is defined on the interval $s \in [0,1]$ as

$$\psi_{j,0}(s) = \begin{cases} 2^{j/2} & s \in [0, 2^{-(j+1)}) \\ -2^{j/2} & s \in [2^{-(j+1)}, 2^{-j}) , \\ 0 & \text{otherwise} \end{cases} \tag{16}$$

$$\psi_{j,k}(s) = \psi_{j,0}(s - 2^{-j}k) = \begin{cases} 2^{j/2} & s \in [2^{-j}k, 2^{-j}k + 2^{-(j+1)}) \\ -2^{j/2} & s \in [2^{-j}k + 2^{-(j+1)}, 2^{-j}(k+1)) . \\ 0 & \text{otherwise} \end{cases} \tag{17}$$

See also Fig. A.1 for an illustration. Notice that each wavelet can be represented as a linear combination of three shifted Heaviside functions, namely:

$$\psi_{j,k}(s) = 2^{\frac{j}{2}} \left( H\left(s - 2^{-j}k\right) - 2H\left(s - \left(2^{-j}k + 2^{-(j+1)}\right)\right) + H\left(s - 2^{-j}(k+1)\right) \right), \tag{18}$$

where $H(s) = \mathbb{1}_{s>0}$ denotes the Heaviside function. The goal is then to show that Mamba can indeed reproduce a shifted Heaviside as one of its basis functions $g_{j,k}^M$, by opportunely choosing its free parameters $\lambda_{j,k}, \Delta_{j,k}$ and $B_{j,k}$ determining $g_{j,k}^M$.
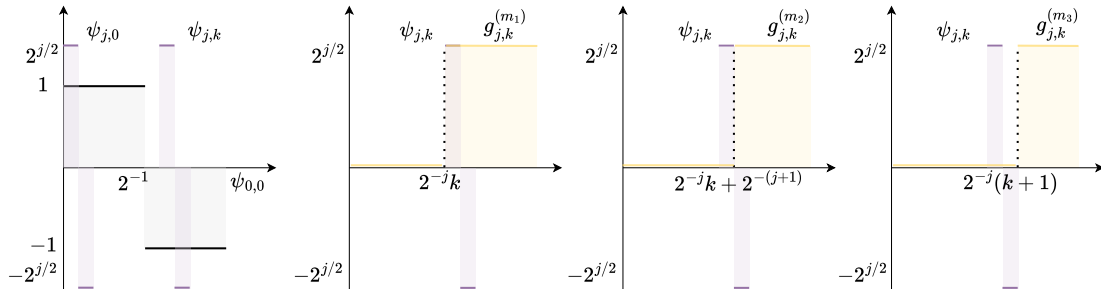


*Figure A.1.* (Left) Example of Haar Wavelets; (Middle-Right) Shape of the 3 Mamba basis $g_{j,k}^{M_1}, g_{j,k}^{M_2}, g_{j,k}^{M_3}$, whose linear combination can arbitrarily approximate the Haar wavelet $\psi_{j,k}$.

To this end, we consider a 1D input. We set $\lambda_{j,k} = 1$ and $B_{j,k} = 2^{\frac{j}{2}}$, which reduces the Mamba basis function to

$$g_{j,k}^{\mathrm{M}}(s; 1, x) = 2^{\frac{j}{2}} e^{-\int_s^1 \Delta(x(r))\, dr}. \tag{19}$$

We then leave it to the last free parameter $\Delta(x(r))$ to perform most of the heavy-lifting. By extending the input signal $x$ to include the absolute time $t$ as an input, we can directly write $\Delta(x(t)) = \Delta(t)$ (this requirement will be relaxed in Prop. 1). Notice that we can use the Mamba basis function to recover a (scaled) Heaviside centered at any given instant $\hat{t}_k$, provided $\Delta(t)$ behaves as follows:

$$\Delta(s) \to \hat{\Delta}_k(s) = \begin{cases} \infty & \text{if } s < \hat{t}_k \\ 0 & \text{else} \end{cases} \implies g_{j,k}^{\mathrm{M}}(s; 1, x) = \begin{cases} 0 & \text{if } s < \hat{t}_k \\ 2^{\frac{j}{2}} & \text{else} \end{cases} \to 2^{\frac{j}{2}} H(s - \hat{t}_k). \tag{20}$$

Remember from (7) that $\Delta(s) \equiv \Delta(s; b_\Delta, w_\Delta) := \mathrm{SoftPlus}(w_\Delta s + b_\Delta)$. Note that

$$\lim_{x \to -\infty} \mathrm{Softplus}(x) = 0, \qquad \lim_{x \to \infty} \mathrm{Softplus}(x) = \infty. \tag{21}$$

Then, by choosing $b_\Delta = -w_\Delta \hat{t}_k$, and pushing $w_\Delta \to -\infty$, we can get arbitrarily close to approximating $\hat{\Delta}_k(s)$. Concretely, we have

$$g_{j,k}^{\mathrm{M}}(s; 1, x) := 2^{\frac{j}{2}} e^{-\int_s^1 \Delta(r; b_\Delta = -w_\Delta \hat{t}_k, w_\Delta)\, dr} \xrightarrow{w_\Delta \to -\infty} 2^{\frac{j}{2}} H\left(s - \hat{t}_k\right). \tag{22}$$

By substituting $\hat{t}_k \in \{2^{-j}k, 2^{-j}k + 2^{-(j+1)}, 2^{-j}(k+1)\}$, the resulting mamba basis functions $g_{j,k}^{\mathrm{M1}}, g_{j,k}^{\mathrm{M2}}, g_{j,k}^{\mathrm{M3}}$ can approximate arbitrarily well $H(s - 2^{-j}k), H(s - (2^{-j}k + 2^{-(j+1)})), H(s - 2^{-j}(k+1))$, respectively, which are precisely the shifted Heaviside functions appearing in (18). This allows us to finally write that $\forall \epsilon, \exists w_\Delta$ such that

$$\left| \psi_{j,k}(t) - \left( g_{j,k}^{\mathrm{M1}} + g_{j,k}^{\mathrm{M3}} - 2 g_{j,k}^{\mathrm{M2}} \right) \right| < \epsilon. \tag{23}$$

$\square$

Additionally, to remove the requirement of explicitly augmenting the Mamba input with time positional encoding, we show that Mamba can autonomously recover said time position by considering a constant input $x_t \equiv 1\ \forall t$. This further relaxes the assumption in Thm. 1 to using an additional Mamba layer receiving all-ones as input.

**Proposition 1.** *A 1-layer Mamba can recover absolute time, given an all-ones input.*

*Proof.* The proof follows by considering a Mamba layer (8) with $B \equiv 1, \lambda \equiv 0$. Substituting this into (8), and providing an all-one input $x_t \equiv 1\ \forall t$ gives

$$h^{\mathrm{M}}(t) = \int_0^t x(s)\, ds = t. \tag{24}$$

This applies analogously to the discrete view, $h_t^{\mathrm{M}} = \sum_{s=1}^t x(s) = t$. $\square$

Before proving Cor. 1, we recall the function approximation problem. Given a function $\rho \in L^2([0,1])$ (as a function of time $t \in [0,1]$), and a set of bases parameterized by the SSM $\{g_n\}$ (e.g. Mamba or S4D), the goal is to project $\rho$ into the best size-$N$ basis $\mathcal{G}_N = \{g_1, \ldots, g_N\}$ that minimizes the approximation error,

$$\arg\min_{\mathcal{G}_N} \|\rho - \mathrm{proj}_{\mathcal{G}_N}(\rho)\|_{L^2} = \arg\min_{\mathcal{G}_N} \int_0^1 \left( \rho(t) - \mathrm{proj}_{\mathcal{G}_N}(\rho)(t) \right)^2 dt, \tag{25}$$

where $N$ denotes the hidden state size and

$$\mathrm{proj}_{\mathcal{G}_N}(\rho) := \sum_{g_n \in \mathcal{G}_N} \langle \rho, g_n \rangle g_n. \tag{26}$$

We note that the $N$ basis $\mathcal{G}_N$ is optimally chosen for the target function $\rho$ (instead of using a fixed set of $N$ bases independent of $\rho$). The inner products $\langle \rho, g_n \rangle = \int_0^t \rho(t) \overline{g_n(t)} dt$ for the various $n$ determine the coefficients of the projection of the target function $\rho$ onto the basis $\mathcal{G}_N$.

We are now ready to prove Cor. 1.

**Corollary 1.** *For a piecewise-constant function $\rho(t)$ with $m \geq 1$ discontinuities, there exist $N$ Mamba basis functions (8) such that the $L^2$ approximation error $\|\rho - \sum_{n=1}^{N} g_n^M\|_{L^2}$ is of order $\mathcal{O}(2^{-\frac{N}{3m}})$. On the other hand, S4D basis functions can achieve an approximation error of $\mathcal{O}(N^{-1})$.*

*Proof.* As discussed in Sec. 4.1, a S4D basis function (9) can represent any Fourier basis function $f_n = e^{i2\pi n}$. When approximating a target function with discontinuities, the Fourier coefficients decay very slowly (Eckhoff, 1993),

$$\lim_{n \to \infty} |\langle \rho, f_n \rangle| = O(n^{-1}). \tag{27}$$

This implies that using $N$ Fourier bases (and thus $N$ S4D basis functions) for approximating discontinuous functions yields an approximation error $O(N^{-1})$, proving the second part of Cor. 1.

By Thm. 1, $N$ Mamba basis functions can approximate arbitrarily close any $N/3$ number of Haar wavelets. A set of $N/3$ Haar wavelets can achieve an approximation error of $O(2^{-N/3m})$ (Vetterli, 2001) when targeting a piecewise-constant function with $m \geq 1$ discontinuities. We provide a concrete argument for completeness. Suppose the function $\rho$ is piecewise constant with $m = 1$ discontinuity. Then using an *optimal* set of $N$ Haar wavelets $\mathcal{H}_N \subset \{\psi_{j,k}\}_{j \in \mathbb{N}, 0 \leq k \leq 2^j - 1}$, we can achieve an approximation error of

$$\min_{\mathcal{H}_N} \|\rho - \text{proj}_{\mathcal{H}_N}(\rho)\|_{L^2} = \sum_{j=0}^{\infty} \sum_{k=0}^{2^j - 1} \langle \rho, \psi_{j,k} \rangle^2 - \sum_{\psi_{j,k} \in \mathcal{H}_N} \langle \rho, \psi_{j,k} \rangle^2 \tag{28}$$

$$= \sum_{j=0}^{\infty} \langle \rho, \psi_{j,k(\rho)} \rangle^2 - \sum_{j=0}^{N-1} \langle \rho, \psi_{j,k(\rho)} \rangle^2 \tag{29}$$

$$= \sum_{j=N}^{\infty} \langle \rho, \psi_{j,k(\rho)} \rangle^2 \tag{30}$$

$$= O(2^{-N}). \tag{31}$$

The second equality holds by noting that $\langle \rho, \psi_{j,k} \rangle = \int_0^1 \rho(t) \psi_{j,k}(t) dt = 0$ if $\rho$ is constant in the interval $[2^{-j}(k-1), 2^{-j}k]$. Thus, for a piecewise-constant $\rho$ with one discontinuity, only one wavelet at each scale $j$ has $\langle \rho, \psi_{j,k(x)} \rangle \neq 0$. Given that we can choose adaptively the best $N$ wavelets, we pick the ones with nonzero coefficients across all $J = N$ scales. Then the approximation error is bounded the coefficient the highest scale, which has magnitude $O(2^{-N/2})$ and squared error $O(2^{-N})$. Similar analysis shows that for piecewise-constant functions with $m$ discontinuities, the optimal adaptive wavelet basis with $N$ wavelet functions achieves an approximation error of $O(2^{-N/m})$. $\square$

**Remark 6.** *While the assumption of a piecewise-constant target function seems restrictive in Cor. 1, we note that this can be relaxed to a continuous target function $\rho$ by additionally assuming the input signals $x$ are piecewise-constant. The equivalence is due to the intermediate value theorem. Note that piecewise-constant input signals are ubiquitous in language modelling tasks, where $x_t$ takes values in a finite-size vocabulary.*

### A.3. Mamba can solve KEEP $n$-TH

Here we provide a direct application of the result from App. A.2, to explain the ability of Mamba to exactly solve the KEEP $n$-TH task when equipped with Positional Encoding.

**Corollary 2.** *There exists an S6 layer that solves KEEP $n$-TH on input augmented with time Positional Encoding.*

*Proof.* To memorize the $n$-th position of a piecewise-constant signal $x(s) = \sum_{i=1}^{t} x_i \mathbb{1}_{[i-1,i)}(s)$, it suffices to represent two Heaviside functions $h_1 = H(s - n), h_2 = H(s - (n-1))$ (recall $H(s - n) := \mathbb{1}_{s>n}$), by noting that $h_1 - h_2$ is one at the interval $[n-1, n)$ and zero elsewhere, we have $\langle x, h_1 - h_2 \rangle = x_n$, as desired.

As shown in (22), the Mamba basis function can represent any Heaviside function. Specifically, we can let $w_\Delta \to -\infty$, and

$$\Delta_1([x_t; t]) = \text{SoftPlus}(w_\Delta t - w_\Delta n) = \begin{cases} \infty & t < n \\ 0 & t \geq n \end{cases}, \tag{32a}$$

$$\Delta_2([x_t; t]) = \text{SoftPlus}(w_\Delta t - w_\Delta(n-1)) = \begin{cases} \infty & t < n-1 \\ 0 & t \geq n-1 \end{cases}. \tag{32b}$$

Suppose $\lambda \equiv 1, B = 1$ in the Mamba basis function (8). Using the above $\Delta_1, \Delta_2$, we can obtain the Mamba basis functions representing the desired Heavisides,

$$g_1^{\text{M}}(s; t, x) = e^{-\int_s^t \Delta_1(x(r))\, dr} \xrightarrow{w_\Delta \to -\infty} H(s - n). \tag{33a}$$

$$g_2^{\text{M}}(s; t, x) = e^{-\int_s^t \Delta_2(x(r))\, dr} \xrightarrow{w_\Delta \to -\infty} H(s - (n-1)). \tag{33b}$$

$\square$

# B. Sensitivity of Mamba

In this section, we prove Lem. 1 and 2. These follow directly from the general sensitivity formula (11), which we derive in more details below: Given a generic, input-dependent recurrence relationship as in (4), we have that the sensitivity of the state at time $t$ with respect to its input sequence at the $j$-th instant $x_j \in \mathbb{R}$ is given by

$$
\begin{aligned}
\frac{\partial \boldsymbol{h}_t}{\partial x_j} &= \frac{\partial}{\partial x_j} \left( \sum_{s=1}^{t} \left( \prod_{r=s+1}^{t} \overline{\boldsymbol{\Lambda}}_r \right) \overline{\boldsymbol{B}}_s x_s \right) \\
&= \sum_{s=1}^{t} \left[ \frac{\partial}{\partial x_j} \left( \prod_{r=s+1}^{t} \overline{\boldsymbol{\Lambda}}_r \right) \overline{\boldsymbol{B}}_s x_s + \left( \prod_{r=s+1}^{t} \overline{\boldsymbol{\Lambda}}_r \right) \frac{\partial}{\partial x_j} \left( \overline{\boldsymbol{B}}_s x_s \right) \right] \\
&= \sum_{s=1}^{t} \left[ \frac{\partial \overline{\boldsymbol{\Lambda}}_j}{\partial x_j} \left( \prod_{\substack{r=s+1, \\ r \neq j}}^{t} \overline{\boldsymbol{\Lambda}}_r \right) \overline{\boldsymbol{B}}_s x_s \delta_{s<j} + \left( \prod_{r=s+1}^{t} \overline{\boldsymbol{\Lambda}}_r \right) \frac{\partial}{\partial x_j} \left( \overline{\boldsymbol{B}}_s x_s \right) \right] \\
&= \left( \prod_{r=j+1}^{t} \overline{\boldsymbol{\Lambda}}_r \right) \frac{\partial}{\partial x_j} \left( \overline{\boldsymbol{B}}_j x_j \right) + \sum_{s=1}^{j-1} \frac{\partial \overline{\boldsymbol{\Lambda}}_j}{\partial x_j} \left( \prod_{\substack{r=s+1, \\ r \neq j}}^{t} \overline{\boldsymbol{\Lambda}}_r \right) \overline{\boldsymbol{B}}_s x_s \\
&= \left( \prod_{r=j+1}^{t} \overline{\boldsymbol{\Lambda}}_r \right) \left( \frac{\partial}{\partial x_j} \left( \overline{\boldsymbol{B}}_j x_j \right) + \frac{\partial \overline{\boldsymbol{\Lambda}}_j}{\partial x_j} \sum_{s=1}^{j-1} \left( \prod_{r=s+1}^{j-1} \overline{\boldsymbol{\Lambda}}_r \right) \overline{\boldsymbol{B}}_s x_s \right).
\end{aligned}
\tag{34}
$$

Since $\overline{\boldsymbol{\Lambda}}_r = -\operatorname{diag}([\overline{\lambda}_1(r), \dots, \overline{\lambda}_n(r)])$ is diagonal and $\overline{\boldsymbol{B}}_j = [\overline{B}_1(j), \dots, \overline{B}_n(j)]$, we obtain from (34) the sensitivity of the $n$-th component of the hidden state $\boldsymbol{h}_{t,n} \equiv h_t$ with respect to the input $x_j$ as

$$
\frac{\partial h_t}{\partial x_j} = \left( \prod_{r=j+1}^{t} \overline{\lambda}_n(r) \right) \left( \frac{\partial}{\partial x_j} (\overline{B}_n(j) x_j) + \frac{\partial \overline{\lambda}_n(j)}{\partial x_j} \sum_{s=1}^{j-1} \left( \prod_{r=s+1}^{j-1} \overline{\lambda}_n(r) \right) \overline{B}_n(s) x_s \right).
\tag{35}
$$

With these ingredients, we are now ready to prove the main results in this section.

**Lemma 1.** *Consider the hidden states arising from the S4D and S6 SSMs defined in (6) and (7). The sensitivity of the $n$-th component of their states at time $t$ with respect to the input at time $j \ll t$ is given by, respectively,*

$$
\left| \frac{\partial h_t^{S4D}}{\partial x_j} \right| = \tilde{c}(\lambda_n, B_n, x_{\leq j}) \, e^{-\lambda_n (t - (j+1))},
$$

$$
\left| \frac{\partial h_t^{M}}{\partial x_j} \right| = \tilde{c}(\Delta, \lambda_n, B_n, x_{\leq j}) \, e^{-\lambda_n \sum_{r=j+1}^{t} \Delta(x_r)},
$$

*where $\tilde{c}(\Delta, \lambda_n, B_n, x_{\leq j})$ depends on the input subsequence $x_{\leq j}$, independent of the sequence length $t$.*

*Proof.* Recall the Mamba discretization (5) chooses input-dependent SSM parameters as

$$
\overline{\lambda}_n(r) = e^{-\lambda_n \Delta(x_r)}, \qquad \overline{B}_n(j) = B_n(x_j) \Delta(x_j).
\tag{36}
$$

Substituting (36) to (35), we see that the first factor becomes

$$
\prod_{r=j+1}^{t} \overline{\lambda}_n(r) = e^{-\lambda_n \sum_{r=j+1}^{t} \Delta(x_r)},
\tag{37}
$$

while the second factor becomes

$$
\begin{aligned}
&\frac{\partial}{\partial x_j} (\overline{B}_n(j) x_j) + \frac{\partial \overline{\lambda}_n(r)}{\partial x_j} \sum_{s=1}^{j-1} \left( \prod_{r=s+1}^{j-1} \overline{\lambda}_n(r) \right) \overline{B}_n(s) x_s \\
&= \frac{\partial}{\partial x_j} (B_n(x_j) \Delta(x_j) x_j) - \lambda e^{-\lambda \Delta(x_j)} \frac{\partial \Delta(x_j)}{\partial x_j} \sum_{s=1}^{j-1} e^{-\lambda_n \sum_{r=s+1}^{j-1} \Delta(x_r)} B_n(x_s) \Delta(x_s) x_s.
\end{aligned}
\tag{38}
$$

16

We are interested in the behavior of the sensitivity for $j$ fixed and $t \to \infty$ (i.e., the sensitivity of the current state with respect to early input in long-range sequences). Notice that the second factor does not scale with $t$ (since $j$ is fixed), and can be bound in terms of the parameters defining the transformations $\boldsymbol{B}(x)$ and $\Delta(x)$, as well as $\lambda_n$ and $x$, as such:

$$|(38)| \leq \tilde{c}(B_n, \Delta, \lambda_n, x_{\leq j}). \tag{39}$$

On the other hand, the first factor in (37) shows in general *exponential* dependence on $t$. Putting both together, we have

$$\left| \frac{\partial h_t^{\mathrm{M}}}{\partial x_j} \right| \leq \tilde{c}(\Delta, \lambda_n, B_n, x_{\leq j}) \, e^{-\lambda_n \sum_{r=j+1}^t \Delta(x_r)}. \tag{40}$$

The proof for S4D is immediate by taking $\Delta(x_t) = 1, B_n(x_t) = B_n$ for all $t$, resulting in

$$\left| \frac{\partial h_t^{\mathrm{S4D}}}{\partial x_j} \right| \leq \tilde{c}(\lambda_n, B_n, x_{\leq j}) \, e^{-\lambda_n (t - (j+1))}. \tag{41}$$

$\square$

**Lemma 2.** *Consider the discrete-time S6 in (7) where $\overline{\boldsymbol{B}}_t = [B_1(x_t), \ldots, B_n(x_t)]^\top \in \mathbb{R}^N$. Suppose there exists a constant $c \geq 0$ such that*

$$\lim_{t \to \infty} \lambda_n \sum_{r=1}^t \Delta(x_r) \leq c. \tag{12}$$

*Then the sensitivity of the $n$-th component of the state at time $t$ with respect to any input $x_j$ is lower bounded by*

$$\lim_{t \to \infty} \left| \frac{\partial h_t^M}{\partial x_j} \right| \geq e^{-c} \left| \frac{\partial}{\partial x_j} B_n(x_s) \, \Delta(x_s) \, x_s \right|. \tag{13}$$

*Proof.* Recall for the scalar input sequence, each component of the Mamba hidden state is given by

$$h_t^{\mathrm{M}} = \sum_{s=1}^t e^{-\lambda_n \sum_{r=s+1}^t \Delta(x_r)} B_n(x_s) \, \Delta(x_s) \, x_s. \tag{42}$$

Then the condition in (12) implies

$$\lim_{t \to \infty} h_t^{\mathrm{M}} \geq \sum_{s=1}^t e^{-c} B_n(x_s) \, \Delta(x_s) \, x_s. \tag{43}$$

Straightforward computation of $\lim_{t \to \infty} \left| \frac{\partial h_t^{\mathrm{M}}}{\partial x_j} \right|$ completes the proof. $\square$

For comparison, we provide a similar sensitivity analysis for the softmax Attention layer as well.

*Proof of Sensitivity of Softmax Attention.*

$$
\begin{aligned}
\frac{\partial y_t^{\mathrm{Attn}}}{\partial x_j} &= \frac{\partial}{\partial x_j} \left( \sum_{s=1}^t \frac{e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_s}}{\sum_{r=1}^t e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_r}} \boldsymbol{W}_V x_s \right) \\
&= \frac{e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_j} (x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K \boldsymbol{W}_V x_j + \boldsymbol{W}_V)}{\sum_{r=1}^t e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_r}} - \frac{e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_j} x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K}{\left( \sum_{r=1}^t e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_r} \right)^2} \sum_{s=1}^t e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_s} \boldsymbol{W}_V x_s \\
&= \frac{e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_j}}{\sum_{r=1}^t e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_r}} \left( x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K \boldsymbol{W}_V x_j + \boldsymbol{W}_V - x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K \frac{\sum_{s=1}^t e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_s} \boldsymbol{W}_V x_s}{\sum_{r=1}^t e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_r}} \right) \\
&= \frac{e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_j}}{\sum_{r=1}^t e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_r}} \left( \boldsymbol{W}_V + x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K \left( \boldsymbol{W}_V x_j - \frac{\sum_{s=1}^t e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_s} \boldsymbol{W}_V x_s}{\sum_{r=1}^t e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_r}} \right) \right).
\end{aligned}
\tag{44}
$$

17

Even in this case, the second term can be bound by a constant $C$ in terms of $\|\boldsymbol{W}_Q\|, \|\boldsymbol{W}_K\|, \|\boldsymbol{W}_V\|$ and $\|\boldsymbol{x}\|$ (notice that the rightmost term, where a fraction of two sums over $t$ appear, behaves as a weighted average of $\boldsymbol{W}_V x_s$, and hence does *not* scale with $t$). Only the denominator at the first factor remains as a function of $t$, providing

$$\left| \frac{\partial h_t^{\texttt{Attn}}}{\partial x_j} \right| \le C\frac{1}{t} \left( \min_r e^{x_t \boldsymbol{W}_Q^\top \boldsymbol{W}_K x_r} \right)^{-1}. \tag{45}$$

$\square$

# C. Proofs of Section 5

## C.1. Proofs of Mamba Solving the MQAR Task

**The MQAR Task**   As a reminder, the MQAR task receives an input in the form

$$\boldsymbol{x} = [\underbrace{k_1 \ v_1 \ \ldots \ k_i \ v_i \ \ldots \ k_\kappa \ v_\kappa}_{\kappa \text{ key-value pairs}} \ | \ \underbrace{v_\times \ \ldots \ v_\times \ k_{i_1} \ v_\times \ \ldots \ v_\times \ k_{i_j} \ v_\times \ \ldots \ v_\times \ k_{i_\kappa} \ v_\times \ \ldots}_{\text{shuffled queries (=keys), interwoven with noise}}] \in \mathbb{R}^T, \quad (46)$$

where the keys $k_i$ are randomly taken from a key set of size $\kappa$, and the values $v_i$ and the various noise tokens $v_\times$ are randomly taken from a vocabulary of size $|V|$. The goal is to output a sequence that, at the location of each query, reports the value matching the corresponding key, that is

$$\boldsymbol{y} = [\times \ \ldots \ \times \ | \ \times \ \ldots \ \times \ v_{i_1} \ \times \ \ldots \ \times \ v_{i_j} \ \times \ \ldots \ \times \ v_{i_\kappa} \ \times \ \ldots \ \times], \quad (47)$$

while the other components of the output (denoted with $\times$) are ignored.

For the rest of the proofs, we often make use of the Johnson-Lindenstrauss (JL) Lemma to reduce the embedding dimensionality. For completeness, we state and prove JL lemma below.

**Lemma 4** (JL Lemma). *Given a set of standard bases $\{\boldsymbol{e}_1, \ldots, \boldsymbol{e}_d\}$ and $\epsilon \in (0, 0.5)$, there exists a random projection matrix $M : \mathbb{R}^d \to \mathbb{R}^p$ where $p = O(\epsilon^{-2} \log d)$, $M_{i,j} \stackrel{i.i.d.}{\sim} \frac{1}{\sqrt{p}} \mathrm{Unif}\{\pm 1\}$ such that for all pairs $(i,j)$,*

$$|\langle M\boldsymbol{e}_i, M\boldsymbol{e}_j \rangle| \leq \epsilon. \quad (48)$$

*Proof.* Let $\sqrt{p} M_{i,j} \coloneqq Z_{i,j}$, where $Z_{i,j}$ is the symmetric Bernoulli variable. For any $i, j \in [d]$, we have

$$\langle M\boldsymbol{e}_i, M\boldsymbol{e}_j \rangle = \langle M_{[:,i]}, M_{[:,j]} \rangle = \frac{1}{p} \sum_{l=1}^{p} Z_{l,i} Z_{l,j} \equiv \frac{1}{p} \sum_{l=1}^{p} Z_l,$$

where the last equality follows from the fact that the product of two independent Bernoulli variables is Bernoulli. In other words, the dot product of the two projected vectors is a sum of $k$ i.i.d. symmetric Bernoullis. By Hoeffding's inequality (e.g. Vershynin (2018, Theorem 2.2.2)),

$$\mathbb{P}\left[\frac{1}{p} \sum_{l=1}^{p} Z_l \geq \epsilon\right] \leq \exp\left(-\frac{p\epsilon^2}{2}\right). \quad (49)$$

Thus

$$\mathbb{P}[|\langle M\boldsymbol{e}_i, M\boldsymbol{e}_j \rangle| \geq \epsilon] = 2\mathbb{P}\left[\frac{1}{p} \sum_{l=1}^{p} Z_l \geq \epsilon\right] \leq 2\exp\left(-\frac{p\epsilon^2}{2}\right). \quad (50)$$

Therefore, except with probability less than $2\exp\left(-\frac{p\epsilon^2}{2}\right)$, it holds that $|\langle M\boldsymbol{e}_i, M\boldsymbol{e}_j \rangle| \leq \epsilon$. Let $p = \frac{4}{\epsilon^2} \ln \frac{d}{\delta}$ for some $\delta \in (0, 1)$. By a union bound, this holds for all $\binom{d}{2}$ pairs of $(\boldsymbol{e}_i, \boldsymbol{e}_j)$ except with probability

$$\binom{d}{2} 2\exp\left(-\frac{p\epsilon^2}{2}\right) < d^2 \exp\left(-\frac{p\epsilon^2}{2}\right) = \exp\left(-2\ln\frac{d}{\delta}\right) = \delta^2. \quad (51)$$

Since there is a probability grater than $1 - \delta^2$ that $|\langle M\boldsymbol{e}_i, M\boldsymbol{e}_j \rangle| \leq \epsilon$ holds for all $i, j$, this guarantees the existence of such $M$ by the probabilistic method.

$\square$

**Theorem 2.** *There exists a 1-layer Mamba model without gating that solves* MQAR *with $\kappa$ pairs using embedding size $d = O(\kappa + \log |V|)$, and state size $N = \kappa$.*

*Proof.* The proof is divided into two steps: We first construct a 1-layer Mamba model without gating that solves MQAR using standard basis vectors $d = O(\kappa + |V|)$; based on such construction, we then apply the JL Lemma (Lem. 4) together with a shifting trick to complete the proof for $d = O(\kappa + \log |V|)$.

19

**Step 1: Construction With** $d = (\kappa + O(|V|))$**.** We consider one-hot encoding of keys and values. The main idea is to (i) use the size-2 convolution to combine key-value pairs and filter out other uninformative pairs (particularly, preserve key-value pairs and discard value-value pairs); (ii) use appropriate input-dependent SSM matrices to store and retrieve the key-value information in the hidden state; (iii) use the output layer to project to the value embedding subspace. Crucially, we observe that the hidden state of Mamba at time $t$ is a $d \times N$ matrix (14f). Thus, we leverage this matrix structure such that each column of the state corresponds to a given key ($N = \kappa$), and holds the value associated with it, that is:

$$\boldsymbol{h} = [\quad \boldsymbol{v}_{j_1} \quad | \quad \boldsymbol{v}_{j_2} \quad | \quad \ldots \quad | \quad \boldsymbol{v}_{j_\kappa} \quad ] \in \mathbb{R}^{d \times N}. \tag{52}$$

We first describe how our proposed solution operates, and then prove that indeed our solution solves the MQAR task exactly.

- **Embedding** The task of the embedding layer is to clearly distinguish values and keys. We achieve this by letting it map to orthogonal directions: let the embedding dimension be $d = |V| + \kappa$, and $\boldsymbol{e}_i \in \mathbb{R}^d$ denote the standard basis vector. We impose

$$k_i \mapsto \boldsymbol{k}_i := k \cdot \boldsymbol{e}_i, \quad \text{and} \quad v_i \mapsto \boldsymbol{v}_i := v \cdot \boldsymbol{e}_{\kappa+i} \tag{53}$$

  for some parameters $k, v > 0$.

- **Convolution** We use a size-2 convolution to combine information of a key-value pair, as it is essential to the task solution. We remind that for a sequence $(\ldots, \boldsymbol{x}_{t-1}, \boldsymbol{x}_t, \ldots)$, the action of the convolution layer with size-2 kernel with left padding $\boldsymbol{x}_0$ is given by

$$\hat{\boldsymbol{x}}_t = \text{conv}(\boldsymbol{x}_{t-1}, \boldsymbol{x}_t) := \sigma\left(\boldsymbol{c}_0 \odot \boldsymbol{x}_{t-1} + \boldsymbol{c}_1 \odot \boldsymbol{x}_t - \boldsymbol{b}\right), \qquad \boldsymbol{x}_0 := \boldsymbol{0}, \tag{54}$$

  for certain parameters $\boldsymbol{c}_0, \boldsymbol{c}_1, \boldsymbol{b} \in \mathbb{R}^d$ and a nonlinearity[1] $\sigma = \text{ReLU}$. For our construction, it suffices to pick $\boldsymbol{c}_0 = c_0 \boldsymbol{1}, \boldsymbol{c}_1 = c_1 \boldsymbol{1}$, and $\boldsymbol{b} = b \boldsymbol{1}$, with $c_0, c_1, b > 0$. We want such nonlinear convolution to preserve information from the $(\boldsymbol{k}_i, \boldsymbol{v}_j)$ pairs, so we impose $c_0 k > b$ and $c_1 v > b$. We also need to prevent the value in $(\boldsymbol{v}_i, \boldsymbol{k}_j)$ from getting associated to the wrong key, but we want to preserve key information to be able to extract it at retrieval time, so we ask $c_0 v \leq b$ and $c_1 k > b$. The $(\boldsymbol{k}_i, \boldsymbol{k}_j)$ pair represents an edge-case, and we will show later how to deal with this. Finally, we want to ignore contributions from $(\boldsymbol{v}_i, \boldsymbol{v}_j)$ pairs (as they refer to "noise" tokens), but we will delegate this task to input-selectivity in the SSM layer. To summarize, we need to satisfy

$$c_0 k > b, \qquad c_1 v > b, \qquad c_1 k > b, \qquad c_0 v \leq b. \tag{55}$$

  A feasible parameter combination satisfying (55) is given by:

$$v = 1, \qquad k = 2, \qquad c_0 = 1, \qquad c_1 = 2, \quad \text{and} \quad b = 1. \tag{56}$$

  Let us show how such nonlinear convolution acts on the four different types of input pairs, under the assumptions (55):

$$\text{conv}((\boldsymbol{k}_i, \boldsymbol{v}_j)) = \text{ReLU}(c_0 \boldsymbol{k}_i + c_1 \boldsymbol{v}_j - \boldsymbol{b}) = (c_0 k - b)\boldsymbol{e}_i + (c_1 v - b)\boldsymbol{e}_{\kappa+j} \tag{57a}$$
$$\text{conv}((\boldsymbol{v}_i, \boldsymbol{k}_j)) = \text{ReLU}(c_0 \boldsymbol{v}_i + c_1 \boldsymbol{k}_j - \boldsymbol{b}) = (c_1 k - b)\boldsymbol{e}_j \tag{57b}$$
$$\text{conv}((\boldsymbol{k}_i, \boldsymbol{k}_j)) = \text{ReLU}(c_0 \boldsymbol{k}_i + c_1 \boldsymbol{k}_j - \boldsymbol{b}) = (c_0 k - b)\boldsymbol{e}_i + (c_1 k - b)\boldsymbol{e}_j \tag{57c}$$
$$\text{conv}((\boldsymbol{v}_i, \boldsymbol{v}_j)) = \text{ReLU}(c_0 \boldsymbol{v}_i + c_1 \boldsymbol{v}_j - \boldsymbol{b}) = (c_1 v - b)\boldsymbol{e}_{\kappa+j}. \tag{57d}$$

- **Mamba SSM** The SSM layer organizes a hidden state *matrix* where its columns are indexed by the keys and store information of the associated values. To this end, we let

$$\boldsymbol{\Lambda} = \boldsymbol{0}, \quad \Delta_t = \boldsymbol{1}, \quad \boldsymbol{B}(\boldsymbol{x}) = \boldsymbol{C}(\boldsymbol{x}) = [I_\kappa | \boldsymbol{0}_{|V|}]\boldsymbol{x} \equiv W_k \boldsymbol{x}, \tag{58}$$

  where $W_k$ projects the input to the key embedding subspace. Consequently, the SSM layer covers three roles at once:

---

[1]In the original Mamba definition, we have $\sigma \equiv \text{SiLU}$; for ease of illustration we consider instead $\sigma \equiv \text{ReLU}$, but notice that similar considerations still hold in this case, in light of the fact that $\text{SiLU}(x) \to \text{ReLU}(x)$ for $x \to \pm\infty$: it suffices then to opportunely scale the inputs.

1. Ensure that the right key get associated to the right column in the hidden state. This role is covered by the input matrix $\boldsymbol{B}(\boldsymbol{x}) = W_B \boldsymbol{x}$. By picking $W_B = W_k = [I_\kappa | \mathbf{0}_{|V|}]$, we can see that only keys in the convolved input pair are used for populating the hidden state. This choice of the input matrix also ensures that information from $(\boldsymbol{v}_i, \boldsymbol{v}_j)$ pairs does not affect the hidden state.

2. Propagate information down the sequence without corrupting it (i.e., memorization). The state matrix $A(\boldsymbol{x})$ can take care of this, provided we fix it to the all-one matrix. This can be achieved by prescribing $\boldsymbol{\Lambda} \equiv \mathbf{0}$ in (14a). Notice that with this choice, the SSM layer simply performs a cumulative sum $\boldsymbol{h}_t = \sum_{s=1}^{t} \hat{\boldsymbol{x}}_s \otimes \boldsymbol{B}_s$.

3. Retrieve the correct column when needed. This task is taken care of by the output matrix $C(\boldsymbol{x}) = W_C \boldsymbol{x}$. As with the input matrix, it suffices to have $W_C = W_k = [I_\kappa | \mathbf{0}_{|V|}]$: also in this case, when an input containing a key is encountered, the $C(\boldsymbol{x})$ matrix will retrieve information only from the column corresponding to that specific key.

- **Output** The final output layer simply needs to correctly classify $\boldsymbol{v}_i$ as the most likely value from the retrieved vector $\boldsymbol{y}_t$. To this end, it suffices to consider

$$W_o = \begin{bmatrix} \boldsymbol{v}_1^\top \\ \vdots \\ \boldsymbol{v}_{|V|}^\top \end{bmatrix}, \tag{59}$$

as $\boldsymbol{v}_i$ will return the maximum scalar product with $\boldsymbol{y}_t$.

Let us illustrate with a step-by-step example how the above construction yields the required result. Consider a generic input sequence $\boldsymbol{x}$ for the MQAR task. Under the embedding layer prescribed in (53), the (embeded) input will be in the form

$$\boldsymbol{x} = \begin{bmatrix} \boldsymbol{k}_{i_1} & \boldsymbol{v}_{j(i_1)} & \boldsymbol{k}_{i_2} & \boldsymbol{v}_{j(i_2)} & \dots & \boldsymbol{k}_{i_\kappa} & \boldsymbol{v}_{j(i_\kappa)} \mid \boldsymbol{v}_\times & \dots & \boldsymbol{k}_{k_1} & \dots & \boldsymbol{k}_{k_\kappa} & \dots \end{bmatrix} \in \mathbb{R}^{d \times T}, \tag{60}$$

so that the $i_m$-th key is associated to the $j(i_m)$-th value. The order in which the keys appear at query time is also random and denoted by the $k_m$ subscript. Moreover, keys at query time are randomly interwoven with value tokens, denoted as $\boldsymbol{v}_\times$. After convolution, this input gets mapped to

$$\hat{\boldsymbol{x}} = \begin{bmatrix} c_{1,k}\boldsymbol{e}_{i_1} & c_{0,k}\boldsymbol{e}_{i_1} & c_{1,k}\boldsymbol{e}_{i_2} & c_{0,k}\boldsymbol{e}_{i_2} & \dots & c_{1,k}\boldsymbol{e}_{i_\kappa} & c_{0,k}\boldsymbol{e}_{i_\kappa} & \mathbf{0} & \dots & c_{1,k}\boldsymbol{e}_{k_1} & \dots & c_{1,k}\boldsymbol{e}_{k_\kappa} & \dots \\ + & + & + & + & & + & + & + & & + & & + & \\ \mathbf{0} & c_{1,v}\boldsymbol{e}_{j(i_1)} & \mathbf{0} & c_{1,v}\boldsymbol{e}_{j(i_2)} & \dots & \mathbf{0} & c_{1,v}\boldsymbol{e}_{j(i_\kappa)} & \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{0} & \dots \end{bmatrix}, \tag{61}$$

where we denote $c_{0,k} = c_0 k - b$, $c_{1,k} = c_1 k - b$, $c_{1,v} = c_1 v - b$ to slim notation, and separate the components of $\hat{\boldsymbol{x}}$ into key-related (top) and value-related (bottom) for illustrative purposes. Notice how, for the initial part of the input, distinct keys (the various $\boldsymbol{e}_{i_m}$ at the top) are only associated with their respective values (the various $\boldsymbol{e}_{j(i_m)}$ at the bottom), or with $\mathbf{0}$. Before moving on to the SSM layer, let us remind its action (14f):

$$\boldsymbol{y}_t = \left(e^{\boldsymbol{\Lambda} \odot (\Delta_t \otimes \mathbf{1}_N)} \odot \boldsymbol{h}_{t-1} + (\Delta_t \odot \hat{\boldsymbol{x}}_t) \otimes \boldsymbol{B}_t\right) \boldsymbol{C}_t = \sum_{s=1}^{t} \left(\prod_{m=s}^{t-1} e^{\boldsymbol{\Lambda} \odot (\Delta_m \otimes \mathbf{1}_N)} \odot (\hat{\boldsymbol{x}}_s \cdot \boldsymbol{B}_s^\top)\right) \boldsymbol{C}_t \tag{62}$$

$$\xRightarrow{\boldsymbol{\Lambda} \equiv \mathbf{0}, \Delta_m \equiv 1} \quad \boldsymbol{y}_t = \left(\boldsymbol{h}_{t-1} + \hat{\boldsymbol{x}}_t \cdot \boldsymbol{B}_s^\top\right) \boldsymbol{C}_t = \sum_{s=1}^{t} \left(\hat{\boldsymbol{x}}_s \cdot \boldsymbol{B}_s^\top\right) \boldsymbol{C}_t.$$

Notice how the hidden state naturally admits a matrix structure $\boldsymbol{h}_t \in \mathbb{R}^{d \times N}$ due to the outer product $\boldsymbol{x}_s \cdot \boldsymbol{B}_s^\top$, whose columns are updated by $\boldsymbol{x}_t$, where $\boldsymbol{B}_s$ determines which columns are affected. If $\boldsymbol{B}_t \propto \boldsymbol{e}_i$ (i.e., only nonzero at component $i$), then $\boldsymbol{x}_t$ will contribute to only the $i$-th column in $\boldsymbol{h}_t$. Similarly, taking $\boldsymbol{C}_t \propto \boldsymbol{e}_i$ ensures that only the $i$-th column from the hidden state is retrieved. This is precisely what we achieve with the choice outlined above. In light of this, the hidden state after the initial part of the input (where the key-value pairs are listed) admits the following form at time $2\kappa$

$$\boldsymbol{h}_{2\kappa} = \begin{bmatrix} \hat{k}\boldsymbol{e}_1 & \hat{k}\boldsymbol{e}_2 & \dots & \hat{k}\boldsymbol{e}_\kappa \\ + & + & \dots & + \\ \hat{v}\boldsymbol{e}_{j(1)} & \hat{v}\boldsymbol{e}_{j(2)} & \dots & \hat{v}\boldsymbol{e}_{j(\kappa)} \end{bmatrix}, \quad \text{with} \quad \hat{k} = c_{1,k}^2 + c_{0,k}^2, \quad \hat{v} = c_{1,v}c_{0,k}, \tag{63}$$

and remains unchanged until the first key $\boldsymbol{k}_{k_1}$ is encountered at query time, as $\boldsymbol{B}(\boldsymbol{x}) = \boldsymbol{B}((\boldsymbol{v}_i, \boldsymbol{v}_i)) = \mathbf{0}$ for all the tokens in between. For each key $\boldsymbol{k}_i$ encountered at query time, $\boldsymbol{C}(\boldsymbol{x})$ then takes care of extracting the corresponding $i$-th column of

$h$, which holds a vector proportional to the embedding of its associated value, $e_{j(i)}$. Finally, the output matrix computes scalar products of all the values embedding with the extracted vector, which will then be maximum for $e_{j(i)}$, thus accurately solving the MQAR task.

There is an edge case to this construction, which occurs when, at retrieval time, two keys appear adjacent to each other in a key-key pair $(k_i, k_l)$. In this case, $k_l$ represents the actual query, while $k_i$ acts as noise: the convolution will in fact contain information from two keys, implying that $C(x)$ will have two nonzero components at that point. However, with the correct parameter choice, we can have the query information dominate the noise, and still recover the desired solution. We have by (57c)

$$C(\text{conv}(k_i, k_l)) = (c_0 k - b)e_i + (c_1 k - b)e_l. \tag{64}$$

When tested against the hidden state at that instant, then, the output matrix will return a linear combination of the values associated to the two keys:

$$C(\text{conv}(k_i, k_l)) \, h_t = (c_0 k - b)\left(\hat{k}e_i + \hat{v}e_{j(i)}\right) + (c_1 k - b)\left(\hat{k}e_l + \hat{v}e_{j(l)}\right). \tag{65}$$

To ensure the key-value pair of $k_l$ dominates that of $k_i$, we need to impose

$$c_0 k - b \ll c_1 k - b \quad \Longleftarrow \quad c_0 \ll c_1, \tag{66}$$

which is already satisfied by the parameter choice (56).

**Step 2: Dimensionality Reduction to** $d = O(\kappa + \log|V|)$**.** Having proved the construction with $d = O(\kappa + |V|)$, we now apply JL Lemma to reduce the embedding dimension and suitably adjust the Mamba architecture weights to ensure the desired output. Concretely:

- **Embedding** We use the same key embeddings as above $\{e_1, \ldots, e_\kappa\}$ while reducing the value embedding dimensionality to satisfy almost-orthogonality. By JL Lemma (c.f. Lem. 4), given the value embeddings $\{e_j\}_{j=1+\kappa}^{|V|+\kappa}$ and $\epsilon \in (0, 0.5)$, there exists $M_v : \mathbb{R}^{|V|} \to \mathbb{R}^p$ where $p = O(\log|V|)$ and each of its entry is in $\{-\frac{1}{\sqrt{p}}, \frac{1}{\sqrt{p}}\}$ such that $\langle M_v e_{j_1}, M_v e_{j_2} \rangle \leq \epsilon$ for any $j_1 \neq j_2$. Let $M := I_\kappa \oplus M_v \in \mathbb{R}^{(\kappa+p) \times (\kappa+|V|)}$ be the direct sum of the identity matrix preserving the one-hot key embeddings and the JL matrix projecting the value embeddings. Let the normalized value embedding be $\bar{v}_j := M e_j$ Let the *shifted* value embedding be

$$v_j := M e_j + \beta[0_\kappa | 1_p]^\top = \sum_{i=1}^p (M_{i+\kappa, j} + \beta)e_{i+\kappa} \in \mathbb{R}^{\kappa+p}, \tag{67}$$

  where each component of $v_j$ is in the range $[-\frac{1}{\sqrt{p}} + \beta, \frac{1}{\sqrt{p}} + \beta] \equiv [v_{\min}, v_{\max}]$. Note that by letting $\beta > -\frac{1}{\sqrt{p}}$, we can ensure all components of $v_j$ are nonnegative.

- **Convolution** We use a size-2 convolution as (54), to retain the value from $(k_i, v_j)$ pairs and discard the value from $(v_i, k_j)$ pairs. Ideally we want

$$\text{conv}(k_i, v_j) := \text{ReLU}(c_0 k_i + c_1 v_j - b) \propto (c_0 - b)k_i + (c_1 - b)v_j \tag{68a}$$

$$\text{conv}(v_i, k_j) := \text{ReLU}(c_0 v_i + c_1 k_j - b) \perp \text{span}(\{v_1, \ldots, v_{|V|}\}). \tag{68b}$$

  Given the shifted value embeddings, we must impose the following constraints to achieve (68):

$$c_0 k > b, \qquad c_1 v_{\min} > b, \qquad c_0 v_{\max} \leq b. \tag{69}$$

  A feasible parameter combination satisfying (69) is given by

$$\beta = 2\,(\implies [v_{\min}, v_{\max}] \subseteq [1, 3] \text{ since } p \geq 1), \qquad k = 10, \qquad c_0 = 1, \qquad c_1 = 10, \quad \text{and} \quad b = 3. \tag{70}$$

  Consequently, we have the desired convolution outputs

$$\text{conv}(k_i, v_j) = (c_0 k - b)e_i + \sum_{i=1}^p \left(c_1(M_{i+\kappa, j} + \beta) - b\right)e_{i+\kappa} = (c_0 k - b)e_i + c_1 \bar{v}_j + \underbrace{(c_1 \beta - b)\sum_{i=1}^p e_{i+\kappa}}_{:=s} \tag{71a}$$

$$\text{conv}(v_i, k_j) = (c_1 k - b)e_j. \tag{71b}$$

22

Note that the edge case for $\mathrm{conv}(\boldsymbol{k}_i, \boldsymbol{k}_j)$ is taken care of since $c_0 \ll c_1$, while the $(\boldsymbol{v}_i, \boldsymbol{v}_j)$ pairs will be ignored in the SSM layer, following the same argument as in step 1.

- **Mamba SSM** The SSM layer is the same as in step 1 (58), propagating information through the hidden state (c.f. (63)). At the query time for key $\boldsymbol{k}_i, i = 1, \ldots, \kappa$, the output $\boldsymbol{y}_t = \left((c_0 k - b)^2 + (c_1 k - b)^2\right) \boldsymbol{e}_i + (c_0 k - b) \left(c_1 \bar{\boldsymbol{v}}_{j(i)} + \boldsymbol{s}\right)$ contains the (scaled) normalized value embedding shifted by a constant vector $\boldsymbol{s}$.

- **Output** The output layer undoes the shift and classifies based on the normalized value embeddings. Recall $M_v \in \mathbb{R}^{|V| \times p}$ is the JL matrix from the value embedding projection. We set the output linear layer with the weight matrix $W_o = [\boldsymbol{0} | M_v^\top] \in \mathbb{R}^{|V| \times (\kappa + p)}$ and the bias vector $\boldsymbol{b}_o = -(c_0 k - b) W_o \boldsymbol{s}$. The final output is given by

$$W_o \boldsymbol{y}_t + \boldsymbol{b}_o = (c_0 k - b) W_o \left(c_1 \bar{\boldsymbol{v}}_{j(i)} + \boldsymbol{s}\right) - (c_0 k - b) W_o \boldsymbol{s} = (c_0 k - b) c_1 M_v^\top M_v \boldsymbol{e}_{j(i)}, \tag{72}$$

which yields the maximum at component $j(i)$ since $\langle M_v \boldsymbol{e}_{j(i)}, M_v \boldsymbol{e}_l \rangle \le \epsilon$ for any $l \ne j(i)$ by JL Lemma. This completes the proof.

$\square$

**Theorem 3.** *There exists a 1-layer Mamba-2 model without gating that solves* MQAR *with $\kappa$ pairs using embedding size $d = O(\log \kappa + \log |V|)$, and state size $N = \log \kappa$.*

*Proof.* The idea is to execute the Mamba solution outlined in the proof of Thm. 2, but in a *leaner* manner due to the additional degrees of freedom in Mamba-2: particularly, we leverage the fact that the convolutions for the value (14b), key (14d), and query (14e) can be chosen independently (instead of using the same convolution in Mamba). The proof is divided into two steps: We first present the construction using standard basis vectors with $d = \kappa + |V|$; We then reduce the embedding dimension by applying JL lemma.

**Step 1: Construction With $d = O(\kappa + |V|)$.**

- **Embedding - Same as Mamba** The role of the embedding layer is to clearly distinguish values and keys. We achieve this by letting it map to independent directions: let the embedding dimension be $d = |V| + \kappa$, and $\boldsymbol{e}_i \in \mathbb{R}^d$ denote the standard basis vector. We let

$$k_i \mapsto \boldsymbol{k}_i := \boldsymbol{e}_i, \quad \text{and} \quad v_i \mapsto \boldsymbol{v}_i := \boldsymbol{e}_{\kappa+i}. \tag{73}$$

- **Convolution** Differently from Mamba that uses the same convolution kernel to compute the SSM input $\boldsymbol{u}$ and parameters $\boldsymbol{B}, \boldsymbol{C}$, Mamba-2 has the additional flexibility of using three independent convolutions for computing $\boldsymbol{u}, \boldsymbol{B}, \boldsymbol{C}$ (see details in Tab. 2). We now exploit this flexibility by setting:

$$\begin{aligned} \hat{\boldsymbol{x}}_t^B &= \mathrm{conv}_B(\boldsymbol{x}_{t-1}, \boldsymbol{x}_t) = \sigma(\boldsymbol{c}_0 \odot \boldsymbol{x}_{t-1} + \boldsymbol{c}_1 \odot \boldsymbol{x}_t - \boldsymbol{b}) := \boldsymbol{x}_{t-1}, \\ \hat{\boldsymbol{x}}_t &= \mathrm{conv}_u(\boldsymbol{x}_{t-1}, \boldsymbol{x}_t) := \boldsymbol{x}_t, \\ \hat{\boldsymbol{x}}_t^C &= \mathrm{conv}_C(\boldsymbol{x}_{t-1}, \boldsymbol{x}_t) := \boldsymbol{x}_t. \end{aligned} \tag{74}$$

Consequently, the output from $\mathrm{conv}_B$ *shifts* the input sequence to the right by one position, whereas the outputs from $\mathrm{conv}_u, \mathrm{conv}_C$ are the same as the input.

- **Mamba-2 SSM - Same as Mamba** The role of the SSM layer is to associate key-value pairs, propagate information through the state, and retrieve the correct value given a query(=key). Unlike Thm. 2, the convolved input (74) contains only key or value information but never mixes them. A simple choice is to set $\boldsymbol{B}, \boldsymbol{C}$ as the identity matrices, but this requires the state size be the same as the embedding size $N = d$. To further reduce the state size to $N = \kappa$, we use the same choice as Mamba (58) by letting

$$\boldsymbol{\lambda} = \boldsymbol{0}, \boldsymbol{B}(\boldsymbol{x}_t) = \boldsymbol{C}(\boldsymbol{x}_t) = [\boldsymbol{I}_\kappa | \boldsymbol{0}_{|V|}] \boldsymbol{x}_t \equiv W_k \boldsymbol{x}_t, \tag{75}$$

- **Output - Same as Mamba** Even in this case, as a classifier it suffices to pick

$$W_o = \begin{bmatrix} \boldsymbol{v}_1^\top \\ \vdots \\ \boldsymbol{v}_{|V|}^\top \end{bmatrix}. \tag{76}$$

23

With the definitions above, we can simplify the outcome of the Mamba-2 layer application. This in fact reduces to

$$\boldsymbol{y}_t = \sum_{s=1}^{t} \left( \hat{\boldsymbol{x}}_s \, \boldsymbol{B}(\hat{\boldsymbol{x}}_s^B)^\top \right) \boldsymbol{C}(\hat{\boldsymbol{x}}_s^C) = \sum_{s=1}^{t} \left( \boldsymbol{x}_s \left( W_k \boldsymbol{x}_{s-1} \right)^\top \right) W_k \boldsymbol{x}_t = \sum_{s=1}^{t} \boldsymbol{x}_s (\boldsymbol{x}_{s-1}^\top W_k^\top W_k \boldsymbol{x}_t) \tag{77}$$

where $\boldsymbol{x}_0 = \boldsymbol{0}$ is the zero padding vector. Note that the output $\boldsymbol{y}_t$ for $t \geq 2\kappa$ contains the desired key-value association, $\boldsymbol{x}_s(W_k \boldsymbol{x}_{s-1})^\top = \boldsymbol{v}_{j(i)} \boldsymbol{k}_i^\top$ for $s = 2, 4, \ldots, 2\kappa$.

At query time where $\boldsymbol{x}_t = \boldsymbol{k}_i$, thanks to our construction of the Embedding layer, $\langle \boldsymbol{k}_i, \boldsymbol{x}_{s-1} \rangle = 0$ for all $s$, *except* when $\boldsymbol{x}_{s-1} = \boldsymbol{k}_i$: in that case we have instead $\langle \boldsymbol{k}_i, \boldsymbol{x}_{s-1} \rangle = 1$. Whatever the index $s-1$ at which this occurs, the associated value $\boldsymbol{x}_s$ is precisely the one we are seeking, i.e., the value embedding $\boldsymbol{v}_{j(i)}$ immediately following the key matching the corresponding query $\boldsymbol{k}_i$. Thus, $\boldsymbol{y}_t = \boldsymbol{v}_{j(i)}$. Now, applying the output matrix amounts to computing scalar products between $\boldsymbol{v}_{j(i)}$ and all possible value vectors, which will return 1 only at the desired value per the orthogonal embedding construction.

**Step 2: Dimensionality Reduction to** $d = O(\log \kappa + \log |V|)$**.** We now reduce the embedding dimension $d = |V| + \kappa$ to $d = O(\log |V| + \log \kappa)$. To this end, we apply JL Lemma (c.f. Lem. 4) to construct nearly orthogonal embedding vectors, while keeping the convolution, SSM, and output layers the same as step 1. By JL lemma, fixed $\epsilon \in (0, 0.5)$, there exists a matrix $M \in \mathbb{R}^{d \times p}$ for $p = O(\log d)$ such that $|\langle M\boldsymbol{e}_i, M\boldsymbol{e}_j \rangle| \leq \epsilon$ for all $i, j \in [d]$. We apply JL lemma separately for the key embedding subspace (with a JL matrix $M_k \in \mathbb{R}^{\log \kappa \times \kappa}$) and the value embedding subspace (with another JL matrix $M_v \in \mathbb{R}^{\log |V| \times |V|}$). We collect the final JL matrix via a direct sum, $M = M_k \oplus M_v$. Let $\boldsymbol{k}_i := M\boldsymbol{e}_i$, $\boldsymbol{v}_j := M\boldsymbol{e}_{\kappa+j}$. Then for $i \neq j$, we have $\langle \boldsymbol{k}_i, \boldsymbol{k}_j \rangle \leq \epsilon$, $\langle \boldsymbol{v}_i, \boldsymbol{v}_j \rangle \leq \epsilon$. On the other hand, $\langle \boldsymbol{k}_j, \boldsymbol{k}_j \rangle = \sum_{i=1}^{p} M_{i,j}^2 \approx 1$. Similarly, we see that $\langle \boldsymbol{v}_j, \boldsymbol{v}_j \rangle \approx 1$. It remains to show such low-dimensional embeddings, followed by convolutions, SSM, and output layer, yields the desired solution. The action of the convolution layer on the low-dimensional embeddings is the same as in step 1, shifting the attention-keys by one position to the right while keeping the attention-queries and attention-values unchanged. Also the action of the SSM layer is the same as step 1 (77), where the hidden state is a sum of key-value association matrices $\boldsymbol{v}_{j(i)} \boldsymbol{k}_i^\top$, except with embedding dimension $O(\log |V| + \log \kappa)$. At query time, upon encountering $\boldsymbol{k}_i$, the retrieved output is $\boldsymbol{y}_t \propto \boldsymbol{v}_{j(i)} + \epsilon(\sum_{l \neq j(i)} \boldsymbol{v}_l)$ for $\epsilon \ll 1$, in light of the fact that the low-dimensional embeddings are nearly orthogonal by construction. Then, applying the output matrix $W_o \boldsymbol{y}_t$ yields the maximum component at $j(i)$, as desired. $\square$

**Remark 7.** *The above construction with $N = \log \kappa$ works for generic inputs where the values are drawn randomly from the vocabulary with sufficient size. Such construction may fail in the adversarial case where most values are the same and the number of keys $\kappa$ is large. Concretely, suppose the input is $[k_1, v_1, k_2, v, \ldots, k_\kappa, v]$, such that all the values at time $t = 4, 6, \ldots, 2\kappa$ are the same token $v$, and $v_1 \neq v$. Then at the retrieval time for the query token $k_1$, the signal value is $\boldsymbol{v}_1$, whereas the noisy values from other keys contribute to $\epsilon(\kappa - 1)\boldsymbol{v}$. Then if $\epsilon(\kappa - 1) > 1$, the model might fail to retrieve the desired value. This can be counteracted by decreasing $\epsilon$; as per JL Lemma, though, this might come at the cost of scaling the embedding dimension — without however impacting its logarithmic behavior (remember $p = O(\epsilon^{-2} \log d)$).*

**Theorem 4.** *There exists a 1-layer Mamba model with an S4D mixer that solves* MQAR *with $\kappa$ pairs, using embedding size $d = O(\kappa \log |V|)$, and state size $N = 1$.*

*Proof.* The proof is divided into two steps: we first construct a 1-layer Mamba model using S4D as SSM layer, that solves MQAR using $d = O(\kappa|V|)$; we then apply JL Lemma with a shifting trick to complete the proof for $d = O(\kappa \log |V|)$.

**Step 1: Construction With** $d = O(\kappa|V|)$**.** On a high level, the idea is to organize the hidden state of the SSM in chunks, each collecting a vector representing the value associated to a specific key, similarly to the proof in Thm. 2. However, unlike the Mamba layer, the S4D layer does not have access to input-dependent matrices $B(\boldsymbol{x}), C(\boldsymbol{x})$, implying that the $N$ columns of the S4D state are the same up to scaling, and hence cannot encapsulate additional information regarding the input. We then decide to work with a single-column as hidden state, and instead partition it along the embedding dimension $d$. Ideally, we want the hidden state before retrieval to be a long column vector,

$$\boldsymbol{h}_t = \begin{bmatrix} \boldsymbol{v}_1^\top & | & \boldsymbol{v}_2^\top & | & \ldots & | & \boldsymbol{v}_\kappa^\top \end{bmatrix}^\top \in \mathbb{R}^{\kappa|V|}. \tag{78}$$

This can be achieved by specifying each component of the full Mamba architecture as follows.

- **Embedding** The goal of the embedding layer is to organize keys and values in a form that is suitable to assemble a

hidden state as in (78). More in detail, we let a key $k_i$ and a value $v_i$ be mapped to, respectively

$$k_i \mapsto \boldsymbol{k}_i = k \cdot [0, \ldots, 0| \quad \ldots \quad | \underbrace{1, \ldots, 1}_{|V|i:|V|(i+1)} | \quad \ldots \quad |0, \ldots, 0| \ldots]^\top \tag{79}$$

$$v_i \mapsto \boldsymbol{v}_i = v \cdot [0, \ldots, \underbrace{1}_{i}, \ldots, 0| \quad \ldots \quad |0, \ldots, \underbrace{1}_{\kappa|V|+i}, \ldots, 0]^\top \tag{80}$$

$$= v \cdot [\quad \boldsymbol{e}_i^\top \quad | \quad \ldots \quad | \quad \boldsymbol{e}_i^\top \quad ]^\top, \tag{81}$$

for some fixed parameters $k > 0, v > 0$. Notice that, in light of this, any combination $\boldsymbol{k}_i + \boldsymbol{v}_j$ can form a dictionary: it has a maximum value at a component uniquely defined by the pair $i, j$.

- **Convolution** The (short) convolution layer is responsible for filtering out irrelevant information from the sequence, and retaining only the one pertaining $(\boldsymbol{k}_i, \boldsymbol{v}_j)$ pairs. To this end, we limit ourselves to a convolution with kernel size 2: this acts on any pair $(\boldsymbol{x}_t, \boldsymbol{x}_{t+1})$ of the input sequence by mapping it to

$$(\boldsymbol{x}_t, \boldsymbol{x}_{t+1}) \mapsto \sigma\left(\boldsymbol{c}_0 \odot \boldsymbol{x}_t + \boldsymbol{c}_1 \odot \boldsymbol{x}_{t+1} - \boldsymbol{b}\right). \tag{82}$$

By carefully picking the parameters $k, v, \boldsymbol{c}_0, \boldsymbol{c}_1, \boldsymbol{b}$, we can ensure that only a $(\boldsymbol{k}_i, \boldsymbol{v}_j)$ pair "survives" the operation, and everything else gets mapped to the null vector. This for example can be achieved by setting $k = 10, v = 1, \boldsymbol{c}_0 = 10 \cdot \boldsymbol{1}, \boldsymbol{c}_1 = \boldsymbol{1}, \boldsymbol{b} = k\boldsymbol{c}_0$. With this choice, we see that the convolution maps

$$(\boldsymbol{k}_i, \boldsymbol{v}_j) \mapsto [0, \ldots, 0| \quad \ldots \quad |0, \ldots, \underbrace{1}_{|V|i+j}, \ldots, 0| \quad \ldots \quad |0, \ldots, 0]^\top \tag{83}$$

$$= [\quad \boldsymbol{0}^\top \quad | \quad \ldots \quad | \quad \boldsymbol{e}_j^\top \quad | \quad \ldots \quad | \quad \boldsymbol{0}^\top \quad ]^\top \tag{84}$$

$$(\boldsymbol{k}_i, \boldsymbol{k}_j) \mapsto [0, \ldots, 0]^\top \tag{85}$$

$$(\boldsymbol{v}_i, \boldsymbol{k}_j) \mapsto [0, \ldots, 0]^\top \tag{86}$$

$$(\boldsymbol{v}_i, \boldsymbol{v}_j) \mapsto [0, \ldots, 0]^\top, \tag{87}$$

- **S4D SSM** The SSM layer simply needs to accumulate and propagate the combined values down the sequence. We remind that from (5) and (6), the output of S4D is given by

$$\boldsymbol{h}_t^{\text{S4D}} = \sum_{s=1}^{t} \boldsymbol{x}_s \cdot \left(\overline{\boldsymbol{\Lambda}}^{t-(s+1)} \boldsymbol{B}_s\right)^\top. \tag{88}$$

We make use of a "trivial" SSM where $\overline{\boldsymbol{\Lambda}} = \boldsymbol{B} = 1$, resulting in a hidden state which, after collecting the initial $(\boldsymbol{k}_i, \boldsymbol{v}_{j(i)})$ pairs, is constant in the form:

$$\boldsymbol{h}_t = [0, \ldots, \underbrace{1}_{j_1}, \ldots, 0|0, \ldots, \underbrace{1}_{|V|+j_2}, \ldots, 0| \quad \ldots \quad |0, \ldots, \underbrace{1}_{\kappa|V|+j_\kappa}, \ldots, 0]^\top \tag{89}$$

$$= [\quad \boldsymbol{e}_{j_1}^\top \quad | \quad \boldsymbol{e}_{j_2}^\top \quad | \quad \ldots \quad | \quad \boldsymbol{e}_{j_\kappa}^\top \quad ]^\top.$$

- **Gate** The role of the gating mechanism (14g) is to retrieve the part of the hidden state which refer to the requested key. The gate branch acts on a linear transformation of the original sequence: by picking this transformation as the identity. we ensure that, when a key is encountered, only the corresponding value is retrieved from the hidden state, in fact:

$$\tilde{\boldsymbol{y}}_t = \boldsymbol{k}_i \odot \boldsymbol{h}_t = [0, \ldots, 0| \ldots |0, \ldots, \underbrace{1}_{|V|i+j(i)}, \ldots, 0|0, \ldots, 0| \ldots]. \tag{90}$$

- **Output** The final output layer simply needs to test the retrieved vector $\boldsymbol{y}_t$ against all values $\boldsymbol{v}_i$: only the correct one will return a scalar product different from 0. It suffices to pick

$$W_o^{\text{S4D}} = \begin{bmatrix} \boldsymbol{v}_1^\top \\ \vdots \\ \boldsymbol{v}_{|V|}^\top \end{bmatrix} = \begin{bmatrix} \boldsymbol{e}_1^\top & | & \cdots & | & \boldsymbol{e}_1^\top \\ \vdots & | & \vdots & | & \vdots \\ \boldsymbol{e}_{|V|}^\top & | & \cdots & | & \boldsymbol{e}_{|V|}^\top \end{bmatrix}. \tag{91}$$

**Step 2: Dimensionality Reduction to $d = O(\kappa \log(|V| + 1))$.** Having proved the construction with $d = O(|V|)$, we now apply JL Lemma to reduce the embedding dimension and suitably adjust the Mamba architecture weights to ensure the desired output. Concretely:

- **Embedding** We use JL lemma to identify $d = |V| + 1$ almost-orthogonal vectors within a $p \sim O(\log(|V| + 1))$-dimensional space. Namely, there exists a random matrix $M : \mathbb{R}^d \to \mathbb{R}^p$ such that any two (different) vectors in the set $\{Me_1, \ldots, Me_{|V|}, Me_{|V|+1}\}$ are almost-orthogonal, in the sense that $|\langle Me_i, Me_j \rangle| < \epsilon$ for some small $\epsilon \in (0, 0.5)$. Without loss of generality, we ensure that the last of these vectors is parallel to the all-one vector $\mathbf{1}_p$: notice this is always possible via an opportune rotation, which does not affect the scalar product of the recovered vectors (and hence, their almost-orthogonality). We let the embedding layer perform the following map

$$k_i \mapsto \boldsymbol{k}_i = k \cdot [0, \ldots, 0| \quad \ldots \quad |\underbrace{1, \ldots, 1}_{pi:p(i+1)}| \quad \ldots \quad |0, \ldots, 0| \ldots]^\top \in \mathbb{R}^{\kappa p}, \qquad i = 1 \ldots \kappa \tag{92}$$

$$v_i \mapsto \boldsymbol{v}_i = [(Me_i + \beta \mathbf{1})^\top| \quad \ldots \quad |(Me_i + \beta \mathbf{1})^\top]^\top \in \mathbb{R}^{\kappa p}, \qquad i = 1 \ldots |V|, \tag{93}$$

where $M : \mathbb{R}^d \to \mathbb{R}^p$, $M_{i,j} \in \{-\frac{1}{\sqrt{p}}, \frac{1}{\sqrt{p}}\}$ is the (rotated) projection matrix recovered with JL Lemma, and $\beta > \frac{1}{\sqrt{p}}$, so that each component of the value embedding $\boldsymbol{v}_i$ is nonnegative, and falls in the range of $[\beta - \frac{1}{\sqrt{p}}, \beta + \frac{1}{\sqrt{p}}]$.

- **Convolution** We use a size-2 convolution as in step 1, requiring it to retain only the $(\boldsymbol{k}_i, \boldsymbol{v}_j)$ pair information while sending other pairs $(\boldsymbol{v}_i, \boldsymbol{k}_j), (\boldsymbol{v}_i, \boldsymbol{v}_j), (\boldsymbol{k}_i, \boldsymbol{k}_j)$ to zero. To this end, we let

$$\beta = 1, \qquad k = 10, \qquad c_0 = 10, \qquad c_1 = 1, \qquad b = kc_0. \tag{94}$$

Note that by the choice of $\beta$ and $p \geq 1$, the range of components in $\boldsymbol{v}_i$ is limited to $[0, 2]$. Consequently, we have

$$\mathrm{conv}(\boldsymbol{k}_i, \boldsymbol{v}_j) = \mathrm{ReLU}(c_0 \boldsymbol{k}_i + c_1 \boldsymbol{v}_j - b) = [\quad \mathbf{0}^\top \quad | \quad \ldots \quad | \quad (Me_j + \beta \mathbf{1})^\top \quad | \quad \ldots \quad | \quad \mathbf{0}^\top \quad ]^\top; \tag{95}$$

$$\mathrm{conv}(\boldsymbol{v}_i, \boldsymbol{k}_j) = \mathrm{conv}(\boldsymbol{k}_i, \boldsymbol{k}_j) = \mathrm{conv}(\boldsymbol{v}_i, \boldsymbol{v}_j) = 0. \tag{96}$$

- **S4D SSM** The SSM layer proceeds similarly as the construction in step 1, yielding the state at $t > 2\kappa$ as

$$\boldsymbol{h}_t = [(Me_{j_1} + \beta \mathbf{1})^\top|(Me_{j_2} + \beta \mathbf{1})^\top \quad | \quad \ldots \quad |(Me_{j_\kappa} + \beta \mathbf{1})^\top]^\top. \tag{97}$$

- **Gate and Output** Also the gating layer and the output layer proceed similarly as the construction in step 1. After gating, when a key $\boldsymbol{k}_i$ is encountered, we have

$$\hat{\boldsymbol{y}}_t = \boldsymbol{h}_t \odot \boldsymbol{k}_i = \cdot[0, \ldots, 0| \quad \ldots \quad |\underbrace{(Me_{j_i} + \beta \mathbf{1})^\top}_{pi:p(i+1)}| \quad \ldots \quad |0, \ldots, 0| \ldots]^\top. \tag{98}$$

And finally, after applying the output matrix, we obtain

$$W_o^{\mathrm{S4D}} \hat{\boldsymbol{y}}_t = \begin{bmatrix} (Me_1)^\top & \bigg| & \cdots & \bigg| & (Me_1)^\top \\ \vdots & \bigg| & \vdots & \bigg| & \vdots \\ (Me_{|V|})^\top & \bigg| & \cdots & \bigg| & (Me_{|V|})^\top \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ \vdots \\ (Me_{j_i} + \beta \mathbf{1}) \\ \vdots \\ \mathbf{0} \end{bmatrix} \approx \begin{bmatrix} \epsilon + \epsilon\beta \\ \vdots \\ 1 + \epsilon\beta \\ \vdots \\ \epsilon + \epsilon\beta \end{bmatrix}, \tag{99}$$

in light of the fact that both $(Me_i)^\top \cdot Me_j \approx \epsilon$ if $i \neq j$ otherwise $(Me_i)^\top \cdot Me_i \approx 1$ per JL construction, and $(Me_i)^\top \cdot \mathbf{1} \leq \epsilon$ per the assumption that $Me_{|V|+1}$ is parallel to $\mathbf{1}_p$. This allows us to recover the correct value, completing the proof.

$\square$

## C.2. Proofs of Mamba Solving the INDUCTION HEADS Task

**The INDUCTION HEADS Task**  As a reminder, for the INDUCTION HEADS task, the input is a sequence of tokens $[x_1, \ldots, x_t]$ from a finite vocabulary $V$; The output is a sequence of tokens $[y_1, \ldots, y_t]$ from the augmented vocabulary $V \cup \{\times\}$, where $y_i$ equals the input token right after the latest previous occurrence of the input token $x_i$, i.e., $y_i = x_{j(i)+1}$ where $j(i) = \max\{j : j < i, x_j = x_i\}$; otherwise $y_i = \times$. An example input and output drawn from the vocabulary $V = \{1, 2, 3, 4\}$ with sequence length 8 is shown below:

$$
\begin{aligned}
t &= [\quad 1, \quad 2, \quad 3, \quad 4, \quad 5, \quad 6, \quad 7, \quad 8] \\
x &= [\quad 2, \quad 1, \quad 3, \quad 2, \quad 4, \quad 3, \quad 2, \quad 4]. \\
y &= [\quad \times, \quad \times, \quad \times, \quad 1, \quad \times, \quad 2, \quad 4, \quad 3]
\end{aligned}
$$

Note that the input token 2 appears three times, at instants $t = 1, 4, 7$. Thus, at $t = 7$, the *latest* previous occurrence of $j(7) = 4$, which yields the output $y_7 = x_{j(7)+1} = x_{4+1} = 4$.

**Lemma 3.**  *There exists a 1-layer Mamba model with the Mamba-$\Delta^\top$ SSM mixer (15) that solves* INDUCTION HEADS *with vocabulary $V$ using embedding size $d = 2|V|$ and state size $N = |V|$.*

*Proof.*  We follow a procedure similar to the MQAR construction, in that we leverage the matrix structure in the hidden state such that its columns are indexed by the key token and store the associated value token. However, differently from the MQAR task, in the INDUCTION HEADS task there is no distinction between the key and value set, but rather all tokens are drawn from the same vocabulary $V$ – i.e., each token $x_i$ acts as key in the $(x_i, x_{i+1})$ pair, and as value in the $(x_{i-1}, x_i)$ pair. To resolve this, we use a doubling-embedding trick in Mamba, together with suitable choices of convolution. Moreover, the INDUCTION HEADS task requires finding the *latest* previous occurrence; we will achieve this by leveraging the input-dependent state matrix.

The main idea is to *double* the embedding size in Mamba, which enables the convolution layer to perform *concatenations* of the adjacent embedding pairs. Concretely: we let the state size $N = |V|$, and the embedding size $d = 2|V|$. We design the architecture as follows.

- **Embedding** We use $2|V|$-dimensional standard basis vectors to embed the vocabulary $V = \{1, 2, \ldots, |V|\}$, i.e.,

$$
v_i \mapsto \begin{bmatrix} e_{v_i} \\ e_{v_i} \end{bmatrix} \in \mathbb{R}^d \equiv \mathbb{R}^{2|V|}. \tag{100}
$$

- **Convolution** We use size-2 convolution (with left-padding $\mathbf{0}$) combining the pair $(\boldsymbol{x}_{i-1}, \boldsymbol{x}_i)$ by summing the first $|V|$-dimensions of $\boldsymbol{x}_{i-1}$ with the last $|V|$-dimensions of $\boldsymbol{x}_i$, effectively *concatenating* $(e_{x_{i-1}}, e_{x_i})$. We let

$$
\hat{\boldsymbol{x}}_i \equiv \mathrm{conv}(\boldsymbol{x}_{i-1}, \boldsymbol{x}_i) = \boldsymbol{c}_0 \odot \boldsymbol{x}_{i-1} + \boldsymbol{c}_1 \odot \boldsymbol{x}_i, \qquad \text{where } \boldsymbol{c}_0 = \begin{bmatrix} \mathbf{1} \\ \mathbf{0} \end{bmatrix}, \boldsymbol{c}_1 = \begin{bmatrix} \mathbf{0} \\ \mathbf{1} \end{bmatrix}. \tag{101}
$$

  Note that we describe the proof for *linear* convolution here to simplify notation, but it holds also for nonlinear convolution (14b), noting that each embedding vector $\boldsymbol{x}_i$ and the convolution weights $\boldsymbol{c}_0, \boldsymbol{c}_1$ are nonnegative, effectively reducing the nonlinearity $\sigma = \mathrm{ReLU}$ (or SiLU) to be the identity map. The same reasoning applies to the design of $\boldsymbol{B}, \boldsymbol{C}$, as we discuss next.

- **Mamba-$\Delta^\top$ SSM** Since the convolved output $\hat{\boldsymbol{x}}_i$ contains $x_{i-1}$ in its first $|V|$ dimensions and $x_i$ in its last $|V|$ dimensions, we choose the state matrix $\overline{\boldsymbol{\Lambda}}$ and the input matrix $\boldsymbol{B}$ to depend on the first $|V|$ dimensions of $\hat{\boldsymbol{x}}$ (i.e., extracting the key), and the output matrix $\boldsymbol{C}$ to depend on the last $|V|$ dimensions (i.e., extracting the query). To this end, we let $\boldsymbol{\Lambda} = -\mathbf{1} \in \mathbb{R}^{d \times N}, w_\Delta \gg 0$, and

$$
\Delta(\hat{\boldsymbol{x}}_i) := \mathrm{SoftPlus}(\mathrm{Linear}(\hat{\boldsymbol{x}}_i)), \text{ where } \mathrm{Linear}(\hat{\boldsymbol{x}}_i) := w_\Delta [\boldsymbol{I}_{|V|} \mid \mathbf{0}_{|V|}] \hat{\boldsymbol{x}}_i \in \mathbb{R}^N, \tag{102a}
$$

$$
\overline{\boldsymbol{\Lambda}}(\hat{\boldsymbol{x}}_i) := e^{\boldsymbol{\Lambda} \odot (\mathbf{1}_d \otimes \Delta(\hat{\boldsymbol{x}}_i))} = e^{-\mathbf{1} \otimes \Delta(\hat{\boldsymbol{x}}_i)} \in \mathbb{R}^{d \times N} \equiv \mathbb{R}^{|V| \times |V|}, \tag{102b}
$$

$$
\boldsymbol{B}(\hat{\boldsymbol{x}}_i) := \mathrm{Linear}(\hat{\boldsymbol{x}}_i) = [\boldsymbol{I}_{|V|} \mid \mathbf{0}_{|V|}] \hat{\boldsymbol{x}}_i \in \mathbb{R}^N, \tag{102c}
$$

$$
\boldsymbol{C}(\hat{\boldsymbol{x}}_i) := \mathrm{Linear}(\hat{\boldsymbol{x}}_i) = [\mathbf{0}_{|V|} \mid \boldsymbol{I}_{|V|}] \hat{\boldsymbol{x}}_i \in \mathbb{R}^N. \tag{102d}
$$

27

- **Output** $W_o = [\mathbf{0}_{|V|} \mid I_{|V|}] \in \mathbb{R}^{|V| \times d}$

We now show the correctness of such construction. Consider the generic input and output sequences:

$$
\begin{aligned}
t &= [\;\; 1, \quad 2, \quad 3, \quad 4, \quad 5, \quad 6, \quad 7, \quad 8, \; \ldots] \\
x &= [\;\; v_2, \quad v_1, \quad v_3, \quad v_2, \quad v_4, \quad v_3, \quad v_2, \quad v_4, \; \ldots] \\
y &= [\;\; \times, \quad \times, \quad \times, \quad v_1, \quad \times, \quad v_2, \quad v_4, \quad v_3, \; \ldots].
\end{aligned}
\tag{103}
$$

After the embedding and convolution layers, the SSM input is a sequence of $d$-dimensional vectors for $d = 2|V|$,

$$
\hat{\boldsymbol{x}} = \left[
\begin{array}{c|c|c|c|c|c|c|c|c}
\mathbf{0} & \boldsymbol{e}_{v_2} & \boldsymbol{e}_{v_1} & \boldsymbol{e}_{v_3} & \boldsymbol{e}_{v_2} & \boldsymbol{e}_{v_4} & \boldsymbol{e}_{v_3} & \boldsymbol{e}_{v_2} & \cdots \\
\boldsymbol{e}_{v_2} & \boldsymbol{e}_{v_1} & \boldsymbol{e}_{v_3} & \boldsymbol{e}_{v_2} & \boldsymbol{e}_{v_4} & \boldsymbol{e}_{v_3} & \boldsymbol{e}_{v_2} & \boldsymbol{e}_{v_4} & \cdots
\end{array}
\right] \in \mathbb{R}^{2|V| \times t},
\tag{104}
$$

where $\hat{\boldsymbol{x}}_i$ stores the $(x_{i-1}, x_i)$ pair.

The action of the SSM layer organizes the hidden state matrix of size $d \times N \equiv 2|V| \times |V|$ by the input matrix $\boldsymbol{B}(\hat{\boldsymbol{x}}_i)$ taking outer-product with $\hat{\boldsymbol{x}}_i$, followed by retrieving the desired column via the output matrix $\boldsymbol{C}(\hat{\boldsymbol{x}}_i)$. Now by the choice of SSM parameters, we have

$$
\left[\begin{array}{c} \boldsymbol{B}(\hat{\boldsymbol{x}}) \\ \hline \boldsymbol{C}(\hat{\boldsymbol{x}}) \end{array}\right] = \left[
\begin{array}{c|c|c|c|c|c|c|c|c}
\mathbf{0} & \boldsymbol{e}_{v_2} & \boldsymbol{e}_{v_1} & \boldsymbol{e}_{v_3} & \boldsymbol{e}_{v_2} & \boldsymbol{e}_{v_4} & \boldsymbol{e}_{v_3} & \boldsymbol{e}_{v_2} & \cdots \\
\hline
\boldsymbol{e}_{v_2} & \boldsymbol{e}_{v_1} & \boldsymbol{e}_{v_3} & \boldsymbol{e}_{v_2} & \boldsymbol{e}_{v_4} & \boldsymbol{e}_{v_3} & \boldsymbol{e}_{v_2} & \boldsymbol{e}_{v_4} & \cdots
\end{array}
\right] \in \mathbb{R}^{2|V| \times t},
\tag{105}
$$

where we stack them together to visualize that $\boldsymbol{B}(\hat{\boldsymbol{x}})$ amounts to shifting $\boldsymbol{C}(\hat{\boldsymbol{x}})$ to the right by one position, due to the design of $\mathrm{conv}_B$.

We now verify the desired behavior in the SSM layer. Suppose temporarily the state matrix is $\overline{\boldsymbol{\Lambda}} = \mathbf{1} \in \mathbb{R}^{d \times N}$. Then the hidden state at time $s \leq t$ would be a cumulative sum,

$$
\boldsymbol{h}_s = \sum_{i=1}^{s} \hat{\boldsymbol{x}}_i \boldsymbol{B}(\hat{\boldsymbol{x}}_i)^\top = \sum_{i=2}^{s} \begin{bmatrix} \boldsymbol{e}_{x_{i-1}} \\ \boldsymbol{e}_{x_i} \end{bmatrix} \boldsymbol{e}_{x_{i-1}}^\top.
\tag{106}
$$

Thus, the $j$-th column of the hidden state matrix would store the sum of all $\hat{\boldsymbol{x}}_i$ where the key $\boldsymbol{B}(\hat{\boldsymbol{x}}_i) = x_{i-1} = \boldsymbol{e}_j$. Yet the INDUCTION HEADS task requires storing the *latest* associated value only (not *all* associated values). To this end, we leverage the input-dependence of the state matrix, and particularly of $\Delta(x_t)$. Recall the Mamba-$\Delta^\top$ layer is given by

$$
\boldsymbol{h}_s = e^{\boldsymbol{\Lambda} \odot (\mathbf{1} \otimes \Delta(\hat{\boldsymbol{x}}_s))} \odot \boldsymbol{h}_{s-1} + \hat{\boldsymbol{x}}_s \boldsymbol{B}(\hat{\boldsymbol{x}}_s)^\top,
\tag{107}
$$

$$
\boldsymbol{y}_s = \boldsymbol{h}_s \boldsymbol{C}(\hat{\boldsymbol{x}}_s).
\tag{108}
$$

We design $\Delta(\hat{\boldsymbol{x}}_t) \in \mathbb{R}^N \equiv \mathbb{R}^{|V|}$ such that when the input contains the key information, the corresponding key column in the state is *erased* (while the other columns remain the same). Without loss of generality, suppose $\boldsymbol{B}(\hat{\boldsymbol{x}}_s) = \boldsymbol{e}_j$. By the definition of $\boldsymbol{B}$ (102c), this implies that $[\boldsymbol{I}_{|V|} \mid \mathbf{0}_{|V|}]\hat{\boldsymbol{x}}_s = \boldsymbol{e}_j$. By the choice of $w_\Delta \gg 0$ and the definition of $\Delta$ (102a), we have

$$
\Delta(\hat{\boldsymbol{x}}_i) = \mathrm{SoftPlus}(w_\Delta [\boldsymbol{I}_d \mid \mathbf{0}_d]\hat{\boldsymbol{x}}_i) = w_\Delta \boldsymbol{e}_j.
\tag{109}
$$

Therefore (102b) yields the state matrix as

$$
\overline{\boldsymbol{\Lambda}}_s = e^{-\mathbf{1} \otimes \Delta(\hat{\boldsymbol{x}}_s)} = e^{-\mathbf{1} \otimes (w_\Delta \boldsymbol{e}_j)} = \mathbf{1} \otimes \exp[0, \ldots, \underbrace{-w_\Delta}_{j}, \ldots, 0]^\top \overset{w_\Delta \to \infty}{=} \mathbf{1} \otimes (\mathbf{1} - \boldsymbol{e}_j) \in \mathbb{R}^{d \times N},
\tag{110}
$$

which has an all-zeros $j$-th column and all-ones columns elsewhere. Consequently, the $j$-th column of the hidden state $\boldsymbol{h}_s[:, j]$ is erased by the action $e^{-\mathbf{1} \otimes \Delta(\hat{\boldsymbol{x}}_s)} \odot \boldsymbol{h}_{s-1}$, and then updated with the current input containing the latest value by the action $\hat{\boldsymbol{x}}_s \otimes \boldsymbol{B}(\hat{\boldsymbol{x}}_s)$, as desired. We remark that such erasure operation is akin to the construction of the S6 layer for solving the KEEP $n$-TH task in Cor. 2, in which we have Mamba approximate a Heaviside by tweaking $\Delta(x_t)$, so to erase information from all tokens before a given one. We also see that such *selective* erasure works consistently well for long sequences when $t \to \infty$, since it preserves all other columns (except the $j$-th one) by setting $\Delta(\boldsymbol{u}_s)[l] = 0$ for $l \neq j$, and thereby satisfying the condition in Lem. 2 to avoid sensitivity decay.

Finally, the SSM output is given by $\boldsymbol{y}_s = \boldsymbol{h}_s \boldsymbol{C}_s = \boldsymbol{h}_s \boldsymbol{e}_{x_s}$, which retrieves the $x_s$-th column of the state that stores the token immediately after the *latest* previous occurrence of $x_s$, i.e. $\boldsymbol{y}_s = \begin{bmatrix} \boldsymbol{e}_{x_{j(s)}} \\ \boldsymbol{e}_{x_{j(s)+1}} \end{bmatrix}$. We then apply the output matrix $W_o \boldsymbol{y}_s = \boldsymbol{e}_{x_{j(s)+1}}$ to obtain the target value, which completes the proof.

$\square$

# D. Additional Experiment Details

## D.1. Training Details

For all experiments, unless otherwise noted, we train with the Adam optimizer (Kingma & Ba, 2014), for 600 epochs using an initial learning rate $\eta = 0.03$ and cosine annealing down to $\eta = 1 \times 10^{-6}$. The training set is composed of $10^5$ randomly generated samples with a fixed seed and the batch size is 16, which results in up to $3.75 \cdot 10^6$ gradient updates. We perform early stopping if the validation loss reaches below $10^{-6}$ or if six hours have elapsed since the beginning of training. The validation set and test sets have $10^3$ and $10^5$ samples respectively and are generated with the same function but using different seeds. Reported accuracies are always obtained from the test set after the last epoch of training.

## D.2. Task KEEP $n$-TH

In this section, we provide additional details and ablation of the KEEP $n$-TH task used for Sec. 4, with experimental set-up and partial results reported in Tab. 1.

**Model** All the results in Tab. A.1 - Tab. A.4 use 1-layer models. The MAMBA and S4D models (with or without PE) are simplified architectures, which consists of embedding layer, SSM layer, and output linear layer, *without convolution nor gating* from the Mamba mixer block. The simplification is intended to investigate the role of the SSM layer alone (i.e. S6 versus S4D), without confounding factors from other components in the mixer block.

**Experiment Set-up** For each input sequence $\boldsymbol{x} = (x_1, \ldots, x_T)$, we draw $x_i \overset{i.i.d}{\sim} \text{Unif}(\{1, \ldots, |V|\})$ randomly from a vocabulary with size $|V| = 128$. The target output is the $n$-th token in the input sequence, i.e. $y_t = x_n$ for $n < t \leq T$. The predicted output at time $t = n, \ldots, T$ are taken to compute (cross-entropy) loss for training, and accuracy for evaluation.

**Discussion** When equipped with PE, Mamba manages to achieve 100% accuracy on KEEP $n$-TH, regardless of sequence length. This is thanks to its ability to dynamically adjust (via $\Delta(x_t)$) for how long the hidden state retains memory of the target token (in position $n = 5$). On the other hand, S4D is lacking such ability, and already fails at the task for $T = 20$ (Tab. A.1 and A.2). Then again Transformers do not need to retain memory, as they can look back to the whole sequence at each step, in light of their attention mechanism, and have no issue solving the task for any $T$. When removing PE from Mamba, however (Tab. A.3 and A.4), the model loses its way to discriminate the specific token that must be retrieved, and performance drops to that of S4D, as expected.

Table A.1: Ablation on KEEP $n$-TH: MAMBA+PE, S4D+PE, TRANSFORMERS with varying sequence length $T = 10, 20$, embedding dimension $d$, and state size $N$.

| | | T=10 | | | T=20 | | |
|---|---|---|---|---|---|---|---|
| | | acc. | # epch. | # prm. | acc. | # epch. | # prm. |
| MAMBA (+PE) | d=8, N=8 | 1.00 (0.00) | 107 | 2.2k | 1.00 (0.00) | 37 | 2.2k |
| | d=8, N=32 | 1.00 (0.00) | 241 | 2.8k | 1.00 (0.00) | 39 | 2.8k |
| | d=8, N=64 | 1.00 (0.00) | 161 | 3.5k | 1.00 (0.00) | 47 | 3.5k |
| | d=32, N=8 | 1.00 (0.00) | 56 | 9.2k | 1.00 (0.00) | 15 | 9.2k |
| | d=32, N=32 | 1.00 (0.00) | 19 | 11.5k | 1.00 (0.00) | 14 | 11.5k |
| | d=32, N=64 | 1.00 (0.00) | 88 | 14.5k | 1.00 (0.00) | 19 | 14.5k |
| | d=64, N=8 | 0.99 (0.00) | 553 | 18.7k | 0.99 (0.01) | 208 | 18.7k |
| | d=64, N=32 | 1.00 (0.00) | 429 | 23.3k | 1.00 (0.00) | 12 | 23.3k |
| | d=64, N=64 | 1.00 (0.00) | 345 | 29.4k | 1.00 (0.00) | 12 | 29.4k |
| S4D (+PE) | d=8, N=8 | 0.99 (0.00) | 600 | 2.1k | 0.43 (0.04) | 600 | 2.1k |
| | d=8, N=32 | 0.99 (0.00) | 600 | 2.5k | 0.46 (0.05) | 600 | 2.5k |
| | d=8, N=64 | 0.99 (0.00) | 600 | 3.0k | 0.43 (0.03) | 600 | 3.0k |
| | d=32, N=8 | 0.94 (0.00) | 600 | 8.7k | 0.72 (0.01) | 600 | 8.7k |
| | d=32, N=32 | 0.93 (0.01) | 600 | 10.3k | 0.72 (0.01) | 600 | 10.3k |
| | d=32, N=64 | 0.93 (0.01) | 600 | 12.4k | 0.73 (0.01) | 600 | 12.4k |
| | d=64, N=8 | 0.79 (0.00) | 600 | 17.5k | 0.14 (0.01) | 600 | 17.5k |
| | d=64, N=32 | 0.79 (0.00) | 600 | 20.6k | 0.14 (0.00) | 600 | 20.6k |
| | d=64, N=64 | 0.79 (0.00) | 600 | 24.8k | 0.14 (0.00) | 600 | 24.8k |
| TRANS-FORMER | l=1, d=16 | 1.00 (0.00) | 16 | 5.8k | 1.00 (0.00) | 19 | 6.0k |
| | l=1, d=32 | 1.00 (0.00) | 12 | 15.1k | 1.00 (0.00) | 12 | 15.4k |
| | l=1, d=64 | 1.00 (0.00) | 10 | 42.3k | 1.00 (0.00) | 9 | 42.9k |

Table A.2: Ablation on KEEP $n$-TH: MAMBA+PE, S4D+PE, TRANSFORMERS with varying sequence length $T = 30, 40, 50$, embedding dimension $d$, and state size $N$.

| | | T=30 | | | T=40 | | | T=50 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | acc. | # epch. | # prm. | acc. | # epch. | # prm. | acc. | # epch. | # prm. |
| MAMBA (+PE) | d=8, N=8 | 0.94 (0.04) | 456 | 2.2k | 0.85 (0.15) | 325 | 2.2k | 0.22 (0.15) | 600 | 2.2k |
| | d=8, N=32 | 1.00 (0.00) | 285 | 2.8k | 0.58 (0.37) | 416 | 2.8k | 0.28 (0.10) | 600 | 2.8k |
| | d=8, N=64 | 1.00 (0.00) | 180 | 3.5k | 0.59 (0.31) | 420 | 3.5k | 0.98 (0.02) | 346 | 3.5k |
| | d=32, N=8 | 1.00 (0.00) | 20 | 9.2k | 1.00 (0.00) | 50 | 9.2k | 1.00 (0.00) | 241 | 9.2k |
| | d=32, N=32 | 1.00 (0.00) | 20 | 11.5k | 1.00 (0.00) | 21 | 11.5k | 1.00 (0.00) | 225 | 11.5k |
| | d=32, N=64 | 1.00 (0.00) | 19 | 14.5k | 1.00 (0.00) | 23 | 14.5k | 1.00 (0.0) | 130 | 14.5k |
| | d=64, N=8 | 1.00 (0.00) | 17 | 18.7k | 1.00 (0.00) | 164 | 18.7k | 0.99 (0.00) | 600 | 18.7k |
| | d=64, N=32 | 1.00 (0.00) | 24 | 23.3k | 1.00 (0.00) | 66 | 23.3k | 0.99 (0.00) | 600 | 23.3k |
| | d=64, N=64 | 1.00 (0.00) | 29 | 29.4k | 1.00 (0.00) | 226 | 29.4k | 0.98 (0.01) | 600 | 29.4k |
| S4D (+PE) | d=8, N=8 | 0.14 (0.04) | 600 | 2.1k | 0.03 (0.00) | 600 | 2.1k | 0.03 (0.00) | 600 | 2.1k |
| | d=8, N=32 | 0.46 (0.05) | 600 | 2.5k | 0.15 (0.01) | 600 | 2.5k | 0.04 (0.01) | 600 | 2.5k |
| | d=8, N=64 | 0.43 (0.03) | 600 | 3.0k | 0.13 (0.02) | 600 | 3.0k | 0.04 (0.00) | 600 | 3.0k |
| | d=32, N=8 | 0.72 (0.01) | 600 | 8.7k | 0.09 (0.00) | 600 | 8.7k | 0.08 (0.00) | 600 | 8.7k |
| | d=32, N=32 | 0.72 (0.01) | 600 | 10.3k | 0.09 (0.00) | 600 | 10.3k | 0.09 (0.00) | 600 | 10.3k |
| | d=32, N=64 | 0.73 (0.01) | 600 | 12.4k | 0.09 (0.00) | 600 | 12.4k | 0.09 (0.00) | 600 | 12.4k |
| | d=64, N=8 | 0.14 (0.01) | 600 | 17.5k | 0.09 (0.00) | 600 | 17.5k | 0.09 (0.00) | 600 | 17.5k |
| | d=64, N=32 | 0.14 (0.00) | 600 | 20.6k | 0.10 (0.00) | 600 | 20.6k | 0.10 (0.00) | 600 | 20.6k |
| | d=64, N=64 | 0.14 (0.00) | 600 | 24.8k | 0.09 (0.00) | 600 | 24.8k | 0.09 (0.00) | 600 | 24.8k |
| TRANS-FORMER | l=1, d=16 | 1.00 (0.00) | 16 | 5.8k | 1.00 (0.00) | 23 | 6.3k | 1.00 (0.00) | 17 | 6.4k |
| | l=1, d=32 | 1.00 (0.00) | 12 | 15.1k | 1.00 (0.00) | 12 | 16.0k | 1.00 (0.00) | 13 | 16.4k |
| | l=1, d=64 | 1.00 (0.00) | 9 | 42.3k | 1.00 (0.00) | 9 | 44.2k | 1.00 (0.00) | 9 | 44.9k |

Table A.3: Ablation on KEEP $n$-TH: MAMBA and S4D with varying sequence length $T = 10, 20$, embedding dimension $d$, and state size $N$.

| | | T=10 | | | T=20 | | |
|---|---|---|---|---|---|---|---|
| | | acc. | # epch. | # prm. | acc. | # epch. | # prm. |
| MAMBA | d=8, N=8 | 0.20 (0.01) | 600 | 2.3k | 0.05 (0.00) | 600 | 2.3k |
| | d=8, N=32 | 0.18 (0.03) | 600 | 2.9k | 0.05 (0.00) | 600 | 2.9k |
| | d=8, N=64 | 0.17 (0.02) | 600 | 3.6k | 0.05 (0.00) | 600 | 3.6k |
| | d=32, N=8 | 0.21 (0.01) | 600 | 9.3k | 0.10 (0.00) | 600 | 9.3k |
| | d=32, N=32 | 0.25 (0.03) | 600 | 11.6k | 0.11 (0.00) | 600 | 11.6k |
| | d=32, N=64 | 0.28 (0.03) | 600 | 14.7k | 0.11 (0.00) | 600 | 14.7k |
| | d=64, N=8 | 0.19 (0.00) | 600 | 18.8k | 0.11 (0.00) | 600 | 18.8k |
| | d=64, N=32 | 0.20 (0.00) | 600 | 23.4k | 0.11 (0.00) | 600 | 23.4k |
| | d=64, N=64 | 0.21 (0.00) | 600 | 29.6k | 0.11 (0.00) | 600 | 29.6k |
| S4D | d=8, N=8 | 0.08 (0.00) | 600 | 2.2k | 0.04 (0.00) | 600 | 2.2k |
| | d=8, N=32 | 0.08 (0.00) | 600 | 2.6k | 0.04 (0.00) | 600 | 2.6k |
| | d=8, N=64 | 0.08 (0.00) | 600 | 3.2k | 0.04 (0.00) | 600 | 3.2k |
| | d=32, N=8 | 0.20 (0.00) | 600 | 8.8k | 0.10 (0.00) | 600 | 8.8k |
| | d=32, N=32 | 0.20 (0.00) | 600 | 10.4k | 0.10 (0.00) | 600 | 10.4k |
| | d=32, N=64 | 0.20 (0.00) | 600 | 12.5k | 0.10 (0.00) | 600 | 12.5k |
| | d=64, N=8 | 0.20 (0.00) | 600 | 17.7k | 0.11 (0.00) | 600 | 17.7k |
| | d=64, N=32 | 0.20 (0.00) | 600 | 20.8k | 0.11 (0.00) | 600 | 20.8k |
| | d=64, N=64 | 0.20 (0.00) | 600 | 24.9k | 0.11 (0.00) | 600 | 24.9k |

Table A.4: Ablation on KEEP $n$-TH: MAMBA and S4D with varying sequence length $T = 30, 40, 50$, embedding dimension $d$, and state size $N$.

| | | T=30 | | | T=40 | | | T=50 | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | acc. | # epch. | # prm. | acc. | # epch. | # prm. | acc. | # epch. | # prm. |
| | d=8, N=8 | 0.04 (0.00) | 600 | 2.3k | 0.04 (0.00) | 600 | 2.3k | 0.03 (0.00) | 600 | 2.3k |
| | d=8, N=32 | 0.04 (0.00) | 600 | 2.9k | 0.04 (0.00) | 600 | 2.9k | 0.03 (0.00) | 600 | 2.9k |
| | d=8, N=64 | 0.04 (0.00) | 600 | 3.6k | 0.04 (0.00) | 600 | 3.6k | 0.03 (0.00) | 600 | 3.6k |
| | d=32, N=8 | 0.09 (0.00) | 600 | 9.3k | 0.09 (0.00) | 600 | 9.3k | 0.08 (0.00) | 600 | 9.3k |
| MAMBA | d=32, N=32 | 0.09 (0.00) | 600 | 11.6k | 0.09 (0.00) | 600 | 11.6k | 0.08 (0.00) | 600 | 11.6k |
| | d=32, N=64 | 0.09 (0.00) | 600 | 14.7k | 0.09 (0.00) | 600 | 14.7k | 0.08 (0.00) | 600 | 14.7k |
| | d=64, N=8 | 0.09 (0.00) | 600 | 18.8k | 0.09 (0.00) | 600 | 18.8k | 0.09 (0.00) | 600 | 18.8k |
| | d=64, N=32 | 0.09 (0.00) | 600 | 23.4k | 0.09 (0.00) | 600 | 23.4k | 0.09 (0.00) | 600 | 23.4k |
| | d=64, N=64 | 0.09 (0.00) | 600 | 29.6k | 0.09 (0.00) | 600 | 29.6k | 0.09 (0.00) | 600 | 29.6k |
| | d=8, N=8 | 0.03 (0.00) | 600 | 2.2k | 0.03 (0.00) | 600 | 2.2k | 0.03 (0.00) | 600 | 2.2k |
| | d=8, N=32 | 0.03 (0.00) | 600 | 2.6k | 0.03 (0.00) | 600 | 2.6k | 0.03 (0.00) | 600 | 2.6k |
| | d=8, N=64 | 0.03 (0.00) | 600 | 3.2k | 0.03 (0.00) | 600 | 3.2k | 0.03 (0.00) | 600 | 3.2k |
| | d=32, N=8 | 0.08 (0.00) | 600 | 8.8k | 0.09 (0.00) | 600 | 8.8k | 0.09 (0.00) | 600 | 8.8k |
| S4D | d=32, N=32 | 0.08 (0.00) | 600 | 10.4k | 0.09 (0.00) | 600 | 10.4k | 0.09 (0.00) | 600 | 10.4k |
| | d=32, N=64 | 0.08 (0.00) | 600 | 12.5k | 0.09 (0.00) | 600 | 12.5k | 0.09 (0.00) | 600 | 12.5k |
| | d=64, N=8 | 0.09 (0.00) | 600 | 17.7k | 0.09 (0.00) | 600 | 17.7k | 0.09 (0.00) | 600 | 17.7k |
| | d=64, N=32 | 0.09 (0.00) | 600 | 20.8k | 0.09 (0.00) | 600 | 20.8k | 0.09 (0.00) | 600 | 20.8k |
| | d=64, N=64 | 0.09 (0.00) | 600 | 24.9k | 0.09 (0.00) | 600 | 24.9k | 0.09 (0.00) | 600 | 24.9k |

## D.3. Task MQAR

**Models** All the results in Fig. A.2 use 1-layer models. The Mamba and Mamba-2 models explicitly disable the gating branch. This simplification is intended to verify our constructions without gating in Thm. 2 and Thm. 3. The Mamba-S4D model retains the full Mamba architecture, but only swapping the S6 layer with the S4D layer; not to be confused with the original S4D model proposed in (Gu et al., 2022a).

**Experiment Set-up** We generate the data described in App. C.1 as follows. For each input sequence of the form

$$\boldsymbol{x} = [k_1, v_1, \ldots, k_\kappa, v_\kappa, \ldots, \mid k_{i_1}, \ldots, k_{i_2}],$$

we draw the key token $k_i \overset{i.i.d}{\sim} \text{Unif}(\{1, \ldots, \kappa\})$, and value token $v_j \overset{i.i.d}{\sim} \text{Unif}(\{1, \ldots, |V|\})$. The target output sequence consists of masked tokens except at the the query chunk where the input query is a key (e.g., at $k_{i_1}, k_{i_2}$ in the example above). We compute loss during training (and accuracy for evaluation) only at the query positions, informed from the target output sequence.

**Discussion** Here we expand on the results from Fig. 2 in the main text, by reporting the accuracy of Mamba, Mamba-2, and Mamba-S4D trained on MQAR for varying model sizes. In Fig. A.2, we sweep over values of the value vocabulary size $|V|$, to show how our theoretical bounds hold while varying this parameter. The bounds are still reasonably tight, and the observations drawn from Fig. 2 still hold in this case. The extra caveat (which does not however invalidate our claim) is that by increasing $|V|$ we are making the task more difficult to solve, and our training procedure for the simplest Mamba-S4D fails to achieve satisfactory performance for the model sizes considered. Notice also that, for Mamba-2, the simplest task $\kappa = 4$ can achieve 100% accuracy even *below* the theoretical curve proposed in our theorems. This can be attributed to the following factors. On the one hand, our bounds rely on JL Lemma, which provides only *asymptotical* behaviors which might not be verified in practice for $\kappa$ so small. On the other hand, it is perfectly feasible that at this regime the architecture can recover a more efficient solution than the one theorized.
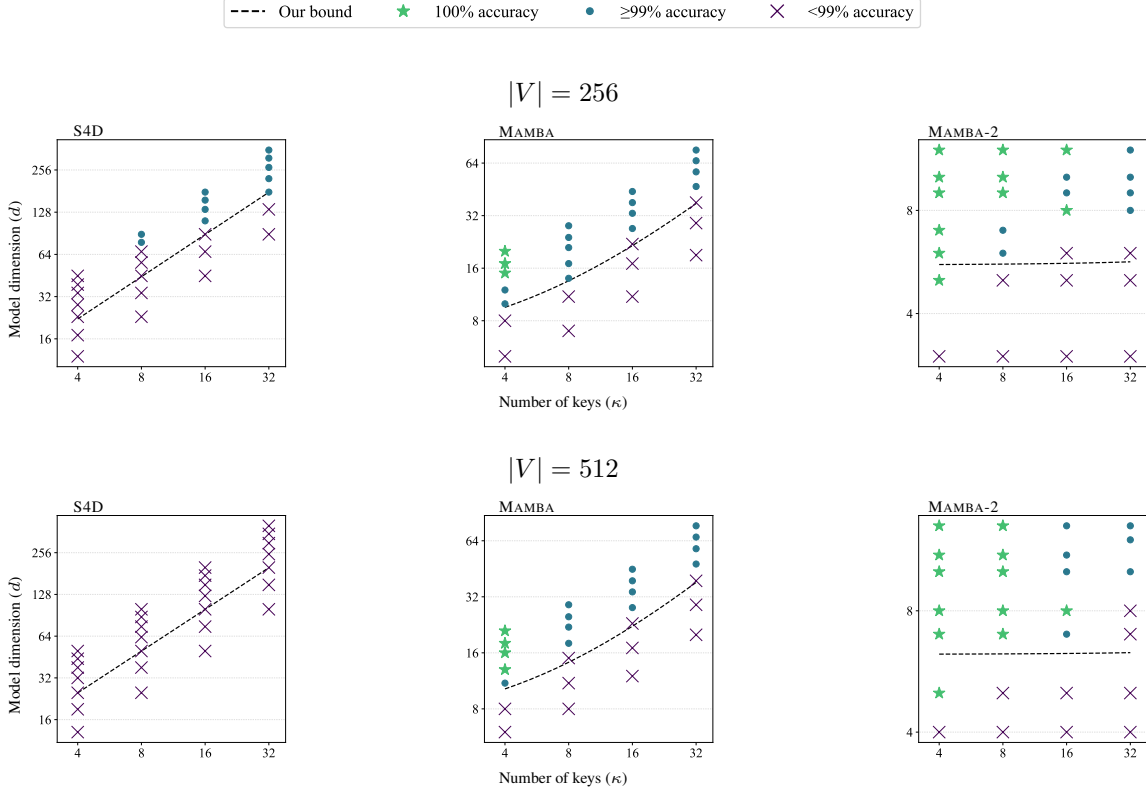


*Figure A.2.* Trained models accuracy on MQAR (best of 7 seeds) across $\kappa$, and $d$. For S4D $N = 4$, for Mamba $N = 2 \times \kappa$, and for Mamba2 $N = 8 \times \ln \kappa$. $T = 100$ and $|V| \in \{256, 512\}$ for all runs.

With Fig. A.3, we further complement our results by sweeping over values of the $N$ state size parameter. We remind that, according to Thm. 2 to 4, our theorized MQAR solutions require a value of at least $N = 1$, $N = \kappa$ and $N \sim \log \kappa$ for S4D, Mamba and Mamba-2, respectively. Indeed, in Fig. A.3 we observe that varying $N$ does not have a particular impact on the final accuracy of S4D. For Mamba, on the other hand, we see that for $N < \kappa$ the training procedure fails to recover an exact solution to MQAR. Similarly, for Mamba-2, no solution is recovered for $N < 4 \log \kappa$. These results further validate the tightness of our theoretical solutions.
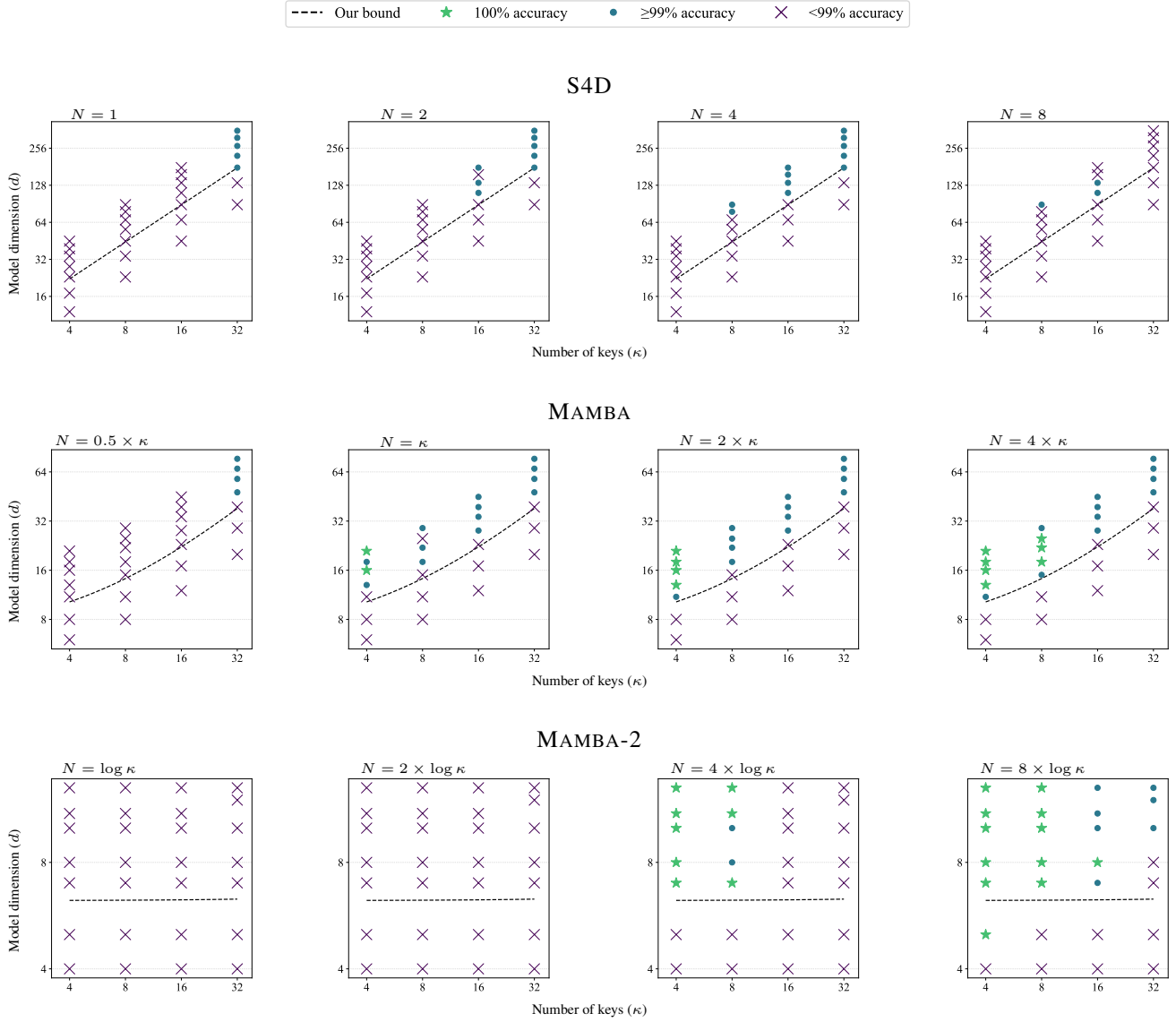
*Figure A.3.* Trained models accuracy on MQAR (best of 7 seeds) across $N$, $\kappa$, and $d$. For Mamba and Mamba-2 $|V| = 512$, while for S4D $|V| = 256$ (as it failed to reach satisfactory accuracy for larger $|V|$). $T = 100$ for all runs.

## D.4. Task INDUCTION HEADS

**Model** All the results in Fig. A.2 use 1-layer Mamba that explicitly disables the gating branch. This simplification is intended to verify our constructions without gating in Lem. 3.

**Experiment Set-up** We generate the data described in App. C.2 as follows. The data generation requires four scalar parameters: vocabulary size $|V|$, sequence length $T > 2|V| + 1$, hard case probability $p \in [0, 1]$, and special range ratio $\gamma \in (0, 0.1]$. We first draw $X \sim \text{Bernoulli}(p)$: if $X = 0$, we sample from the standard setting, otherwise the hard setting. The standard setting generates the input sequence $\boldsymbol{x} = (x_1, \ldots, x_T)$ by randomly drawing $x_i \overset{i.i.d.}{\sim} V = \{1, \ldots, |V|\}$ for $i = 1, \ldots, T$. The hard setting is intended to evaluate the long-range memorization capability (i.e., placing repeated tokens at the beginning and the end of the sequence), which consists of the following steps.

1. Randomly pick a special token $v^* \in V$

2. Generate the input sequence by randomly drawing $\boldsymbol{x} = (x_1, \ldots, x_T)$ where $x_i \overset{i.i.d.}{\sim} V \setminus v^*$ for $i = 1, \ldots, T$

3. Randomly draw a position from $r \in \text{Unif}(\{1, \ldots, \gamma T\})$

4. Place the special token $v^*$ at positions $r, T - r$.

In the experiments for Fig. 3, we use $T = 100, |V| \in \{5, 10, 20, 40\}, p = 0.75, \gamma = 0.1$. In Fig. A.4, we further ablate the choice of state size $N \in \{|V|, 2|V|, 4|V|\}$.

**Discussion** We design the hard setting to better differentiate the capabilities from Mamba and our proposed Mamba-$\Delta^\top$. Specifically, solving the standard setting of the INDUCTION HEADS task requires memorizing the latest previous occurrence, or forgetting the earlier previous occurrences. Note that the input sequence generated from the standard setting consists of many repeated tokens (by the requirement $T > 2|V| + 1$), and the expected time for reappearance of any token is $|V|$. Thus, for small and medium-size $|V|$, Mamba can solve for these cases by using the state matrix with negative eigenvalues to discount the remote past pairs, and thereby correctly output the latest previous occurrence. However, solving the hard setting additionally requires the model to memorize long-range information due to the special token (occurring at the beginning part and the end part of the sequence). We see that the Mamba solution with negative eigenvalues is *at odds with* the long-range memorization, as shown in Lem. 1. On the other hand, Mamba-$\Delta^\top$ can satisfy both selective forgetting and long-range memorization via the input-dependence state matrix that erases outdated information specific to the input key, while retaining other information in the hidden state, as illustrated in the proof of Lem. 3 (see details in App. C.2).

Below we expand on the results in Sec. 5.2 by reporting a sweep on the hidden state size $N$ for the models used in the INDUCTION HEADS task experiments, complementing the findings shown in Fig. 3.
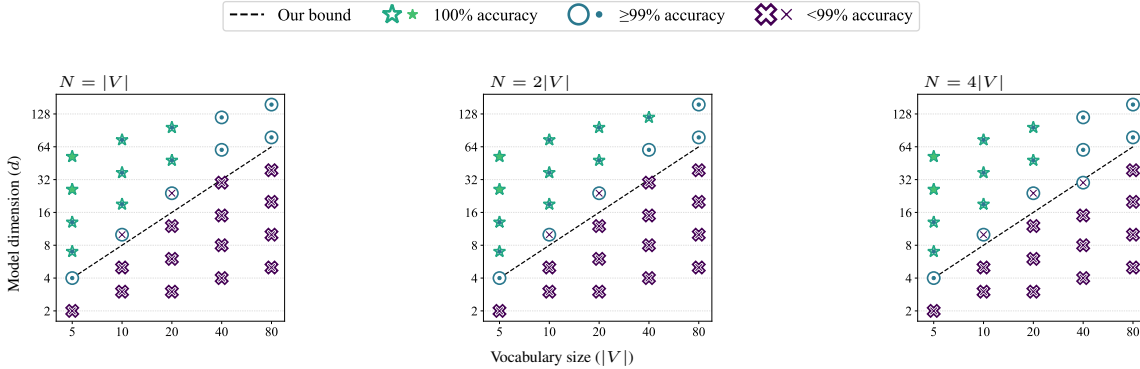


*Figure A.4.* Trained models accuracy on INDUCTION HEADS task (best of 5 seeds), varying $|V|$ and $d$, with $N \in \{|V|, 2|V|, 4|V|\}$ (left, middle, right). Mamba-$\Delta^\top$'s performance (outlined) is equal or better than Mamba's (filled) and only hits $100\%$ above the theoretical bound from Lem. 3 (black).