

# QUANTIFYING MEMORY UTILIZATION WITH EFFECTIVE STATE-SIZE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

As the space of causal sequence modeling architectures continues to grow, the need to develop a general framework for their analysis becomes increasingly important. With this aim, we draw insights from classical signal processing and control theory, to develop a quantitative measure of *memory utilization*: the internal mechanisms through which a model stores past information to produce future outputs. This metric, which we call *effective state-size* (ESS), is tailored to the fundamental class of systems with *input-invariant* and *input-varying linear operators*, encompassing a variety of computational units such as variants of attention, convolutions, and recurrences. Unlike prior work on memory utilization, which either relies on raw operator visualizations (e.g. attention maps), or simply the total *memory capacity* (i.e. cache size) of a model, our metrics provide highly interpretable and actionable measurements. In particular, we show how ESS can be leveraged to improve initialization strategies, inform novel regularizers and advance the performance-efficiency frontier through model distillation. Furthermore, we demonstrate that the effect of context delimiters (such as end-of-speech tokens) on ESS highlights cross-architectural differences in how large language models utilize their available memory to recall information. Overall, we find that ESS provides valuable insights into the dynamics that dictate memory utilization, enabling the design of more efficient and effective sequence models.

## 1 INTRODUCTION

In recent years, the success of autoregressive sequence modeling in the context of deep learning has largely been driven by advancements in highly parallelizable causal architectures, such as the transformer (Vaswani et al., 2023). However, despite their strong performance and hardware efficiency, understanding the inner workings of these neural networks remains a challenging task due to their non-linearity and the diversity of fundamental building blocks used. To this end, we leverage a new class of model abstractions, allowing for the development of a unified framework for the analysis of these computational units.

In particular, we note that the majority of sequence models of practical interest can formally be expressed as either linear systems ( $y = Tu$ ) or systems with *input-varying linear operators* ( $y = f(u)u$ ), the latter of which we abbreviate as LIV<sup>1</sup>. LIVs generalize the notion of adaptive, or *data-controlled* operators to a broader class than previously described in Massaroli et al. (2021); Poli et al. (2023). The input-varying linear operator framework decouples the input-varying *featurization*  $u \mapsto T := f(u)$  and the linear mapping  $y = Tu$  required to construct and apply the operator respectively.

This decomposition enables a wide array of deep learning primitives to be uniformly formulated as linear systems, including models like convolutions [38; 45; 34; 54], linear<sup>2</sup> recurrences [15; 18; 27; 59; 29; 26; 73; 47; 16], and attention variants [66; 33; 63].

<sup>1</sup>We abbreviate them as LIVs instead of IVLs to conform to the classical convention, where linear time-varying and time-invariant systems are abbreviated as LTVs and LTIs respectively, as well as recent works such as Zancato et al. (2024), which call them “linear input-varying systems”.

<sup>2</sup>Here, by *linear* we refer to the linearity of the state transition.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

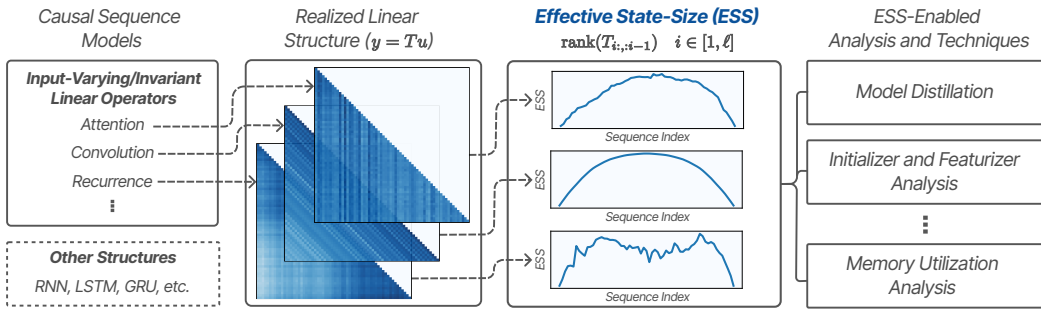


Figure 1: An overview of the effective state-size metric and its various downstream applications.

Current approaches to analyzing the inner workings of LIVs often rely upon simple visualizations of the materialized operator  $T$  (or the aggregation of  $T$  across multiple layers and residuals) [43; 67; 1; 4; 72; 61]. However, these visualizations alone often fail to highlight critical properties that explain how different models construct internal representations of the input data. Moreover, prior attempts in obtaining quantitative metrics, such as through spectral analysis of the operator  $T$  (Min & Li, 2024; Bhojanapalli et al., 2020), are either limited to a specific model class or do not appropriately take into account important conflating factors like the causal masking of  $T$  which significantly distorts the metric (Wu et al., 2024).

In this work, we focus our analysis on the working memory<sup>3</sup> of model architectures and examine two aspects of model memory in particular: *memory capacity* (i.e. cache/state size) and *memory utilization*. Notably, memory capacity alone can be misleading, as models with similar capacities may learn to utilize their available memory to varying degrees. Therefore, we introduce the notion of memory utilization – a measure that provides deeper insight into the differences between architectures with comparable computational efficiency.

Our main technical contributions can be summarized as follows:

- We draw from classical signal processing and control theory, and propose *effective state-size* (ESS) – a metric computed by taking the rank of  $T_{i:,i-1}$  – as a proxy for *memory utilization* in LIVs (Section 3).
- We validate ESS beyond its theoretical interpretation by demonstrating its correlation with performance across a wide range of models and memory-intensive synthetic tasks, including associative recall and selective copying (Section 4).
- We construct initializers motivated by ESS which result in performance improvements relative to default initializers (Section 5.1).
- We explore the use of the ESS metric as a means of enhancing the performance-efficiency trade-off by demonstrating its ability to inform model distillation (Section 5.2).
- We extend the utility of ESS to language, demonstrating how it captures a previously uncharacterized property of LLMs: state modulation (Section 5.3). This phenomenon provides concrete intuition as to why introducing input dependencies into state transitions is crucial for in-context recall.

## 2 RELATED WORK

In this section, we briefly describe previous works that are closely related to the core findings of the paper. For an extended discussion of the related works, refer to Section B.

**Classifying sequence models.** A sequence model can be defined as a mapping from input sequences to output sequences that leverages a state to maintain information across time. Within the scope of deep learning, the early attempts at constructing sequence models fall under the class of non-linear RNNs (Rumelhart & McClelland, 1987; Hochreiter & Schmidhuber, 1997) which exhibit non-linear state space dynamics. However, in recent years, models with linear state transitions

<sup>3</sup>Here, we refer to “memory” in the sense commonly associated with the “state” of dynamical systems, as described by Willems (1989), as opposed to the notion of language models memorizing some fact encountered during training (Allen-Zhu & Li, 2024).

have grown in popularity due to inherent trainability (Pascanu et al., 2013) and efficiency (Martin & Cundy, 2018) limitations in the non-linear regime. These models fall under two categories: linear systems and systems with input-varying linear operators (LIVs). We note that linear systems can further be broken down into time-invariant (LTI) and time-varying (LTV) systems, depending on whether or not the parameters of the realized recurrence change as a function of time (i.e. sequence index). This taxonomy has also been used in prior work (Zancato et al., 2024) to motivate the construction of novel architectures. In modern deep learning, LTI systems (Lindquist & Picci, 1979) can be found in frameworks like S4 (Gu et al., 2022a). However, given the lack of expressivity afforded to LTIs, the current state-of-the-art models fall under the more expressive LIV class (Gu & Dao, 2024; Poli et al., 2023; Vaswani et al., 2023) which parameterizes the state-space dynamics as a function of the input. Since LIVs constitute the minimal superclass containing most sequence models of practical interest, the goal of this work is to develop a quantitative means of analyzing models under the LIV abstraction. To do so, we draw from minimal realization theory (DeWilde & van der Veen, 1998; Akaike, 1974), which has previously been applied within the scope of classical linear systems but has yet to be extended to modern LIV sequence models.

**Quantitatively distinguishing sequence models.** As discussed in Section 1, the current landscape of interpreting sequence models primarily reduces to qualitative measures or limited quantitative ones. Here, we expand on the latter, noting that commonly used measures such as model size (Kaplan et al., 2020) and cache/state size (which we define formally as theoretically realizable state-size in Section 3) are imperfect means of distinguishing sequence models. Namely, while these metrics serve as reasonable proxies for the capacity of a model to learn, they fail to capture how much of that capacity is realized. As such, they ignore important aspects of the model pipeline such as data, initialization, and optimization. To capture these aspects of network training, we apply the notions of numerical and effective rank (Roy & Vetterli, 2007) to the minimal realization problem, which gives rise to the ESS metric. We elaborate on this in Section 3.

### 3 THEORY

In this section, we begin by showing that most modern sequence models can effectively materialize a linear operator  $T$ . We then formally define ESS as a metric derived from  $T$ , providing theoretical insights grounded in minimal realization theory. Finally, we detail the practical computation of ESS, proposing two variants: tolerance-ESS and entropy-ESS.

Using the flattened notation, we let  $T \in \mathbb{R}^{d\ell \times d\ell}$ ,  $u, y \in \mathbb{R}^{d\ell}$  denote the operator, inputs, and outputs respectively,  $\ell$  denote the sequence length and  $d$  denote the channel dimension. Here, we index sequence indices with subscripts, i.e.  $T_{ij} \in \mathbb{R}^{d \times d}$ ,  $u_i \in \mathbb{R}^d$  and channels with superscripts, i.e.  $T^{\alpha\beta} \in \mathbb{R}^{\ell \times \ell}$ ,  $u^\alpha \in \mathbb{R}^\ell$ . For additional details on notation, refer to Section C.1.

**A unified representation of sequence models.** While typically nonlinear, most sequence models of interest can effectively materialize a linear operator  $T$ , where the equation  $y = Tu$  faithfully expresses the computation performed by the model (see Section D.2 for further elaboration):

$$\begin{aligned} T_{ij} &= C_i B_j && \text{linear attention,} && T_{ij} &= C_i A_{i-1} \cdots A_{j+1} B_j && \text{recurrence,} \\ T_{ij} &= K_{i-j} && \text{convolution,} && T_{ij} &= C_i K_{i-j} B_j && \text{gated convolution,} \\ T_{ij} &= \sigma(C_i B_j) && \text{attention.} \end{aligned}$$

We make a distinction between linear systems (such as convolutions) and LIVs (such as attention and gated convolutions). In the former, the operator  $T$  is *input-invariant* whereas the latter are constructed via *causal featurizers* that map past inputs into features, i.e.  $f_B : u_{:i} \mapsto B_i$ , which are then used to construct the elements of  $T$  as outlined above. We begin our formulation by restricting the scope to linear systems.

**Operator-recurrence duality and the minimal recurrent realization of linear systems.** Consider a general linear recurrence formulated as follows:

$$\begin{aligned} s_{i+1} &= A_i s_i + B_i u_i \\ y_i &= C_i s_i + D_i u_i, \end{aligned} \tag{1}$$

where  $(A_i \in \mathbb{R}^{n_{i+1} \times n_i}, B_i \in \mathbb{R}^{n_{i+1} \times d}, C_i \in \mathbb{R}^{d \times n_i}, D_i \in \mathbb{R}^{d \times d})_{i \in [\ell]}$ ;  $s_i$  and  $n_i$  are the state and state-size at sequence index  $i$  respectively. Classic results (DeWilde & van der Veen, 1998, ch. 3) establish that for any given input-invariant operator  $T$  (i.e. linear system), there exist infinite recurrent realizations in the form of Equation (1), motivating the search for the minimal one.

**Theorem 3.1.** *Given any causal input-invariant operator  $T$ , there exist infinite variations of linear recurrences in the form of Equation (1) that realize an equivalent input-output operator.*

A simple extension of the proof of this theorem (in Section C.2.4) demonstrates the following:

**Theorem 3.2.** *The rank of the operator submatrix  $(H_i \equiv T_{i:,i-1})$  determines the minimal state size required to represent the causal operation  $(y = Tu)$  as a recurrence.*

*Proof.* The proof of Theorem 3.1 shows that the operator submatrices  $H_i$  can be decomposed arbitrarily into two state-projection matrices,  $\mathcal{O}_i$  and  $\mathcal{C}_i$ , whose inner product dimension defines the state size of its recurrent realization at sequence index  $i$ . By the rank-nullity theorem,  $\text{rank}(H_i)$  represents the minimum inner product dimension of any such state-projection matrices and thus corresponds to the minimally realizable state size of the operator  $T$  at sequence index  $i$ .  $\square$

We formally refer to the metric  $n_i^* = \text{rank}(H_i)$  as **effective state-size** (ESS)<sup>4</sup> and discuss its interpretation for both linear systems and LIVs below.

**Interpreting effective state-size.** As shown in Theorem 3.2, the ESS of an input-invariant linear system is given by its minimal state-size which is directly interpretable as a measure of model memory utilization. For LIVs, however, the minimal realization process outlined in the proof of Theorem 3.1 is no longer guaranteed to obtain recurrences that preserve causality (i.e. the minimally realized features  $A_i^*$  depend on future inputs  $u_k$ ,  $k > i$ )<sup>5</sup>. Nevertheless, ESS lower bounds the state-size  $n_i$  (refer to Section C.2.2), meaning that for any LIV, an equivalent recurrence must necessarily materialize a state-size at least as large as its ESS. Therefore, we claim that ESS serves as a proxy for memory utilization in not only linear systems, but in LIVs as well. We empirically validate the usefulness of this interpretation of ESS for LIVs in Sections 4 and 5.

**Memory capacity in LIVs.** The memory capacity of LIVs is given by the state-size  $n_i$ . We formally refer to it as **theoretically realizable state-size** (TSS), as it serves as a tight upper bound for ESS. For models without realization-agnostic minimal recurrent formulations, such as softmax attention, we resort to the trivial realization as shown in Equation C.2.5, in which TSS (for a single channel) is equal to the sequence index  $i$ . For models like SSMs and linear attention variants, TSS is equal to the state-size defined by their recurrent formulation. We refer readers to Section D.2 for detailed operator-specific derivations of TSS.

Finally, we note that in LIVs, ESS depends not only on the model’s functional form, but also on the input data, optimization, and more generally anything that impacts the realization of  $T$ ; in contrast, TSS is limited in that it is determined solely by the model’s structure. This enables ESS to measure meaningful differences across architectures (Section 4), including ones that possess the same or similar TSS. Additionally, within a single model, ESS can be used to compare memory utilization on a per-token basis as shown in Section 5.3.

### 3.1 COMPUTING EFFECTIVE STATE-SIZE

Computing ESS requires a few additional considerations due to the numerical errors and approximations involved in practice. We propose two approaches – both of which rely on singular values  $(\Sigma_i)$  from taking the singular value decomposition (SVD) of  $H_i$  – that provide complementary perspectives on the same metric.

<sup>4</sup>Note that semiseparable rank as discussed in prior works (Vandebril et al., 2005; Xia et al., 2010; Dao & Gu, 2024) is more closely related to the theoretically realizable state-size ( $n$ ). We also observe that it can potentially be interpreted as  $\max(\{\text{ESS}_i\}_{i \in [\ell]})$ .

<sup>5</sup>An example of a causality preserving realization of LIVs is the trivial realization shown in Equation (C.2.5).

**Tolerance-ESS.** Here, a tolerance value is manually selected to threshold the singular values of  $H_i$ , determining the ESS metric as follows:

$$\text{tolerance-ESS}_i := |\{\sigma_i^m : \sigma_i^m > \tau, \sigma_i^m \in \Sigma_i\}|. \quad (2)$$

According to the Eckart–Young–Mirsky theorem, the tolerance-ESS metric can be interpreted as the minimum state size necessary for an input-invariant recurrence to approximate the original operator, such that the spectral norm of the approximation error remains below the specified tolerance level ( $\|T_{ij} - T_{ij}^*\|_2 \leq \tau$ ).

**Entropy-ESS.** One drawback of tolerance-ESS is its reliance on the somewhat arbitrary selection of a tolerance value. One can instead compute the effective rank (Roy & Vetterli, 2007), which involves exponentiating the normalized spectral entropy (perplexity) of  $H_i$ :

$$\text{entropy-ESS}_i := \exp\left(-\sum_m p_i^m \log(p_i^m)\right), \quad \text{where } p_i^m = \frac{\sigma_i^m}{\|\sigma_i\|_1}. \quad (3)$$

In contrast to the tolerance-based metric which is discrete, entropy-ESS can assume continuous values ranging from 1 to  $|\Sigma_i|$  and does not require the selection of a tolerance value. However, the normalization applied to the singular values results in the loss of absolute values, which may be significant for per-sequence-index comparisons of state size. Nonetheless, both the tolerance-based and entropy-based forms of ESS are valuable for model analysis. Entropy-ESS is particularly useful for summarizing metrics across the entire tolerance space, whereas tolerance-ESS offers a more precise and readily interpretable depiction of rank concerning approximation error. Unless a tolerance is specified, we use entropy-ESS throughout all our experiments. Our code for computing ESS can be found in Section D.1.1.

**Average vs total ESS.** Note that ESS is computed separately for each sequence index of the input on a per-channel, per-layer basis. Unless otherwise stated, we average across the layer, channel, sequence, and batch size dimensions to produce a single-number summary of the ESS of a given model. We will refer to this measure specifically as the average ESS in analyses that also include a different metric, the total ESS. Total ESS is computed by first summing the ESS across channels and then averaging over the remaining dimensions. Total ESS is particularly useful in the case of a non-recurrent model like softmax attention, whose average TSS (computed analogous to average ESS) depends only on the sequence index  $i$ . This means that in softmax attention, the average TSS remains constant irrespective of model width (i.e. channel dimension); total ESS (and analogously total TSS), in contrast, does change as a function of model width. Since we perform analyses that compare ESS to TSS, examining total ESS can serve as a more informative metric than average ESS in certain cases. For more details, refer to Sections D.1 and D.2.

**Computational complexity of ESS.** One final note is that since SVD scales cubically in the size of a square matrix, the time complexity of computing ESS on a model with  $M$  layers,  $D$  channels, sequence length  $L$ , and batch size  $B$  is given by  $O(L^3BDM)$ . Hence, it is a costly metric to compute for long sequence lengths. Fortunately, in recurrent models with bounded TSS  $n \ll L$ , we can use a truncated SVD to reduce the time complexity to  $O(L^2nBDM)$ . Furthermore, in models like linear attention which have channels that share the same recurrence within the same head, the time complexity is further reduced by a factor of the head dimension.

## 4 EMPIRICAL VALIDATION OF EFFECTIVE STATE-SIZE

To demonstrate the practical utility of ESS beyond its theoretical interpretation discussed in Section 3, we next turn to an empirical analysis. In this section, we examine ESS across a wide range of tasks and models in order to understand how it varies across different regimes, with particular focus placed on its relationship with model performance on memory-intensive tasks.

**Task space.** To explore ESS in an extensive, yet controlled, manner, we iterate on a set of synthetic tasks proposed by Poli et al. (2024) which have been shown to effectively approximate model performance on large-scale language tasks. Specifically, we train models on the multi-query associative recall (MQAR), selective copying, and compression tasks, each of which probes the ability

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

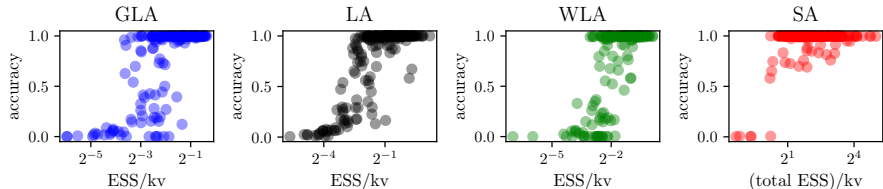


Figure 2: Scatter plots of accuracy vs ESS/kv across featurizers. Within each featurizer plot, all task-model configurations from the sweep corresponding to each featurizer are shown.

of models to effectively utilize their working memory. We note that here, we restrict the presentation of our results to MQAR and refer the reader to Section E.1 for the results on selective copying and compression, which showcase analogous trends.

**Model space.** We explore four models within the scope of this analysis: gated linear attention (GLA), weighted linear attention (WLA), linear attention (LA) and softmax attention (SA). We choose this set of frameworks since, together, they capture a large portion of the space of modern sequence models. The key distinctions between these models are as follows (more details can be found in Section D.2):

- GLA layer: This layer implements the gated linear attention formulation described in Yang et al. (2024a), where the recurrent feature  $A$  (gating term) is input-varying, placing it in the same class as models like Liquid-S4 (Hasani et al., 2022) and Mamba (Gu & Dao, 2024; Dao & Gu, 2024).
- WLA layer: This layer is nearly identical to GLA, but with an input-invariant  $A$  matrix. This lies in the same class as Hyena-S4D (Poli et al., 2023), RetNets (Sun et al., 2023), and gated-convolutions in general.
- LA layer: This layer is based on Katharopoulos et al. (2020);  $A$  is not trainable and is instead fixed as the identity matrix.
- SA layer: This is the canonical attention layer which is similar to linear attention, but with the addition of a softmax non-linearity applied to the attention matrix (Vaswani et al., 2023), enabling unbounded TSS.

**Experimental setup.** In our analysis, we exhaustively sweep across the tasks and models (which are comprised of two sequence-mixing and two channel-mixing layers) detailed above.

Within each task, we also sweep across varying task difficulties. In the case of MQAR, we do so by modulating the number of key-value (kv) pairs the models are tasked to match, as well as the total sequence length of the prompt. Within each model, we sweep across varying TSS. For each task-model configuration, we compute the ESS and accuracy on a validation set every 10 epochs. We will refer to the entire space of tasks and models across which we sweep as the task-model space. Finally, we split our profiling of ESS into two sections: cross task-model analysis (Section 4.1) and within task-model analysis (Section 4.2). For more details on the setup, refer to Section D.3.

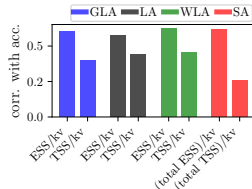


Figure 3: ESS/kv vs TSS/kv as a proxy for model performance as measured by correlation.

#### 4.1 CROSS TASK-MODEL ANALYSIS

Our first goal is to understand how ESS empirically captures memory utilization by studying its correlation with post-training MQAR performance across the entire task-model space. To appropriately analyze ESS across tasks, we normalize it by the memory demands of MQAR, constructing an adjusted form of ESS given by ESS/kv.

**Finding 1:** Measured over the entire task-model space, ESS/kv exhibits a significantly higher correlation with accuracy than TSS/kv (Figures 2, 3, 9a, 9b).

Note that the strong correlation between ESS/kv and accuracy highlights the efficacy of ESS as a proxy for memory utilization. Furthermore, this finding underscores a significant gap in the explanatory power between ESS and TSS, emphasizing the importance of analyzing models beyond just their memory capacity.

#### 4.2 WITHIN TASK-MODEL ANALYSIS

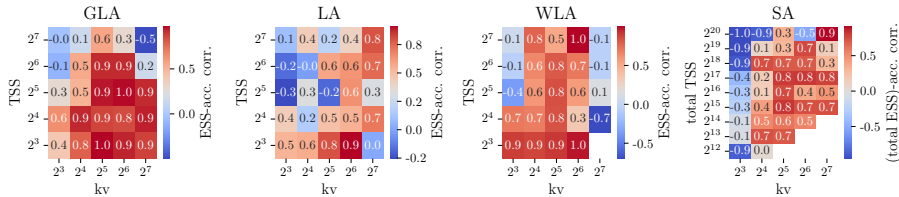


Figure 4: Correlation between ESS and accuracy over the course of model training bucketed by TSS and kv.

Next, to further establish ESS as a proxy for memory utilization, we study how ESS evolves as a function of MQAR performance in a regime where TSS is kept fixed and, therefore, does not correlate with accuracy. We do this by analyzing ESS-accuracy correlation on a per-model, per-task basis over the course of training, uncovering several insights that serve as the basis for our subsequent analysis.

**Finding 2:** For less memory-intensive tasks trained using models with high TSS, we observe a lower correlation between ESS and performance compared to more memory-intensive tasks trained using a lower TSS (Figure 4).

This is in line with the interpretation of ESS as a measure of memory utilization. For easier tasks that are learned by a model with high memory capacity, the model is not incentivized to increase its memory utilization beyond where it resides at initialization. In contrast, for difficult tasks that operate in a memory-constrained regime, the model is forced to increase its memory utilization in order to learn, resulting in strong positive correlations between accuracy and ESS over training.<sup>6</sup> Digging a bit deeper, we find that this form of ESS analysis reveals two failure modes of model learning in recurrent frameworks (which recall have bounded TSS): **state saturation** and **state collapse**.

State saturation refers to the scenario in which a model has insufficient TSS to fully learn a task, resulting in its ESS converging near its TSS. This is reflected in its ESS/TSS (which we refer to as state utilization) residing near 1. We observe this in Figure 5 where we note that models with a TSS of 8 perform worse as the task difficulty scales due to a saturated state. State collapse, on the other hand, refers to the scenario in which a model has sufficient TSS to learn (or partially learn) a task, but its ESS fails to increase during training, resulting in a heavily underutilized state. With respect to state collapse, we observe the following:

**Finding 3:** For GLA and WLA, state collapse occurs in the high kv bucket of task-model space (i.e. kv = 2<sup>7</sup>) whereas for LA it does not (Figure 5). More generally, we find that LA has higher state utilization than GLA and WLA.

While state saturation can only be solved by increasing TSS, state collapse can in principle be solved by increasing ESS. Unlike TSS, which is a fixed hyperparameter of the model, one can modulate ESS by changing various aspects of the model pipeline. Furthermore, even outside of the state collapse regime, given the positive correlation between ESS and performance across the task-model space, increasing ESS is a generally viable approach to improving model performance without sacrificing efficiency. We explore this idea in the results to follow.

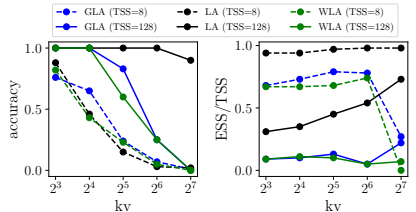


Figure 5: Accuracy and state utilization as a function of kv for low and high TSS models.

<sup>6</sup>In Figure 4, the empty spot in the WLA grid corresponds to a NaN from the entropy-ESS computation. The empty spots in the SA grid correspond to MQAR task constraints discussed in Section D.3.

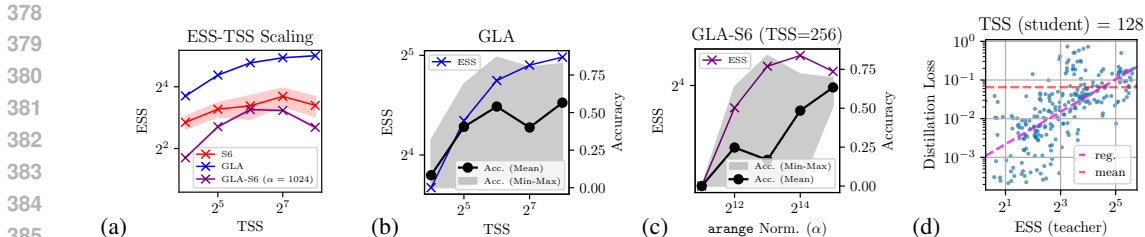


Figure 6: (a) ESS-TSS scaling in the S6, GLA and GLA-S6 featurizers. (b) ESS and accuracy on MQAR as a function of TSS in GLA. (c) ESS and accuracy on MQAR as a function of normalization factor for initialization in GLA-S6. (d) Distillation loss vs ESS of the teacher model.

## 5 APPLICATIONS OF EFFECTIVE STATE-SIZE

In Section 4, we showed that changes in ESS are correlated with changes in performance, both across models and during model training, indicating its importance beyond just interpretability. In this section, we aim to push this insight further by understanding how we can leverage ESS to both improve upon and improve our understanding of the existing performance-efficiency frontier in sequence models. We partition our results based on the stage of model training at which we apply ESS analysis: initialization-phase (Section 5.1), post-training (Section 5.2), and mid-training (Section E.3). Finally, we move beyond synthetics and extend ESS to language in Section 5.3.

### 5.1 INITIALIZATION-PHASE ANALYSIS

Initialization in weight space plays a crucial role in machine learning, significantly impacting model convergence and training stability (Glorot & Bengio, 2010). We extend this concept to the initialization of recurrent models in state space, leaning on the intuition from Figure 2 that suggests higher ESS can enhance performance. Namely, we illustrate how ESS at initialization can be used to inform featurizer selection – the selection of the function that maps the input to the operator  $T = f(u)$  or equivalently the recurrent features  $(A_i(u_i), B_i(u_i), C_i(u_i), D_i(u_i))_{i \in [\ell]}$  – and initialization schemes. In doing so, we uncover design flaws of a prominent model, S6 (Mamba) (Gu & Dao, 2024).

**ESS-informed featurizer selection.** To study the relationship between memory capacity and memory utilization in S6, we remove the short convolutional layer in the Mamba block and stack two of these modified blocks between SwiGLUs (Shazeer, 2020). Under the default MQAR task settings outlined in Poli et al. (2024) (see Tables 2, 3, and 4 for details), we observe that S6 is entirely unable to learn MQAR (accuracy  $\approx 0$ ) across multiple scales of TSS (16 - 256) as shown in Figure 25. This aligns with results from Yang et al. (2024b), which independently demonstrate the poor performance of S6 without the additional short convolutional layer on a different in-context recall task. To investigate the cause, we look into how S6 is preconditioned to utilize its memory by computing its ESS when processing a Gaussian noise input, prior to training.

**Finding 4:** Figure 6a demonstrates that the ESS of S6 layers at initialization scales poorly with respect to TSS, notably failing to increase monotonically. In contrast, GLA layers (Yang et al., 2024a), configured with hyperparameters to match the TSS, model width, number of layers, and hidden-state normalization of the S6 model (see Section D.2 and Table 2), exhibit greater and monotonically increasing ESS-TSS scaling at initialization (Figure 6a). Despite the architectural similarities between the S6 and GLA layers, Figure 6b demonstrates that, unlike S6, GLA achieves accuracy improvements that correlate with increases in both TSS and ESS.

Based on these findings, we conjecture that the poor ESS-TSS scaling of S6 prevents the model from effectively utilizing all of its states, irrespective of increases in memory capacity.

**ESS-informed initialization scheme.** To further investigate the differences between the aforementioned S6 model and GLA model, we construct a composite model termed GLA-S6. This model adopts the feature-sharing structure of GLA (dividing dimensions into heads and sharing



432 computations within a head), but applies the S6 featurization to the  $A$  matrix as follows:

433  
434 GLA (original):  $A = \text{diag}(\text{sigmoid}(Wu)^{1/\beta})$  (4)

435 GLA-S6:  $A = \text{diag}(\exp(-([1/\alpha \ 2/\alpha \ \dots \ n/\alpha]^T \odot \text{softplus}(Wu))))$ . (5)

436  
437 Like S6, GLA-S6 fails to learn MQAR across the same range of TSS (see Figure 25) and exhibits  
438 poor initialization-ESS scaling as shown in Figure 6a. Upon further inspection, we identify the  
439 cause of poor ESS scaling: with each new state introduced, the `arange` term ( $[1 \ 2 \ \dots \ n]$ )  
440 exponentially pushes new entries of  $A$  towards zero, negating the effects of additional states despite  
441 the increase in TSS. Therefore, to ameliorate the poor ESS scaling, we propose a simple solution:  
442 increase the normalization factor.

443 **Finding 5:** By scaling the normalization factor ( $\alpha$ ), Figure 6c shows that GLA-S6 achieves im-  
444 provements in MQAR accuracy post-training, reflecting the impact of increasing its initialization-  
445 ESS, despite the models having identical memory capacities.

446 These experiments demonstrate the efficacy of analyzing ESS at initialization, as they reveal how  
447 different models are preconditioned to utilize their working memory. This analysis helps identify  
448 potentially weak featurization and initialization schemes, enabling us to pinpoint shortcomings in  
449 the S6 featurizer and implement a straightforward fix.

## 450 5.2 POST-TRAINING ANALYSIS

451 Recall from Section 4.1 that we observed a strong correlation between ESS and post-training per-  
452 formance. Building on this insight, a natural question arises: can ESS be used for more than just  
453 performance analysis in the post-training setting? In this section, we answer this question by explor-  
454 ing an additional post-training application: model-order reduction.

455 **ESS-informed model-order reduction.** Model-order reduction refers to the process of improv-  
456 ing model efficiency by reducing state-size while retaining performance. Previous works, such  
457 as Massaroli et al. (2023), have explored the distillation of linear time-invariant (LTI) operators  
458 ( $T_{ij} = T_{i+k,j+k}$ ) into linear recurrences with small state-sizes using backpropagation. Other tech-  
459 niques for model-order reduction such as modal truncation and balanced truncation (Beliczynski  
460 et al., 1992; Gawronski & Juang, 1990) are also applicable to LTIs. In this study, however, we  
461 are concerned with improving the efficiency of general LIVs. Since ESS serves as a lower bound  
462 for the minimally realizable TSS (Section 3), we postulate that ESS can be used as a heuristic for  
463 conducting model-order reduction.

464 To test this, we distill multiple GLA models (with TSS = 256) across various task regimes to  
465 understand how the ESS of the original model (i.e. the teacher model) influences its ability to be  
466 distilled into a smaller student model. We apply the technique outlined in Bick et al. (2024), where  
467 the process can be divided into two steps. 1) matching the operators ( $\min(\|T_{(s)} - T_{(t)}\|_F^2 / \|T_{(t)}\|_F^2)$ )  
468 and 2) matching the output activations ( $\min(\|y_{(s)} - y_{(t)}\|_2^2 / \|y_{(t)}\|_2^2)$ ). More details can be found in  
469 Section D.6. Figure 6d (and more comprehensively Figure 30) shows the relationship between the  
470 ESS of the teacher model and the final activation loss during distillation.

471 **Finding 6:** Higher teacher ESS correlates with greater activation loss. The downstream perfor-  
472 mance after single-layer distillation depends on both the teacher model’s average ESS and student  
473 model’s TSS, with higher teacher ESS and lower student TSS resulting in greater performance  
474 loss (Figure 29).

475 We also note that directly comparing student ESS against teacher ESS provides additional insights  
476 into the effectiveness of the distillation process (Figure 31). These findings position ESS as a useful  
477 heuristic for predicting model compressibility, enabling efficient estimation of the potential for state-  
478 size reduction without extensive experimentation.

## 481 5.3 STATE MODULATION OF LARGE LANGUAGE MODELS

482 In contrast to synthetic tasks like MQAR, selective copying, and compression, we find that strong  
483 recall performance on language depends not only on a model having sufficient ESS, but also on its  
484 ability to dynamically modulate its ESS in response to inputs. We demonstrate that this explains why  
485

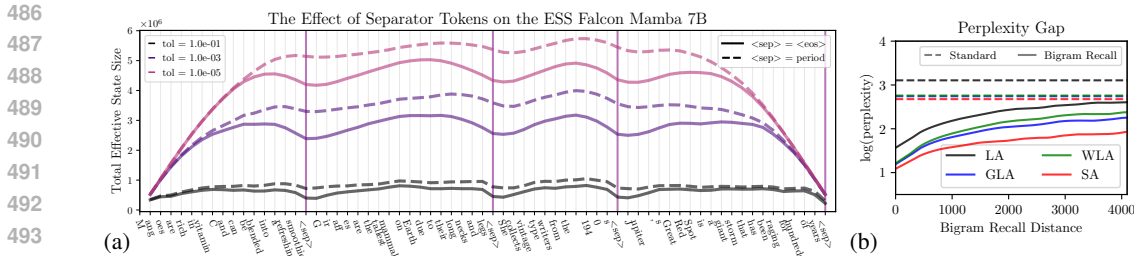


Figure 7: (a) The effect of separator tokens over Falcon Mamba 7B. See Section E.5 for plots of other open-weight models. (b) Comparison of standard perplexity and bigram recall perplexity (Arora et al., 2023).

linear attention, though effective on synthetic experiments (Section 4), is widely known to perform poorly on more complex language tasks (Katharopoulos et al., 2020; Arora et al., 2024).

We begin by evaluating the total ESS (computed across layers and channels) of open-weight pre-trained models. Our analysis shown in Figure 7a (and more broadly in Section E.5) reveals an intriguing phenomenon: ESS undergoes a noticeable dip whenever an end-of-speech (EOS) token is encountered (refer to Section D.8 for experimental details). This behavior aligns with our intuition regarding the role of EOS tokens and provides a quantitative measure of how effectively a model can ‘reset’ or ‘forget’ past contexts when transitioning between distinct segments of text<sup>7</sup>. To investigate these effects in a more controlled environment, we trained four 1B parameter models (LA, WLA, GLA, and SA as described in Section 4) under identical conditions (see Table 9).

**Finding 7:** We observe a clear hierarchy in the degree of state modulation, which can be summarized as follows: SA > GLA > WLA > LA (Figure 40).

SA exhibits the most pronounced state modulation, beginning at a tolerance level of  $1e-2$ , while also realizing the largest ESS. GLA follows, with modulation emerging at a tolerance of  $1e-1$ . WLA shows minimal modulation, only detectable at a tolerance of 1.0, while LA displays no discernible state modulation in response to separator tokens, demonstrating a clear lack of ability to modulate ESS. The importance of state modulation becomes apparent when examining model performance. Figure 7b illustrates that although standard perplexities (computed over a subset of the FineWeb dataset (Penedo et al., 2024)) are similar across SA, WLA, and GLA, significant differences emerge when considering the bigram recall perplexity metric introduced by Arora et al. (2023).

**Finding 8:** The ability of a model to recall information, as measured by bigram recall perplexity across a pre-training dataset (rather than within a narrow task space), reveals a performance hierarchy that closely mirrors the observed state modulation capabilities.

This relationship between bigram recall perplexity and state modulation suggests that state modulation serves as a key mechanism enabling models to effectively manage complex context dependencies commonly found in language, directly impacting their training dynamics and performance on recall-heavy tasks. Further implications and details are discussed in Section E.5.

## 6 CONCLUSION

In this work, we propose effective state-size (ESS), a measure of memory utilization in sequence models derived using dynamical systems theory. We motivate this metric as a valuable tool for analyzing memory utilization in LIVs by demonstrating its strong correlation with performance across a wide range of synthetic tasks. In doing so, we find that ESS offers a versatile framework for understanding both the performance and efficiency of causal sequence models. Leveraging these insights, we are able to construct novel, ESS-informed initializers, regularizers, and distillation strategies that improve beyond the existing performance-efficiency trade-offs found in recurrent models. Finally, we extend the ESS framework to language tasks, introducing the idea of state modulation – a concept which proves crucial for performance on bigram recall tasks. Overall, this work establishes ESS as a foundational tool for understanding and improving sequence model performance, opening new avenues for optimizing memory utilization and, more generally, model efficiency.

<sup>7</sup>We also observe a similar behavior with scope delimiters in code (Section E.6.1).

## 540 REPRODUCIBILITY STATEMENT

541  
542 To ensure reproducibility, we utilized open-source models and tasks, adhering to default task con-  
543 figurations unless otherwise specified. All crucial configurations are detailed in either the main text  
544 or the appendix. Additionally, our code for computing both the tolerance-ESS and entropy-ESS is  
545 provided in the appendix (Section D.1.1).  
546

## 547 REFERENCES

- 548  
549 Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers, 2020. URL  
550 <https://arxiv.org/abs/2005.00928>. (pages 2, 20).  
551
- 552 H. Akaike. Stochastic theory of minimal realization. *IEEE Transactions on Automatic Control*, 19  
553 (6):667–674, 1974. doi: 10.1109/TAC.1974.1100707. (page 3).
- 554 Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-context language learning: Ar-  
555 chitectures and algorithms, 2024. URL <https://arxiv.org/abs/2401.12973>. (page  
556 20).  
557
- 558 Ameen Ali, Itamar Zimmerman, and Lior Wolf. The hidden attention of mamba models, 2024. URL  
559 <https://arxiv.org/abs/2403.01590>. (pages 2, 20).  
560
- 561 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.3, knowledge capacity  
562 scaling laws, 2024. URL <https://arxiv.org/abs/2404.05405>. (page 2).
- 563 Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Al-  
564 subaie, Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwaresh, Areeb Alowisheq,  
565 M Saiful Bari, and Haidar Khan. When benchmarks are targets: Revealing the sensitivity of large  
566 language model leaderboards, 2024. URL <https://arxiv.org/abs/2402.01781>. (page  
567 21).  
568
- 569 Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri  
570 Rudra, and Christopher Ré. Zoology: Measuring and improving recall in efficient language mod-  
571 els, 2023. URL <https://arxiv.org/abs/2312.04927>. (pages 10, 10, 20, 21, 21, 30,  
572 30, 30).
- 573 Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalsina, Silas Alberti, Dylan Zinsley,  
574 James Zou, Atri Rudra, and Christopher Ré. Simple linear attention language models balance the  
575 recall-throughput tradeoff, 2024. URL <https://arxiv.org/abs/2402.18668>. (pages  
576 10, 20).  
577
- 578 Jimmy Ba, Geoffrey Hinton, Volodymyr Mnih, Joel Z. Leibo, and Catalin Ionescu. Using fast  
579 weights to attend to the recent past, 2016. URL <https://arxiv.org/abs/1610.06258>.  
580 (page 21).
- 581 Bartłomiej Beliczynski, Izzet Kale, and Gerald D Cain. Approximation of fir by iir digital filters:  
582 An algorithm based on balanced model reduction. *IEEE Transactions on Signal Processing*, 40  
583 (3):532–542, 1992. (page 9).  
584
- 585 Satwik Bhattamishra, Arkil Patel, Phil Blunsom, and Varun Kanade. Understanding in-context  
586 learning in transformers and llms by learning to learn discrete functions, 2023. URL <https://arxiv.org/abs/2310.03016>. (page 20).  
587
- 588 Srinadh Bhojanapalli, Chulhee Yun, Ankit Singh Rawat, Sashank J. Reddi, and Sanjiv Kumar. Low-  
589 rank bottleneck in multi-head attention models, 2020. URL <https://arxiv.org/abs/2002.07028>. (pages 2, 20).  
590
- 591 Aviv Bick, Kevin Y. Li, Eric P. Xing, J. Zico Kolter, and Albert Gu. Transformers to ssms: Distill-  
592 ing quadratic knowledge to subquadratic models, 2024. URL <https://arxiv.org/abs/2408.10189>. (page 9).  
593

- 594 Nick Cammarata, Shan Carter, Gabriel Goh, Chris Olah, Michael Petrov, Ludwig Schubert, Chelsea  
595 Voss, Ben Egan, and Swee Kiat Lim. Thread: Circuits. *Distill*, 2020. doi: 10.23915/distill.00024.  
596 <https://distill.pub/2020/circuits>. (page 20).  
597
- 598 Chi-Tsong Chen. *Linear System Theory and Design*. Oxford University Press, Inc., USA, 3rd  
599 edition, 1998. ISBN 0195117778. (page 1).
- 600 Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through  
601 structured state space duality, 2024. URL <https://arxiv.org/abs/2405.21060>. (pages  
602 1, 4, 6, 19, 20, 27, 29).  
603
- 604 Soham De, Samuel L. Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Al-  
605 bert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, Guillaume Des-  
606 jardins, Arnaud Doucet, David Budden, Yee Whye Teh, Razvan Pascanu, Nando De Freitas, and  
607 Caglar Gulcehre. Griffin: Mixing gated linear recurrences with local attention for efficient lan-  
608 guage models, 2024. URL <https://arxiv.org/abs/2402.19427>. (page 50).
- 609 P. DeWilde and A.J. van der Veen. *Time-Varying Systems and Computations*. Springer US,  
610 1998. ISBN 9780792381891. URL [https://books.google.co.jp/books?id=](https://books.google.co.jp/books?id=n3bEniJ2Wx8C)  
611 [n3bEniJ2Wx8C](https://books.google.co.jp/books?id=n3bEniJ2Wx8C). (pages 1, 3, 4, 19, 24).  
612
- 613 Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure  
614 attention loses rank doubly exponentially with depth, 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2103.03404)  
615 [abs/2103.03404](https://arxiv.org/abs/2103.03404). (page 20).
- 616 Emilien Dupont, Arnaud Doucet, and Yee Whye Teh. Augmented neural odes, 2019. URL <https://arxiv.org/abs/1904.01681>. (page 20).  
617
- 618 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,  
619 Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep  
620 Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt,  
621 Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and  
622 Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*,  
623 2021. <https://transformer-circuits.pub/2021/framework/index.html>. (page 20).  
624
- 625 Daniel Y. Fu, Tri Dao, Khaled K. Saab, Armin W. Thomas, Atri Rudra, and Christopher Ré. Hungry  
626 hungry hippo: Towards language modeling with state space models, 2023. URL [https://](https://arxiv.org/abs/2212.14052)  
627 [arxiv.org/abs/2212.14052](https://arxiv.org/abs/2212.14052). (page 20).
- 628 Wodek Gawronski and Jer-Nan Juang. Model reduction in limited time and frequency intervals. *In-*  
629 *ternational Journal of Systems Science*, 21(2):349–376, 1990. doi: 10.1080/00207729008910366.  
630 URL <https://doi.org/10.1080/00207729008910366>. (page 9).
- 631 Paolo Glorioso, Quentin Anthony, Yury Tokpanov, James Whittington, Jonathan Pilault, Adam  
632 Ibrahim, and Beren Millidge. Zamba: A compact 7b ssm hybrid model, 2024. URL [https://](https://arxiv.org/abs/2405.16712)  
633 [arxiv.org/abs/2405.16712](https://arxiv.org/abs/2405.16712). (page 50).  
634
- 635 Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural  
636 networks. In Yee Whye Teh and Mike Titterton (eds.), *Proceedings of the Thirteenth Interna-*  
637 *tional Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine*  
638 *Learning Research*, pp. 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR.  
639 URL <https://proceedings.mlr.press/v9/glorot10a.html>. (page 8).
- 640 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2024.  
641 URL <https://arxiv.org/abs/2312.00752>. (pages 1, 3, 6, 8, 19, 20, 27, 28).  
642
- 643 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured  
644 state spaces, 2022a. URL <https://arxiv.org/abs/2111.00396>. (pages 1, 3, 19, 25,  
645 27).  
646
- 647 Albert Gu, Ankit Gupta, Karan Goel, and Christopher Ré. On the parameterization and initialization  
of diagonal state space models, 2022b. URL <https://arxiv.org/abs/2206.11893>.  
(page 19).

- 648 Ramin Hasani, Mathias Lechner, Tsun-Hsuan Wang, Makram Chahine, Alexander Amini, and  
649 Daniela Rus. Liquid structural state-space models, 2022. URL [https://arxiv.org/abs/  
650 2209.12951](https://arxiv.org/abs/2209.12951). (pages 1, 6, 20).  
651
- 652 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja-  
653 cob Steinhardt. Measuring massive multitask language understanding, 2021. URL <https://arxiv.org/abs/2009.03300>. (page 21).  
654
- 655 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):  
656 1735–1780, nov 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL <https://doi.org/10.1162/neco.1997.9.8.1735>. (page 2).  
657  
658
- 659 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,  
660 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
661 models, 2020. URL <https://arxiv.org/abs/2001.08361>. (page 3).  
662
- 663 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are  
664 rnns: Fast autoregressive transformers with linear attention, 2020. URL <https://arxiv.org/abs/2006.16236>. (pages 1, 6, 10, 19, 20, 27).  
665
- 666 Serkan Kiranyaz, Onur Avci, Osama Abdeljaber, Turker Ince, Moncef Gabbouj, and Daniel J.  
667 Inman. 1d convolutional neural networks and applications: A survey, 2019. URL <https://arxiv.org/abs/1905.03554>. (page 1).  
668
- 669 Opher Lieber, Barak Lenz, Hofit Bata, Gal Cohen, Jhonathan Osin, Itay Dalmedigos, Erez Safahi,  
670 Shaked Meir, Yonatan Belinkov, Shai Shalev-Shwartz, Omri Abend, Raz Alon, Tomer Asida,  
671 Amir Bergman, Roman Glozman, Michael Gokhman, Avashalom Manevich, Nir Ratner, Noam  
672 Rozen, Erez Shwartz, Mor Zusman, and Yoav Shoham. Jamba: A hybrid transformer-mamba  
673 language model, 2024. URL <https://arxiv.org/abs/2403.19887>. (pages 34, 50, 50).  
674
- 675 Anders Lindquist and Giorgio Picci. On the stochastic realization problem. *SIAM J. Control Optim.*,  
676 17(3):365–389, May 1979. ISSN 0363-0129. doi: 10.1137/0317028. URL <https://doi.org/10.1137/0317028>. (page 3).  
677
- 678 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>. (pages 32, 33, 34, 35).  
679
- 680 Paul A. Lynn and Wolfgang Fuerst. *Introductory digital signal processing with computer applica-  
681 tions (revised ed.)*. John Wiley & Sons, Inc., USA, 1994. ISBN 0471943746. (page 1).  
682
- 683 Eric Martin and Chris Cundy. Parallelizing linear recurrent neural nets over sequence length, 2018.  
684 URL <https://arxiv.org/abs/1709.04057>. (page 3).  
685
- 686 Stefano Massaroli, Michael Poli, Jinkyoo Park, Atsushi Yamashita, and Hajime Asama. Dissecting  
687 neural odes, 2021. URL <https://arxiv.org/abs/2002.08071>. (page 1).  
688
- 689 Stefano Massaroli, Michael Poli, Daniel Y. Fu, Hermann Kumbong, Rom N. Parnichkun, Aman  
690 Timalisina, David W. Romero, Quinn McIntyre, Beidi Chen, Atri Rudra, Ce Zhang, Christopher  
691 Re, Stefano Ermon, and Yoshua Bengio. Laughing hyena distillery: Extracting compact recur-  
692 rences from convolutions, 2023. URL <https://arxiv.org/abs/2310.18780>. (pages 9,  
693 19).  
694
- 695 Zeping Min and Zhong Li. On the efficiency of transformers: The effect of attention rank, 2024.  
696 URL <https://openreview.net/forum?id=U9sHVjIdYH>. (pages 2, 20).  
697
- 698 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,  
699 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Gan-  
700 guli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion,  
701 Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam  
McCandlish, and Chris Olah. In-context learning and induction heads. *Transformer Cir-  
cuits Thread*, 2022a. [https://transformer-circuits.pub/2022/in-context-learning-and-induction-  
heads/index.html](https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html). (pages 2, 20).

- 702 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,  
703 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli,  
704 Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane  
705 Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish,  
706 and Chris Olah. In-context learning and induction heads. *Transformer Circuits Thread*, 2022b.  
707 <https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html>. (page  
708 20).
- 709 Alan V. Oppenheim, Alan S. Willsky, and S. Hamid Nawab. *Signals & systems (2nd ed.)*. Prentice-  
710 Hall, Inc., USA, 1996. ISBN 0138147574. (page 1).
- 711
- 712 Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pas-  
713 canu, and Soham De. Resurrecting recurrent neural networks for long sequences, 2023. URL  
714 <https://arxiv.org/abs/2303.06349>. (pages 19, 27).
- 715
- 716 Rom N. Parnichkun, Stefano Massaroli, Alessandro Moro, Jimmy T. H. Smith, Ramin Hasani, Math-  
717 ias Lechner, Qi An, Christopher Ré, Hajime Asama, Stefano Ermon, Taiji Suzuki, Atsushi Ya-  
718 mashita, and Michael Poli. State-free inference of state-space models: The transfer function  
719 approach, 2024. URL <https://arxiv.org/abs/2405.06147>. (pages 1, 19, 27).
- 720
- 721 Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural  
722 networks, 2013. URL <https://arxiv.org/abs/1211.5063>. (page 3).
- 723
- 724 Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin  
725 Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the  
726 finest text data at scale, 2024. (pages 10, 36).
- 727
- 728 Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y. Fu, Tri Dao, Stephen Baccus, Yoshua  
729 Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional  
730 language models, 2023. URL <https://arxiv.org/abs/2302.10866>. (pages 1, 3, 6,  
731 19).
- 732
- 733 Michael Poli, Armin W Thomas, Eric Nguyen, Pragaash Ponnusamy, Björn Deiseroth, Kris-  
734 tian Kersting, Taiji Suzuki, Brian Hie, Stefano Ermon, Christopher Ré, Ce Zhang, and Ste-  
735 fano Massaroli. Mechanistic design and scaling of hybrid architectures, 2024. URL <https://arxiv.org/abs/2403.17844>. (pages 5, 8, 20, 30).
- 736
- 737 Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Gener-  
738 alization beyond overfitting on small algorithmic datasets, 2022. URL <https://arxiv.org/abs/2201.02177>. (page 20).
- 739
- 740 Hubert Ramsauer, Bernhard Schäfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Thomas  
741 Adler, Lukas Gruber, Markus Holzleitner, Milena Pavlović, Geir Kjetil Sandve, Victor Greiff,  
742 David Kreil, Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter.  
743 Hopfield networks is all you need, 2021. URL <https://arxiv.org/abs/2008.02217>.  
744 (page 21).
- 745
- 746 David W. Romero, Anna Kuzina, Erik J. Bekkers, Jakub M. Tomczak, and Mark Hoogendoorn.  
747 Ckconv: Continuous kernel convolution for sequential data, 2022. URL <https://arxiv.org/abs/2102.02611>. (page 1).
- 748
- 749 Olivier Roy and Martin Vetterli. The effective rank: A measure of effective dimensionality. In *2007*  
750 *15th European Signal Processing Conference*, pp. 606–610, 2007. (pages 3, 5).
- 751
- 752 David E. Rumelhart and James L. McClelland. *Learning Internal Representations by Error Propa-*  
753 *gation*, pp. 318–362. 1987. (page 2).
- 754
- 755 Noam Shazeer. Glu variants improve transformer, 2020. URL <https://arxiv.org/abs/2002.05202>. (page 8).
- Huitao Shen. Mutual information scaling and expressive power of sequence models, 2019. URL <https://arxiv.org/abs/1905.04271>. (page 20).

- 756 Jimmy T. H. Smith, Andrew Warrington, and Scott W. Linderman. Simplified state space layers for  
757 sequence modeling, 2023. URL <https://arxiv.org/abs/2208.04933>. (pages 1, 19,  
758 25).
- 759 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Reformer: En-  
760 hanced transformer with rotary position embedding, 2023. URL <https://arxiv.org/abs/2104.09864>. (page 28, 28).
- 763 Mingjie Sun, Xinlei Chen, J. Zico Kolter, and Zhuang Liu. Massive activations in large language  
764 models, 2024. URL <https://arxiv.org/abs/2402.17762>. (pages 2, 20).
- 765 Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and  
766 Furu Wei. Retentive network: A successor to transformer for large language models, 2023. URL  
767 <https://arxiv.org/abs/2307.08621>. (pages 6, 19).
- 769 Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan  
770 Salakhutdinov. Transformer dissection: A unified understanding of transformer’s attention via  
771 the lens of kernel, 2019. URL <https://arxiv.org/abs/1908.11775>. (pages 1, 19).
- 772 Neehal Tumma, Mathias Lechner, Noel Loo, Ramin Hasani, and Daniela Rus. Leveraging low-rank  
773 and sparse recurrent connectivity for robust closed-loop control, 2023. URL <https://arxiv.org/abs/2310.03915>. (page 20).
- 774 Raf Vandebril, Marc Van Barel, and Nicola Mastronardi. A note on the representation and definition  
775 of semiseparable matrices. *Numerical Linear Algebra with Applications*, 12(8):839–858, 2005.  
776 doi: <https://doi.org/10.1002/nla.455>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/nla.455>. (page 4).
- 777 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
778 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>. (pages 1, 1, 3, 6, 19, 20).
- 781 Jesse Vig. A multiscale visualization of attention in the transformer model, 2019. URL <https://arxiv.org/abs/1906.05714>. (pages 2, 20).
- 782 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming  
783 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi  
784 Fan, Xiang Yue, and Wenhu Chen. Mmlu-pro: A more robust and challenging multi-task language  
785 understanding benchmark, 2024. URL <https://arxiv.org/abs/2406.01574>. (page  
786 21).
- 787 Jan C. Willems. *Models for Dynamics*, pp. 171–269. Vieweg+Teubner Verlag, Wiesbaden, 1989.  
788 ISBN 978-3-322-96657-5. doi: [10.1007/978-3-322-96657-5\\_5](https://doi.org/10.1007/978-3-322-96657-5_5). URL [https://doi.org/10.1007/978-3-322-96657-5\\_5](https://doi.org/10.1007/978-3-322-96657-5_5). (page 2).
- 791 Xinyi Wu, Amir Ajorlou, Yifei Wang, Stefanie Jegelka, and Ali Jadbabaie. On the role of atten-  
792 tion masks and layernorm in transformers, 2024. URL <https://arxiv.org/abs/2405.18781>. (pages 2, 20).
- 793 Jianlin Xia, Shivkumar Chandrasekaran, Ming Gu, and Xiaoye S. Li. Fast algorithms for hierarchi-  
794 cally semiseparable matrices. *Numerical Linear Algebra with Applications*, 17(6):953–976, 2010.  
795 doi: <https://doi.org/10.1002/nla.691>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/nla.691>. (page 4).
- 800 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming lan-  
801 guage models with attention sinks, 2024. URL <https://arxiv.org/abs/2309.17453>.  
802 (pages 2, 20).
- 803 Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. Gated linear atten-  
804 tion transformers with hardware-efficient training, 2024a. URL <https://arxiv.org/abs/2312.06635>. (pages 1, 6, 8, 19, 20, 33).

810 Songlin Yang, Bailin Wang, Yu Zhang, Yikang Shen, and Yoon Kim. Parallelizing linear trans-  
811 formers with the delta rule over sequence length, 2024b. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2406.06484)  
812 [2406.06484](https://arxiv.org/abs/2406.06484). (page 8).  
813

814 Luca Zancato, Arjun Seshadri, Yonatan Dukler, Aditya Golatkar, Yantao Shen, Benjamin Bowman,  
815 Matthew Trager, Alessandro Achille, and Stefano Soatto. B’mojo: Hybrid state space realizations  
816 of foundation models with eidetic and fading memory, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2407.06324)  
817 [abs/2407.06324](https://arxiv.org/abs/2407.06324). (pages 1, 3).

818 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-  
819 chine really finish your sentence?, 2019. URL <https://arxiv.org/abs/1905.07830>.  
820 (page 21).  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863



---

# Supplementary Material

---

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

## CONTENTS

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>   | <b>1</b>  |
| <b>2</b> | <b>Related Work</b>   | <b>2</b>  |
| <b>3</b> | <b>Theory</b>   | <b>3</b>  |
| 3.1      | Computing Effective State-Size . . . . .                              | 4         |
| <b>4</b> | <b>Empirical Validation of Effective State-Size</b>                   | <b>5</b>  |
| 4.1      | Cross Task-Model Analysis . . . . .                                   | 6         |
| 4.2      | Within Task-Model Analysis . . . . .                                  | 7         |
| <b>5</b> | <b>Applications of Effective State-Size</b>                           | <b>8</b>  |
| 5.1      | Initialization-Phase Analysis . . . . .                               | 8         |
| 5.2      | Post-Training Analysis . . . . .                                      | 9         |
| 5.3      | State Modulation of Large Language Models . . . . .                   | 9         |
| <b>6</b> | <b>Conclusion</b>   | <b>10</b> |
| <b>A</b> | <b>Mid-Training Analysis</b>  | <b>19</b> |
| <b>B</b> | <b>Extended Related Work</b>  | <b>19</b> |
| <b>C</b> | <b>Theoretical Background</b>   | <b>22</b> |
| C.1      | Notation . . . . .  | 22        |
| C.2      | Derivations and Proofs . . . . .                                      | 22        |
| C.2.1    | The Operator Realization of Linear Recurrences . . . . .              | 22        |
| C.2.2    | Factorizing The Operator Realization Submatrix $H_i$ . . . . .        | 22        |
| C.2.3    | The Trivial Recurrence Realization . . . . .                          | 23        |
| C.2.4    | Minimal Recurrent Realization (Proof of Theorem 3.1) . . . . .        | 24        |
| <b>D</b> | <b>Methods</b>  | <b>25</b> |
| D.1      | Computing ESS . . . . .   | 25        |
| D.1.1    | PyTorch Implementation . . . . .                                      | 25        |
| D.2      | Formulation of the Featurizers . . . . .                              | 27        |
| D.3      | Empirical Validation . . . . .  | 30        |
| D.4      | ESS-Informed Featurizer Selection and Initialization Scheme . . . . . | 31        |

|     |  |           |
|-----|--|-----------|
| 918 | D.5 ESS-Informed Regularization . . . . .  | 33        |
| 919 | D.6 ESS-Informed Model-Order Reduction . . . . .   | 33        |
| 920 | D.7 ESS Analysis for Hybrid Networks . . . . .   | 34        |
| 921 | D.8 State Modulation of Large Language Models . . . . .  | 35        |
| 922 |  |           |
| 923 |  |           |
| 924 |  |           |
| 925 | <b>E Extended Experimental Results</b>   | <b>37</b> |
| 926 | E.1 Empirical Validation . . . . .   | 37        |
| 927 | E.1.1 State Collapse Continued . . . . .   | 37        |
| 928 | E.1.2 Entropy-ESS MQAR Results Continued . . . . .   | 37        |
| 929 | E.1.3 Tolerance-ESS MQAR Results . . . . .   | 40        |
| 930 | E.1.4 Selective Copying and Compression Results . . . . .  | 43        |
| 931 | E.1.5 ESS Training Dynamics in MQAR . . . . .  | 45        |
| 932 | E.2 Initialization-Phase Analysis . . . . .  | 46        |
| 933 | E.3 Mid-Training Analysis . . . . .  | 46        |
| 934 | E.4 Post-Training Analysis . . . . .   | 48        |
| 935 | E.4.1 Model-Order Reduction . . . . .  | 48        |
| 936 | E.4.2 Hybridization . . . . .  | 49        |
| 937 | E.5 State Modulation of Large Language Models . . . . .  | 51        |
| 938 | E.6 Miscellaneous . . . . .  | 54        |
| 939 | E.6.1 Effective State-Size on C++ Code . . . . .   | 54        |
| 940 | E.6.2 How the Number of Prompting Shots Affects the Effective State-Size of<br>Language Models . . . . . | 55        |
| 941 |  |           |
| 942 |  |           |
| 943 |  |           |
| 944 |  |           |
| 945 |  |           |
| 946 |  |           |
| 947 |  |           |
| 948 |  |           |
| 949 |  |           |
| 950 |  |           |
| 951 |  |           |
| 952 |  |           |
| 953 |  |           |
| 954 |  |           |
| 955 |  |           |
| 956 |  |           |
| 957 |  |           |
| 958 |  |           |
| 959 |  |           |
| 960 |  |           |
| 961 |  |           |
| 962 |  |           |
| 963 |  |           |
| 964 |  |           |
| 965 |  |           |
| 966 |  |           |
| 967 |  |           |
| 968 |  |           |
| 969 |  |           |
| 970 |  |           |
| 971 |  |           |

## A MID-TRAINING ANALYSIS

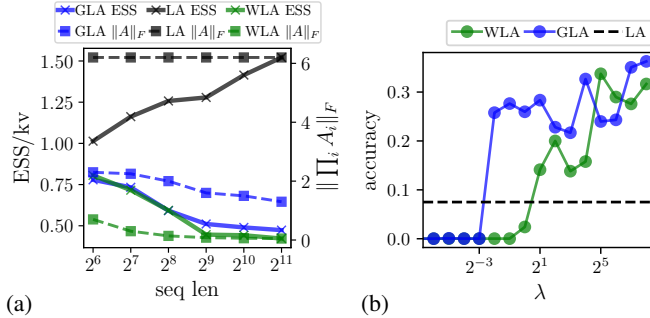


Figure 8: (a) ESS/kv and  $\|\prod_i A_i\|_F$  as a function of sequence length. (b) Accuracy of models as a function of ESS-based regularizer strength.

To motivate the idea of increasing ESS mid-training, we revisit the concept of state collapse – a phenomenon that arises due to trainability issues in recurrent models (Figure 27), as discussed in Section 4.2. Recall that state collapse describes a failure mode of learning in GLA and WLA which, unlike LA, have learnable  $A_i$  matrices (where  $i$  denotes the index along the sequence dimension). To see why this contributes to state collapse, we note that the values of the operator submatrices  $H_i$  are disproportionately influenced by  $A_i$ , due to the presence of terms in the form of  $A_{i-1} \dots A_1$  for each  $i$ . Hence, the closer  $A_i$  lies to the 0-matrix, the faster these terms decay, reducing the numerical rank of  $H_i$ . We demonstrate this empirically in Figure 8a, which shows that for both GLA and WLA, ESS/kv and  $\|\prod_i A_i\|_F$  decrease as a function of sequence length. In contrast, for LA, whose  $A$  matrix is given by the identity, ESS/kv remains large as the sequence length grows.

Given this insight, one approach to addressing state collapse in GLA and WLA is pushing the  $A$  matrices towards the identity by adding the following term to the loss function:  $\lambda\|A - I\|_F$ , where  $\lambda$  denotes the strength of the regularizer and  $I$  denotes the identity. In doing so, we are effectively decaying the model towards LA, increasing its ESS and giving us the following:

**Finding 9:** GLA and WLA trained using the ESS-based regularization scheme described above outperform LA. When trained without it, they perform worse than LA (Figure 8b).

For more commentary on this result, please refer to Section E.3.

## B EXTENDED RELATED WORK

**Causal sequence models.** From classical linear recurrences to modern sequence models like Transformers, a vast array of causal model architectures have emerged (Vaswani et al., 2023; Tsai et al., 2019; Katharopoulos et al., 2020; Poli et al., 2023; Yang et al., 2024a; Gu & Dao, 2024; Dao & Gu, 2024; Sun et al., 2023). In recent years, the ability to process sequences in parallel has become increasingly critical, largely due to advancements in hardware accelerators such as GPUs. This need for parallelism likely explains the growing popularity of models like attention, Mamba, and S4.

We observe that all of these models, which support parallelization across the sequence dimension, can be formulated using a linear system representation ( $y = Tu$ ) as detailed in the introductory and theoretical sections (Sections 1 and 3). For this work, we categorize these models into two types: linear systems and systems with input-varying linear operators. The key distinction between these two frameworks is that in the former, the operator  $T$  in input-invariant models is composed of fixed system parameters, whereas in the latter, the parameters are dynamically generated from the input. Linear systems encompass both linear time-varying (LTV) and linear time-invariant (LTI) systems. Although LTV systems have been relatively unexplored in deep learning, several LTI models have been studied (Gu et al., 2022a;b; Smith et al., 2023; Orvieto et al., 2023; Parnichkun et al., 2024). Convolutional models use kernels  $h$  to construct Hankel matrices  $H$ , whose rank corresponds to the minimal state-size of the model (DeWilde & van der Veen, 1998). Massaroli et al. (2023) explored methods to reduce the order of models by leveraging the Hankel matrix. Notably, the submatrix  $H_i$  (defined in Theory 3.2) exhibits a Hankel structure in LTI models and provides per-sequence-index

1026 information. In this work, however, we do not explore Hankel matrices further, as they are not easily  
1027 generalizable to LTV systems.

1028 In contrast, systems with input-varying linear operators, which we abbreviate as LIVs, are char-  
1029 acterized by an operator  $T$  that is dynamically constructed through a featurizer and is defined by  
1030  $T = f(u)$ . Examples of such models include softmax attention (Vaswani et al., 2023), linear atten-  
1031 tion (Katharopoulos et al., 2020), Liquid-S4 (Hasani et al., 2022), Mamba (Gu & Dao, 2024; Dao &  
1032 Gu, 2024), and gated linear attention (Yang et al., 2024a). Although these models may appear non-  
1033 linear, they can still be represented like a linear system, enabling the application of linear analysis  
1034 techniques. This forms the basis for the effective state-size metric.

1036 **Interpretability.** Analysis tools for sequence models can be categorized into two types: extrinsic  
1037 and intrinsic. Extrinsic tools focus solely on the input and output, treating the model’s internal  
1038 processes as black boxes. This approach is highly generalizable, as it can be applied to any model,  
1039 including those with non-linear recurrences. A notable example by Shen (2019) uses statistical  
1040 measures such as mutual information to compute metrics that capture model “expressivity”. While  
1041 these methods are versatile and applicable to various datasets, their generality makes them less  
1042 effective at capturing the inner workings of causal sequence models, which is the primary focus of  
1043 this work.

1044 Intrinsic tools, conversely, directly visualize the model’s internal mechanisms. A recently popular  
1045 framework known as mechanistic interpretability provides one such example (Power et al., 2022).  
1046 Mechanistic interpretability involves dissecting complex models to understand how specific com-  
1047 ponents contribute to the model’s overall behavior (Cammarata et al., 2020). Unlike our work,  
1048 mechanistic interpretability does not target the operator view of the model but instead emphasizes  
1049 the functional roles and interactions of individual model components.

1050 For our purposes, we are primarily concerned with the visualization and analysis of classical and  
1051 modern causal sequence models through the unifying lens of linear operators. Most analyses of  
1052 these operators rely on visualization techniques (Olsson et al., 2022a; Vig, 2019; Abnar & Zuidema,  
1053 2020; Ali et al., 2024; Xiao et al., 2024; Sun et al., 2024) to gain insights into the model’s internal  
1054 processes. Visualizing the operator  $T$  is advantageous, as it reveals important features like the  
1055 formation of induction heads, strong activations, diagonal and block-diagonal patterns, and Toeplitz  
1056 structures. However, raw visualizations are largely qualitative and oftentimes do not provide the  
1057 quantitative metrics necessary for effectively evaluating a model’s internal mechanisms – a gap we  
1058 aim to address in this work.

1059 Other, more quantitative, intrinsic methods perform some form of spectral analysis on the full oper-  
1060 ator (Dong et al., 2023; Min & Li, 2024; Tumma et al., 2023; Bhojanapalli et al., 2020). A limitation  
1061 of these approaches is that they often disregard the causal masking of  $T$ , which significantly impacts  
1062 the model’s rank and singular values (Wu et al., 2024). As a result, the rank of the causal operator  
1063  $T$  alone lacks a clear interpretation.

1064 The proposed effective state-size metric is a method applicable to both linear systems and LIVs.  
1065 As a quantitative proxy for memory utilization, it offers insights into the inner workings of causal  
1066 sequence models, ensuring generality, usability, and interpretability.

1068 **Synthetic and language benchmarks.** In this work, we build on synthetic tasks from the mech-  
1069 anistic architecture design (MAD) framework introduced in (Poli et al., 2024). MAD defines a set  
1070 of small-scale tasks designed to evaluate key model capabilities, such as in-context recall (Akyürek  
1071 et al., 2024; Bhattamishra et al., 2023; Elhage et al., 2021; Olsson et al., 2022b). Training models  
1072 on these tasks is efficient, making them well-suited for exploring a large space of tasks and models,  
1073 as demonstrated in several prior works (Dupont et al., 2019; Arora et al., 2024; Fu et al., 2023). In  
1074 this work, we investigate the effective state-size across a subset of the MAD tasks: multi-query as-  
1075 sociative recall (MQAR), selective copying, and compression, varying the difficulty of each to gain  
1076 a nuanced understanding of how effective state-size evolves across these task landscapes.

1077 Among the synthetic tasks we examine, MQAR stands out in particular. Proposed by Arora et al.  
1078 (2023), MQAR was designed to bridge the gap between synthetic and real language tasks explained  
1079 by associative recall – the ability of a model to retrieve information based on relationships between  
different elements in its memory. This capability has long been sought after in the construction

1080 of sequence model architectures (Ramsauer et al., 2021; Ba et al., 2016); as such, we evaluate the  
1081 performance of our models on MQAR to measure the benefits of using effective state-size to iterate  
1082 on canonical frameworks used in sequence modeling.

1083 One notable aspect of MQAR observed in Arora et al. (2023) is that the size of the model cache  
1084 needs to scale with the difficulty of the task to maintain performance. While this observation holds,  
1085 our work demonstrates that model cache size is an imperfect measure in this context due to the dis-  
1086 crepancy between memory capacity, as measured by theoretically realizable state size, and memory  
1087 utilization, as measured by effective state-size. At a higher level, this demonstrates how our work  
1088 provides a new perspective on analyzing memory-intensive synthetic tasks.

1089 While the MAD framework and synthetic tasks have shown correlations with model performance  
1090 on large-scale language tasks, language itself poses a unique challenge. Models are tasked with  
1091 predicting the next token given previous tokens – a simple yet general objective. New tasks can be  
1092 created simply by altering the prompts, thereby expanding the range of possible task domains.

1093 Although numerous language evaluation tasks – such as those in Hendrycks et al. (2021); Wang  
1094 et al. (2024); Zellers et al. (2019) – have been proposed, they often probe a narrow task space and  
1095 tend to be brittle. For example, shuffling the order of multiple choices in MMLU can drastically  
1096 change model rankings Alzahrani et al. (2024).  
1097

1098 Unlike narrow benchmarks, perplexity scores can be computed across an entire pre-training dataset,  
1099 covering a much broader task domain. However, small perplexity gaps between models make it a  
1100 challenging metric for evaluation. Recently, Arora et al. found that much of the difference in per-  
1101 plexity between models can be attributed to bigram perplexity – a measure of a model’s ability to  
1102 utilize the context and predict a successor token (second token of a bigram) given a repeated context  
1103 token (first token of a bigram) within a sequence. They demonstrate that most of the average per-  
1104 plexity difference between a gated convolution model and an attention model stems from differences  
1105 in bigram perplexity, suggesting that recall is a key capability for language models.

1106 The effective state-size analysis presented in this work reveals that strong recall performance as  
1107 measured by bigram perplexity in language modeling tasks depends not only on memory capacity,  
1108 but also on a model’s ability to modulate its state-size within a given context.  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

## C THEORETICAL BACKGROUND

### C.1 NOTATION

We adopt the following notation in this paper:

- Inputs, outputs, and operators follow flattened notation. i.e.,  $u, y \in \mathbb{R}^{\ell d}$  and  $T \in \mathbb{R}^{\ell d \times \ell d}$ . In particular, the original inputs and outputs with shape  $\ell \times d$  are flattened in row-major ordering, resulting in  $T$  having  $\ell \times \ell$  sub-blocks, each of which is of size  $d \times d$ .
- Tensor subscripts index sequence indices (time-step) and superscripts index channel/hidden dimensions. I.e., for an input  $u \in \mathbb{R}^{\ell d}$ ,  $u_i \in \mathbb{R}^d$  denotes the input vector at sequence-index  $i$ , and  $u^\alpha \in \mathbb{R}^\ell$  denotes the input vector for channel  $\alpha$ . Similarly,  $T_{ij} \in \mathbb{R}^{d \times d}$  denotes the linear weighing of  $u_j$  on to  $y_i$ .
- Indices within square brackets indicate matrix indices void of semantics (sequence index, channels, etc.). I.e.,  $A_{i[\alpha, \beta]}$  indexes row  $\alpha$  and column  $\beta$  of matrix  $A_i$ .
- Semicolons within subscripts denote a product over ranges ( $A_{1;3} = A_1 A_2 A_3$ ).
- Tensor slices are denoted with colons and are inclusive over the ranges. I.e.,  $u_{0:2} = u_0 u_1 u_2$ .

### C.2 DERIVATIONS AND PROOFS

#### C.2.1 THE OPERATOR REALIZATION OF LINEAR RECURRENCES

Unrolling the recurrence in Equation 1 unveils the following formulation:

$$\begin{aligned}
 s_0 &= 0 \\
 s_1 &= B_0 u_0 \\
 s_2 &= B_1 u_1 + A_1 (B_0 u_0) \\
 s_3 &= B_2 u_2 + A_2 (B_1 u_1 + A_1 (B_0 u_0)) \\
 s_i &= \left( \sum_{j=0}^{i-1} \left[ \prod_{k=i-1}^{j+1} A_k \right] B_j u_j \right), \tag{C.2.1}
 \end{aligned}$$

$$y_i = C_i \left( \sum_{j=1}^{i-1} \left[ \prod_{k=i-1}^{j+1} A_k \right] B_j u_j \right) + D_i u_i, \tag{C.2.2}$$

which corresponds to the operator:

$$T_{ij} = \begin{cases} 0 & i < j \\ D_i & i = j \\ C_i A_{i-1; j+1} B_j & i > j \end{cases}. \tag{C.2.3}$$

#### C.2.2 FACTORIZING THE OPERATOR REALIZATION SUBMATRIX $H_i$

Factorizing the strictly lower triangular submatrices of the operator ( $T_{i:,i-1}$ ) into causal and anti-causal factors, unveils that  $n_i$  (i.e. the TSS) upper bounds the dimensionality of the inner product

between the factors, and thus, also the rank of the submatrix ( $n_i \geq \text{rank}(H_i)$ ):

$$\begin{aligned}
T_{i:,i-1} \equiv H_i &= \begin{bmatrix} C_i & & \\ & \ddots & \\ & & C_{\ell-1} \end{bmatrix} \begin{bmatrix} A_{i-1;1} & \dots & I \\ \vdots & \ddots & \vdots \\ A_{\ell-2;1} & \dots & A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} B_0 & & \\ & \ddots & \\ & & B_{i-1} \end{bmatrix} \\
&= \begin{bmatrix} C_i & & \\ & \ddots & \\ & & C_{\ell-1} \end{bmatrix} \begin{bmatrix} I \\ A_i \\ A_{i+1;i} \\ \vdots \\ A_{\ell-2;i} \end{bmatrix} [A_{i-1;1} \quad A_{i-1;2} \quad \dots \quad A_{i-1} \quad I] \begin{bmatrix} B_0 & & \\ & \ddots & \\ & & B_{i-1} \end{bmatrix} \\
&= \underbrace{\begin{bmatrix} C_i \\ C_{i+1}A_i \\ C_{i+2}A_{i+1;i} \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix}}_{\substack{d(\ell-i) \times n_i \\ \text{anti-causal}}} \underbrace{[A_{i-1;1}B_0 \quad A_{i-1;2}B_1 \quad \dots \quad A_{i-1}B_{i-2} \quad B_{i-1}]}_{\substack{n_i \times d_i \\ \text{causal}}} \equiv \mathcal{O}_i C_i.
\end{aligned} \tag{C.2.4}$$

Besides unveiling the relationship between the rank of the realized operator and the original state-size  $n_i$ , the following insights can be drawn from the decomposition:

- The causal portion  $\mathcal{C}_i$  is the input-state projection matrix at time-step  $i$  (i.e.,  $s_i = \mathcal{C}_i u_{i-1}$ ) corresponding to Equation (C.2.1).
- ESS ( $\text{rank}(H_i)$ ) is simply the minimum rank between the causal and anti-causal projections.
- In conjunction with Theorem 3.2, we observe that the causally determinable minimal state-size (causal ESS) is equivalent to the rank of the causal projection. This insight allows us to construct a more efficient realization of the recurrence:
  - We can minimally factorize the causal projection as  $\mathcal{C}_i = L_i R_i$ , where  $L_i \in \mathbb{R}^{n_i \times r}$  and  $R_i \in \mathbb{R}^{r \times d_i}$ , with  $r = \text{rank}(\mathcal{C}_i)$ .
  - The right factor  $R_i$  becomes the new input-state projection matrix for  $H_i$ , effectively reducing the state dimension to the causal ESS.
  - $A_{i-1}^*$  and  $B_{i-1}^*$  can be determined from  $R_i$  using the process outlined in Theorem 3.1, and  $C_i^* = C_i L_i$ .

### C.2.3 THE TRIVIAL RECURRENCE REALIZATION

Any input-varying and input-invariant causal operator can be trivially realized with the following recurrence:

$$\begin{aligned}
s_{i+1} &= \begin{bmatrix} I_{(d_i)} \\ 0_{(d)} \end{bmatrix} s_i + \begin{bmatrix} 0_{(d_i)} \\ I_{(d)} \end{bmatrix} u_i, \\
y_i &= [T_{i,0} \quad T_{i,1} \quad \dots \quad T_{i,i-1}] s_i + T_{i,i} u_i.
\end{aligned} \tag{C.2.5}$$

In simple terms, the state  $s_i$  stores each input from  $t \in [i-1]$ , which is then mapped to the output with operator features at row  $i$ . Note that in the case where the operator is input varying, the trivial realization upholds the causality of the featurization process (i.e. the features  $(A_i, B_i, C_i, D_i)_{i \in [\ell]}$  of the trivial realization are causally determined). Moreover, the causally determined ESS (see Section C.2.2) for the trivially realized recurrence is equivalent to its TSS, as  $\mathcal{C}_i = I_{d_i}$ .

## C.2.4 MINIMAL RECURRENT REALIZATION (PROOF OF THEOREM 3.1)

**Theorem 3.1** *Given any causal input-invariant operator  $T$ , there exist infinite variations of linear recurrences in the form of Equation (1) that realize an equivalent input-output operator.*

*Proof.* We first categorize the operator into two portions: the memoryless portion, where  $i = j$ , and the dynamical portion, where  $i > j$ . The memoryless portion can be trivially realized by setting  $D_i = T_{ii}$ . For the dynamical portion, we draw inspiration from (DeWilde & van der Veen, 1998, ch. 3) and approach the proof of existence by ansatz. The following steps outline the proof:

1. Section C.2.2 demonstrates that, given a linear recurrence in the form of Equation (1), the operator submatrix can be factorized into causal and anti-causal parts, where the causal part represents the input-state projection matrix. Therefore, we proceed by making the following ansatz: for any operator submatrix  $T_{i:,i-1} \equiv H_i$ ,  $H_i$  can be arbitrarily factorized into  $\mathcal{O}_i \in \mathbb{R}^{d(\ell-i) \times n_i}$  and  $\mathcal{C}_i \in \mathbb{R}^{n_i \times d_i}$ , and that  $\mathcal{C}_i$  represents the input-state projection at time-step  $i$  (i.e.,  $s_i = \mathcal{C}_i u_{i-1}$ ).
2. Construct the dynamic features  $(A_i, B_i, C_i)_{i \in [\ell]}$  such that the assumption above holds. Note that we additionally assume the initial and final states to be 0 without loss of generality, therefore the realization of  $C_0, A_0, A_{\ell-1}$ , and  $B_{\ell-1}$  could be ignored.
  - (a) Set  $\mathcal{C}_i = \mathcal{O}_{i[:,d-1]}$  to obtain  $(C_i)_{i \in [1,\ell]}$ , as given the assumptions above, the first set of rows of  $\mathcal{O}_i$  linearly projects  $s_i$  onto  $y_i - D_i u_i$ , which is identical to  $C_i$  in Equation (1).
  - (b) Set  $B_{i-1} = \mathcal{C}_{i[:, -d]}$  to obtain  $(B_i)_{i \in [\ell-1]}$ , for which the identity can be obtained by deconstructing the input-state projection matrix  $\mathcal{C}_i$  and equating its assumed state  $s_i$  with Equation (1).

$$\begin{aligned} s_i &= A_{i-1} s_{i-1} + B_{i-1} u_{i-1} \\ &= \mathcal{C}_{i[:, : -d-1]} u_{i-2} + \mathcal{C}_{i[:, -d]} u_{i-1}. \end{aligned} \quad (\text{C.2.6})$$

- (c) Using the same state-dynamics equation, we could equate the assumed state-projection matrices with each other, obtaining  $(A_i)_{i \in [1,\ell-1]}$ :

$$\begin{aligned} s_{i+1} &= A_i s_i + B_i u_i \\ \mathcal{C}_{i+1} u_i &= A_i \mathcal{C}_i u_{i-1} + \mathcal{C}_{i+1[:, -d]} u_i \\ \mathcal{C}_{i+1[:, : -d-1]} u_{i-1} &= A_i \mathcal{C}_i u_{i-1} \\ A_i &= \mathcal{C}_{i+1[:, : -d-1]} \mathcal{C}_i^+. \end{aligned} \quad (\text{C.2.7})$$

3. Verify that the realized recurrence maps back to the original operator  $T_{ij}$ , proving that arbitrary factorizations (of which there are infinite variations) of the operator submatrices can be used to construct equivalent operators.

$$\begin{aligned} T_{ij} &= C_i A_{i-1} \cdots A_{j+1} B_j = \mathcal{O}_{i[:, d-1]} \mathcal{C}_{i[:, : -d-1]} \cdots \mathcal{C}_{j+2}^+ \mathcal{C}_{j+2[:, : -d-1]} \mathcal{C}_{j+1}^+ \mathcal{C}_{j+1[:, -d]} \\ &= \mathcal{O}_{i[:, d-1]} \mathcal{C}_{i[:, : -d-1]} I_{[:, : (j+1)d-1]} I_{[:, -d]} \\ &= \mathcal{O}_{i[:, d-1]} \mathcal{C}_{i[:, j d : (j+1)d-1]} = H_{i[:, d-1, j d : (j+1)d-1]} = T_{ij}. \end{aligned} \quad (\text{C.2.8})$$

□

As an example,  $H_i$  can be factorized with SVD as follows:

$$\mathcal{O}_i \mathcal{C}_i = (U_{(r)} D_{(r)}^{1/2}) (D_{(r)}^{1/2} V_{(r)}),$$

where  $U_{(r)} \in \mathbb{R}^{m \times r}$ ,  $D_{(r)} \in \mathbb{R}^{r \times r}$ ,  $V_{(r)} \in \mathbb{R}^{r \times n}$  are the  $r$ -truncated SVD decompositions, and  $r = \text{rank of } H_i \in \mathbb{R}^{m \times n}$ . These factors can then be used to realize a minimal recurrence as outlined above.



## D METHODS

### D.1 COMPUTING ESS

In Section 3 and C.1, we introduced the flattened notation, as it offers a general framework for formulating a wide range of operators and recurrences. As an example, an S5 layer (Smith et al., 2023), which mixes both the channels and sequence simultaneously, can be formulated as  $y = Tu$  (with the operator realization outlined in C.2.1) in the same way an S4 layer can (Gu et al., 2022a), which only mixes the sequence. The difference between these two models lies in the structure of  $T$ : for models that only mix the sequence, such as S4,  $T_{ij}$  is diagonal, whereas for S5, it is not.

Note that since all of the models in our experiments have decoupled channel mixing and sequence mixing (like the S4 layer), we compute the effective state-size independently for each channel using the standard operator formulation  $T \in \mathbb{R}^{\ell \times \ell}$ . This approach is significantly more efficient than computing ESS for the multi-channel (flattened) representation. Furthermore, in the case of attention layers, the computation can be further reduced to only the  $h$  independent heads, as the operator (i.e. the attention matrix) is shared across channels within the same head.

In our experiments, the shape of the unprocessed ESS tensor is given by

$$(\text{batch-size}, \text{layers}, \text{heads or channels}, \text{sequence length} - 1)$$

for a multi-layered model that is processing a batch of sequences. Unless stated otherwise, we compute ESS metrics averaged across all dimensions, with an exception made for softmax attention.

Due to the recurrent realization of softmax attention being constrained to that of the trivial form (Section C.2.3, D.2), the per-channel TSS ( $n_i$ ) of these models depends only on the sequence length  $i$ . In this setting, the average TSS across channels remains constant regardless of the width of the model (even when Q and K expansion factors are applied), and therefore no meaningful variations in TSS are captured by changing the model width. To appropriately capture differences in TSS, we instead sum over the ESS across the channels of each layer and then compute the average over that sum. We denote metrics computed in this manner with a prefix “total”, i.e., “total ESS” and “total TSS” as it captures the total TSS or ESS of a model layer<sup>8</sup>. Additionally, for the analyses presented in Section 4, we average across 8 samples (batch-size), and for the rest, we average across 32 samples.

Regarding the distinction between the entropy and tolerance-based forms of ESS, we note that entropy-ESS is a valuable summary metric because its computation is independent of any specific tolerance value chosen. However, it can potentially be misleading when comparing ESS across sequence indices due to the unequal normalization applied to the singular values. Conversely, when comparing entropy-ESS across different operators, it can be useful as the normalization removes the effect of the norm of the operator. In most of our experiments, we observe consistent trends between entropy-ESS and tolerance-ESS when the metrics are marginalized over the sequence length. Therefore, unless stated otherwise, our figures are presented using the entropy-ESS. In cases where we require ESS comparison across the sequence dimension, we instead plot ESS for multiple tolerance values.

#### D.1.1 PYTORCH IMPLEMENTATION

Below, we provide a PyTorch implementation of various ESS metrics and helper functions that were leveraged in our analyses:

```

1341 import torch
1342
1343 def T2H_i(T, i, d=1):
1344     """
1345     Extract H_i from T.
1346
1347     Args:
1348     - T: Flattened operator with shape [..., d*L, d*L].

```

<sup>8</sup>We note that ESS can capture differences in memory utilization under both metric marginalization approaches.

```

1350     9         - i: Index of H (H_i) to retrieve.
1351     10         - d: Block size for multi-channel flattened operator
1352     representation (default is 1).
1353     11
1354     Returns:
1355     12         - H_i: Submatrix of the operator at index i.
1356     13         """
1357     14         return T[...,d*i:, :d*i]
1358     15
1359     @torch.no_grad()
1360     16 def T2Ss(T, d=1):
1361     17     """
1362     18     Converts an operator into a list of singular values (Ss).
1363     19
1364     Args:
1365     20         - T: Flattened operator with shape [..., d*L, d*L]
1366     21         - d: Block size for multi-channel flattened operator
1367     representation (default is 1).
1368     22
1369     Returns:
1370     23         - Ss: A list of singular values for each sequence index in T.
1371     24         """
1372     25     seqlen = T.size(-2)//d
1373     26     Ss = []
1374     27     for i in range(1, seqlen):
1375     28         H_i = T2H_i(T, i, d)
1376     29         _, S_i, _ = torch.svd(H_i)
1377     30         Ss.append(S_i)
1378     31     return Ss
1379     32
1380     @torch.no_grad()
1381     33 def Ss2ToleranceESS(Ss, tol=1e-4):
1382     34     """
1383     35     Computes the tolerance-ESS from the list of singular values.
1384     36
1385     Args:
1386     37         - Ss: List of singular values.
1387     38         - tol: Tolerance value.
1388     39
1389     Returns:
1390     40         - tolerance-ESS
1391     41         """
1392     42     ranks = []
1393     43     for SV in Ss:
1394     44         rank = torch.sum(SV>=tol, dim = -1)
1395     45         ranks.append(rank)
1396     46     ranks = torch.stack(ranks, dim=-1)
1397     47     return ranks
1398     48
1399     @torch.no_grad()
1400     49 def Ss2EntropyESS(Ss, clip=1e-12):
1401     50     """
1402     51     Computes the entropy-ESS from the list of singular values.
1403     52
1404     Args:
1405     53         - Ss: List of singular values.
1406     54         - clip: clips probabilities below this value avoiding numerical
1407     instabilities when the probabilities are too numerically close to 0.
1408     55
1409     Returns:
1410     56         - entropy-ESS
1411     57         """
1412     58     ranks = []
1413     59     for SV in Ss:
1414     60         p = SV/SV.sum(dim=-1)[..., None]

```

```

1404 71     p = torch.clip(p, clip)
1405 72     H = -torch.sum(p * torch.log(p), dim=-1)
1406 73     rank = torch.exp(H)
1407 74     ranks.append(rank)
1408 75     ranks = torch.stack(ranks, dim=-1)
1409 76     return ranks

```

Example usage (Python-pseudocode):

```

1411 1 >>> out = model(u, output_attentions=True)
1412 2 >>> # T shape: [bs, layers, heads, len, len]
1413 3 >>> T = out.attention_matrix
1414 4 >>> Ss = T2Ss(T) # List of singular values
1415 5 >>> # ESS shape [bs, layers, heads, len-1]
1416 6 >>> ESS = Ss2ToleranceESS(Ss, tol=1e-3)
1417 7 >>> mean_ESS = torch.mean(ESS)

```

We note that calculating the effective rank may cause numerical instability when  $p_i^m$  approaches 0 due to the logarithmic term. This is partially mitigated by clipping the normalized singular values, as shown above.

## D.2 FORMULATION OF THE FEATURIZERS

In this section, we first establish the equivalence between linear attention and state-space models, then proceed with formulating the LIVs discussed in this paper.

**Linear attention and state-space model equivalence.** We begin by demonstrating that linear attention models are state-space models, serving as the foundation for the subsequent formulation of featurizers for other models, such as gated linear attention and weighted linear attention.

A single linear attention head with dimension  $d/h$ , typically formulated as

$$y = qk^T v, \quad (\text{D.2.1})$$

in which  $q, k, v \in \mathbb{R}^{\ell \times d/h}$  are input features. They can be reformulated as recurrences with matrix-valued states  $s_i \in \mathbb{R}^{d/h \times d/h}$  as follows (Katharopoulos et al., 2020):

$$\begin{aligned} s_i &= s_{i-1} + k_i v_i^T \\ y_i &= q_i^T s_i, \end{aligned} \quad (\text{D.2.2})$$

Without loss of generality, applying column-major flattening to the matrix-valued state and treating  $v_i$  as the input  $u_i$ , the recurrence can be formulated like Equation (1), by setting  $A_i = I_{(d/h)^2}$  and:

$$B_{i-1} = \begin{bmatrix} k_i^1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ k_i^{d/h} & 0 & \cdots & 0 \\ 0 & k_i^1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & k_i^{d/h} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & k_i^1 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & k_i^{d/h} \end{bmatrix}, \quad C_i = \begin{bmatrix} q_i^1 & \cdots & q_i^{d/h} & 0 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ 0 & \cdots & 0 & q_i^1 & \cdots & q_i^{d/h} & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \cdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 & \cdots & q_i^1 & \cdots & q_i^{d/h} \end{bmatrix}. \quad (\text{D.2.3})$$

Notice that each individual channel forms a single-input-single-output (SISO) recurrence (like many of the SSM architectures including S4, S6, Mamba2, and more [27; 26; 16; 46; 47]), as there is no mixing across channels. Additionally, each of these SISO recurrences has a state-size of  $d/h$ .

**Formulation of the featurizers.** For the sake of completeness, we additionally characterize the “values” feature in attention-like models, with  $f_u(u)$  as follows:

$$\begin{aligned} s_{i+1} &= A_i s_i + B_i f_u(u_i) \\ y_i &= C_i^T s_i + D_i f_u(u_i). \end{aligned} \quad (\text{D.2.4})$$

The following lists the formulations of the recurrent featurizers based on the results above, along with their per-channel (i.e. average) TSS and total TSS:

- **Linear Attention (LA):**

$$A^k = I, \quad B_{i-1}^k = \text{RoPE}(W_B^k u_i), \quad C_i^k = \text{RoPE}(W_C^k u_i), \quad f_u(u_i) = W_u^k u_i, \quad (\text{D.2.5})$$

where  $W_C^k, W_B^k, W_u^k \in \mathbb{R}^{d/h \times d}$ ;  $d$  and  $h$  represent the number of channels and heads, respectively.  $A$  is a fixed identity matrix. Per-channel TSS is  $n_i^k = d/h$ , and the total TSS is  $d^2/h$ . Each channel  $c \in [d]$  is grouped into heads, where the head index corresponding to the channel is given by  $k = \lfloor ch/d \rfloor$ , and within a head, all corresponding projection matrices ( $W_C^k, W_B^k$ , etc.) are weight tied (shared). Moreover, a rotational positional encoding (RoPE) is by default applied to the  $B$  and  $C$  projections (Su et al., 2023).

- **Gated Linear Attention (GLA):**

$$\begin{aligned} A_{i-1}^k &= \text{diag}(\text{sigmoid}(W_{A_2}^k W_{A_1} u_i)^{1/\beta}), \\ B_{i-1}^k &= W_B^k u_i, \quad C_i^k = W_C^k u_i, \quad f_u(u_i) = W_u^k u_i, \end{aligned} \quad (\text{D.2.6})$$

where besides having projections identical to those in LA,  $W_{A_1} \in \mathbb{R}^{16 \times d}$  and  $W_{A_2}^k \in \mathbb{R}^{d/h \times 16}$ . Like LA, per-channel TSS is  $n_i^k = d/h$ , and the total TSS is  $d^2/h$ . By default,  $\beta$  is set to 16.

- **Weighted Linear Attention (WLA):**

$$\begin{aligned} A^k &= \text{diag}(\text{sigmoid}(\hat{A}^k)^{1/\beta}), \\ B_{i-1}^k &= W_B^k u_i, \quad C_i^k = W_C^k u_i, \quad f_u(u_i) = W_u^k u_i, \end{aligned} \quad (\text{D.2.7})$$

where  $W_C, W_B$ , and  $W_u$  are identical to those in LA, and  $\hat{A}^k \in \mathbb{R}^{d/h}$  is explicitly parameterized and initialized to 0. Like LA, per-channel TSS is  $n_i^k = d/h$ , and the total TSS is  $d^2/h$ .

- **Softmax Attention (SA):**

$$\begin{aligned} \hat{B}_i^k &= \text{RoPE}(W_B^k u_i), \quad \hat{C}_i^k = \text{RoPE}(W_C^k u_i), \\ T^k &= \text{softmax}(\hat{C}^k (\hat{B}^k)^T), \quad f_u(u_i) = W_u^k u_i, \end{aligned} \quad (\text{D.2.8})$$

where  $W_C, W_B$ , and  $W_u$  are identical to those in LA, and  $T$  can be converted into a recurrence using the trivial realization in Equation C.2.5. Therefore, the per-channel TSS is  $i$ , and total TSS is  $id$ . Like LA, a rotational positional encoding (RoPE) is by default applied to the  $\hat{B}$  and  $\hat{C}$  projections (Su et al., 2023).

- **S6 (Gu & Dao, 2024):**

$$\begin{aligned} \Delta^c &= \text{softplus}(W_\Delta^c u_i + b^c), \quad A_{i-1}^c = \text{diag}(\exp(-\hat{A} \Delta^c)), \\ B_{i-1}^c &= \Delta^c W_B u_i, \quad C_i = W_C u_i, \end{aligned} \quad (\text{D.2.9})$$

where  $\hat{A} \in \mathbb{R}^n$  is initialized to  $[1 \ 2 \ \dots \ n]^T$ ,  $c$  is the channel index,  $W_C, W_B \in \mathbb{R}^{n \times d}$ , and  $W_\Delta^c \in \mathbb{R}^{1 \times d}$ . Here, per-channel TSS is  $n$ , and total TSS is  $nd$ . Note that by setting  $n = d/h$ , S6 resembles GLA, with the following exceptions:

- S6 has only one (not  $h$ ) different projection matrices for  $B$  and  $C$ .
- S6 has channel-wise projections for the input-varying discretization applied to  $B$  and  $C$ .

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

- S6 has an explicitly parameterized vector valued  $\hat{A}$ .
- S6 has some minor differences in the non-linearity applied to keep  $0 < A < 1$ . Similar to GLA, S6 also has diagonal  $A$  matrices, whereas in Mamba2 [Dao & Gu \(2024\)](#), the  $A$  matrix is scalar-valued. However, the  $B$  and  $C$  projections in Mamba2 more closely resemble that of GLA.

• **GLA-S6:**

$$A_{i-1}^h = \text{diag}(\exp(-[1/\alpha \quad 2/\alpha \quad \dots \quad n/\alpha]^T \odot \text{softplus}(W_{A_2}^h W_{A_1} u_i))), \quad (\text{D.2.10})$$

$$B_{i-1}^h = W_B^h u_i, \quad C_i^h = W_C^h u_i, \quad f_u(u_i) = W_u^k u_i,$$

GLA-S6 is a combination of S6 and GLA such that the  $B$  and  $C$  projections are identical to that of GLA, while  $A$  is featurized similarly to S6. Namely, it has identical  $W_B$ ,  $W_C$ , and  $W_u$  to those found in LA which means that the per-channel TSS is  $d/h$  and total TSS is  $d^2/h$ . The  $A$  matrix is featurized with the `arange` term ( $[1 \quad 2 \quad \dots \quad n]^T$ ) like S6. We additionally added a normalization hyperparameter  $\alpha$ , which controls the rate at which elements of  $A$  decay to 0.

### 1566 D.3 EMPIRICAL VALIDATION

1567  
1568 Here, we provide details on the task-model sweep presented in Section 4. Table 1 lists the hyper-  
1569 parameters that were exhaustively swept across to generate the task-model space. Note that the  
1570 hyperparameter controlling the task difficulty is task-dependent (for more details, see Poli et al.  
1571 (2024)).

1572 For the MQAR and selective copying tasks, a default vocab size of 8192 (Arora et al., 2023) was  
1573 used for all models. For the compression tasks, the vocab size was varied to modulate task difficulty,  
1574 as shown in Table 1. Any other task settings not specified here are defaulted to those presented in  
1575 Arora et al. (2023). Two important constraints on the tasks from Arora et al. (2023) which we also  
1576 utilize in our experiments are as follows: MQAR task requires that

$$1577 \quad 4 * \text{num kv pairs} \leq \text{seq len}$$

1578 and the selective copying task requires that

$$1579 \quad 2 * \text{num tokens to copy} + 1 < \text{seq len}$$

1580 Any of the task configurations from Table 1 that violate these conditions were not trained. This is  
1581 why the SA plot in Figure 4 has empty spots in the grid.

1582 Finally, we note that all architectures analyzed here consist of 4 layers: 2 sequence mixing layers  
1583 (i.e. one of GLA, LA, WLA or SA) and 2 channel mixing layers (i.e. MLPs).

| 1587 Configuration                           | 1588 Value(s)                        |
|--|--------------------------------------|
| 1589 Tasks                                   | MQAR, selective copying, compression |
| 1590 Num. key-value pairs                    | 8, 16, 32, 64, 128                   |
| 1591 Num. tokens to copy                     | 8, 16, 32, 64, 128                   |
| 1592 Vocab size (compression)                | 8, 16, 32, 64, 128                   |
| 1593 Vocab size (MQAR and selective copying) | 8192                                 |
| 1594 Sequence length                         | 64, 128, 256, 512, 1024, 2048        |
| 1595 Model (featurizer)                      | GLA, LA, WLA, SA                     |
| 1596 Model width                             | 64, 128, 256, 512                    |
| 1597 Number of heads                         | 4, 8                                 |
| 1598 Optimizer                               | AdamW                                |
| 1599 Learning Rate                           | 0.002                                |
| 1600 Weight Decay                            | 0.1                                  |
| 1601 Batch Size                              | 64                                   |
| 1602 Epochs                                  | 70                                   |
| 1603 Steps Per Epoch                         | 2000                                 |
| 1604 Num. Training Samples                   | 128k                                 |
| 1605 Num. Testing Samples                    | 6.4k                                 |

1606  
1607  
1608  
1609 Table 1: Set of hyperparameters for the task-model sweep.

1610 Regarding the post-hoc analysis performed on the sweep, we note the following:

- 1611 • Since the average TSS computed over the channels (which equals  $\frac{\text{model width}}{\text{number of heads}}$  for GLA,  
1612 LA, and WLA) explains more meaningful variation with respect to memory utilization than  
1613 model width and number of heads individually, we consolidate those two dimensions into  
1614 one by analyzing across the average TSS axis. For SA, since average TSS is a function  
1615 of the task rather than model hyperparameters (see Equation C.2.5 and Section D.2), we  
1616 instead compute the sum of TSS over all  $d$  channels, given by the total TSS per layer =  $d*i$ .  
1617 In any cases where the average/total qualifier is not specified, note that we are referring to  
1618 the average ESS or TSS.  
1619

- Since we analyze the recurrent models across the average TSS dimension, we compute average ESS in the plots presented in Section 4.1 in order to compare ESS and TSS as proxies for performance. Similarly, since we analyze the SA models across the total TSS dimension, we compute total ESS for those plots. However, we note that plots for both the average/total ESS and TSS are presented in Section E.1.
- When we marginalize across dimensions, we average across all models in that bucket of task-model space. For example, in Figure 4, for each (TSS, kv) pair, we average over the correlations of all models that correspond to that pair. Note, however, that we never average across tasks (i.e. MQAR, selective copying, compression) or featurizers (i.e. GLA, LA, WLA, SA).
- When we compute cross-model correlations (Figure 2) for SA, we filter out models which have an accuracy  $> 0.95$ . This is done in order to observe meaningful variation as a function of (total ESS)/kv and (total TSS)/kv since many of the SA models obtain an accuracy of 1.
- When we compute within-model correlations (Figure 4) for MQAR, we drop epoch 0 from the computation since we observe a phase at the start of training in which ESS tends to decrease, but accuracy does not change. We elaborate on this phenomenon in Section E.1 and hope to characterize it further in future work.
- Regarding the task-adjusted forms of ESS and TSS which, in the case of MQAR, are computed by normalizing the raw ESS value by the number of kv-pairs in the task, we note that this normalization factor is critical for observing the cross task-model correlations presented in Figure 2. In particular, in Figure 9, we find that correlations across the task-model space break down when examining the unnormalized ESS. This points to the higher-level notion that ESS is expected to scale with the memory demands of the task.
- We interpret the state utilization of a model, which is given by ESS/TSS, as a proxy for what portion of the memory capacity of the network is realized in practice. By definition, state utilization takes on values ranging continuously from 0 to 1. Recall that a state utilization near 1 is indicative of state saturation.
- While for most of the ESS analysis conducted on the sweep we use the entropy-ESS, we note that for the state utilization plot presented in Figure 5, we use the tolerance-ESS with a tolerance level set at  $1e-3$ . We do this because we find that entropy-ESS fails to capture the state collapse phenomenon. This is because state collapse is primarily dictated by the magnitude of the singular values, as opposed to the relative decay rate of the entire spectrum. In particular, if all of the singular values are close to 0, the layer is likely failing to learn an expressive state, resulting in poor performance. Due to the normalization applied to the spectrum, the entropy-ESS metric may potentially present this state as having a high effective rank; however, in practice, we know that this is a misrepresentation of the true dynamics. Tolerance-ESS, in contrast, appropriately captures the dynamics of the state with respect to the norm of the operator. Because of this, whenever we analyze ESS as it pertains to state collapse (e.g. Figure 8a), we present the tolerance-ESS instead.

#### D.4 ESS-INFORMED FEATURIZER SELECTION AND INITIALIZATION SCHEME

| Configuration                          | Value |
|--|-------|
| Model width                            | 128   |
| Num. heads                             | 8     |
| arange Norm. ( $\alpha$ ) <sup>a</sup> | 1000  |
| Logit Norm. ( $\beta$ )                | 16    |
| $K$ -expansion <sup>b</sup>            | 1     |

Table 2: Default GLA hyperparameters.

| Configuration             | Value |
|---------------------------|-------|
| Model width               | 128   |
| State expansion (d_state) | 16    |

Table 3: Default S6 hyperparameters.

<sup>a</sup>For GLA-S6.<sup>b</sup> $K$ -expansion is used to vary TSS in the featurizer experiments.

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

| Configuration         | Value              |
|-----------------------|--------------------|
| Sequence length       | 2048               |
| Num. KV Pairs         | 128                |
| KV Dist. Const.       | 0.1                |
| Optimizer             | AdamW <sup>a</sup> |
| Learning Rate         | 0.002              |
| Weight Decay          | 0.1                |
| Batch Size            | 64                 |
| Epochs                | 70                 |
| Steps Per Epoch       | 2000               |
| Num. Training Samples | 128k               |
| Num. Testing Samples  | 6.4k               |
| Vocabulary Size       | 8192               |

Table 4: Default MQAR task settings employed throughout the featurizer and initialization experiments in Section 5.1.

---

<sup>a</sup>Loshchilov & Hutter (2019)



## D.5 ESS-INFORMED REGULARIZATION

We use the following MQAR configuration for the regularization experiments presented in Section A.

| Configuration         | Value              |
|-----------------------|--------------------|
| Sequence length       | 4096               |
| Num. KV Pairs         | 128                |
| KV Dist. Const.       | 0.1                |
| Optimizer             | AdamW <sup>a</sup> |
| Learning Rate         | 0.002              |
| Weight Decay          | 0.1                |
| Batch Size            | 64                 |
| Epochs                | 70                 |
| Steps Per Epoch       | 2000               |
| Num. Training Samples | 128k               |
| Num. Testing Samples  | 6.4k               |
| Vocabulary Size       | 8192               |
| Model width           | 128                |
| Num. heads            | 8                  |

Table 5: MQAR task settings and model hyperparameters employed throughout the mid-training experiments in Section A.

<sup>a</sup>Loshchilov & Hutter (2019)

Regarding the regularization scheme itself, since we examine models with two sequence mixing layers, we explore the following strategies: regularizing both layers, only regularizing the first layer and only regularizing the second layer. Empirically, we find that only regularizing the second layer performs the best and is thus the result presented in Figure 8b. We elaborate on why this is the most successful strategy in Section E.3.

## D.6 ESS-INFORMED MODEL-ORDER REDUCTION

The teacher models used in the distillation experiments are 2-layer GLA models (Yang et al., 2024a) with dimension = 128 and TSS = 256 (num\_heads = 8 and expand\_k = 16). We checkpointed the models every 10 epochs while training on MQAR across different task difficulties. The task ranges are given as follows:

- Sequence length: [512, 1024, 2048]
- Number of Key-Value Pairs: [64, 128]

Other settings follow the defaults shown in Table 4. For each task difficulty pair, we repeated the training run with three different seeds. For each teacher model checkpoint, both layers were distilled independently with student models of different state-sizes (16, 32, 64, and 128). Distillation settings are shown in Table 6.

The ESS metric in Figures 6d, 30, and 31 was computed by taking the minimum across input samples and model channels, evaluated at the mid-point of the sequence ( $\ell/2$ ). Using the mid-point of the sequence as a summary statistic was done to save compute. The midpoint in particular was chosen as it is the point in the sequence at which  $H_i$  has the greatest dimensions, retaining the largest amount of information from the original operator. Other approaches such as taking the maximum or average across the sequence also show similar trends, but we found taking the minimum to be the clearest.

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

| Configuration               | Value |
|-----------------------------|-------|
| Optimizer                   | AdamW |
| Batch Size                  | 1     |
| Learning Rate               | 0.001 |
| Weight Decay                | 0.0   |
| Training Steps (Operator)   | 800   |
| Dropout (Operator)          | 0.2   |
| Training Steps (Activation) | 3200  |
| Dropout (Activation)        | 0.2   |

Table 6: Distillation settings used for the results presented in Section 5.2.

### D.7 ESS ANALYSIS FOR HYBRID NETWORKS

In our ESS analysis applied to hybrid networks, we restrict our scope to GLA-SA hybrids. In particular, we explore the following two settings:

- 8-layer hybrid networks in which 4 layers are sequence mixers (i.e. one of GLA or SA) and 4 layers are channel mixers (i.e. MLPs). We exhaust all possible hybrid networks (of which there are 16) and perform post-training, per-layer ESS analysis on the networks. We train these hybrid models on MQAR with task-model settings given below in Table 7.
- 16-layer hybrid networks in which 8 layers are sequence mixers (i.e. one of GLA or SA) and 8 layers are channel mixers (i.e. MLPs). Here, we explore all combinations of hybrid networks that follow the Jamba hybridization policy (Lieber et al., 2024) and perform post-training, per-layer ESS analysis on the networks. We train these hybrid models on MQAR with task-model settings given below in Table 8.

| Configuration         | Value              |
|-----------------------|--------------------|
| Sequence length       | 2048               |
| Num. KV Pairs         | 512                |
| KV Dist. Const.       | 0.1                |
| Optimizer             | AdamW <sup>a</sup> |
| Learning Rate         | 0.002              |
| Weight Decay          | 0.1                |
| Batch Size            | 64                 |
| Epochs                | 70                 |
| Steps Per Epoch       | 2000               |
| Num. Training Samples | 128k               |
| Num. Testing Samples  | 6.4k               |
| Vocabulary Size       | 8192               |
| Model width           | 64                 |
| Num. heads            | 4                  |

Table 7: Default MQAR task settings employed throughout the hybridization experiments conducted in the first setting described above.

<sup>a</sup>Loshchilov & Hutter (2019)

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

| Configuration         | Value              |
|-----------------------|--------------------|
| Sequence length       | 4096               |
| Num. KV Pairs         | 1024               |
| KV Dist. Const.       | 0.1                |
| Optimizer             | AdamW <sup>a</sup> |
| Learning Rate         | 0.002              |
| Weight Decay          | 0.1                |
| Batch Size            | 64                 |
| Epochs                | 70                 |
| Steps Per Epoch       | 2000               |
| Num. Training Samples | 128k               |
| Num. Testing Samples  | 6.4k               |
| Vocabulary Size       | 8192               |
| Model width           | 16                 |
| Num. heads            | 2                  |

Table 8: Default MQAR task settings employed throughout the hybridization experiments conducted in the second setting described above.

<sup>a</sup>Loshchilov & Hutter (2019)

Results for these experiments can be found in Section E.4.2.

## D.8 STATE MODULATION OF LARGE LANGUAGE MODELS

**State modulation of open-weight models.** The following randomly generated sentences were used to study the effects of separator tokens on state modulation in open-weights pre-trained language models.

*<bos>Mangoes are rich in vitamin C and can be blended into a refreshing smoothie<sep> Giraffes are the tallest mammals on Earth due to their long necks and legs<sep> She collects vintage typewriters from the 1940s<sep> Jupiter’s Great Red Spot is a giant storm that has been raging for hundreds of years<sep>*

**State modulation on custom-trained 1B models.** For our custom-trained 1B language models, we used longer sentences, as state modulation patterns were less discernible with shorter sequences. A collection of randomly generated sentences is shown below:

*<bos>The deep blue ocean, teeming with an extraordinary array of marine life, from the smallest plankton to the largest whales, stretches out infinitely towards the horizon, a vast and mysterious expanse that has captivated the imaginations of explorers, scientists, and poets for centuries, hiding within its depths secrets yet to be discovered and stories yet to be told<sep> In a bustling city where skyscrapers tower over narrow streets filled with the constant hum of cars and the chatter of pedestrians, a small café, nestled between two imposing buildings, offers a quiet refuge for those seeking a moment of peace, with the comforting aroma of freshly brewed coffee and the soft sound of jazz music playing in the background, creating a cozy ambiance that feels like a world away from the urban chaos outside<sep> The ancient oak tree, with its gnarled branches stretching wide and its thick, sturdy trunk standing firm against the passage of time, has witnessed generations of families grow, seasons change, and countless stories unfold beneath its expansive canopy, becoming a silent guardian of the park, offering shade to those who seek solace and a sense of continuity in a rapidly changing world<sep>*

We note that the specific sentences and their order are not crucial to this analysis. Similar patterns have emerged with various sentence arrangements, provided the sentences are sufficiently long.

Training settings are outlined in Table 9.

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

| Configuration       | Value |
|---------------------|-------|
| Batch Size          | 16    |
| Max Sequence Length | 32k   |
| Training Steps      | 160k  |
| Optimizer           | AdamW |
| Learning Rate       | 0.001 |
| Weight Decay        | 0.1   |
| Num. Layers         | 24    |
| Dimension           | 2048  |

Table 9: 1B LLM settings.

The perplexity scores shown in Figure 7b were computed on 16k randomly sampled sequences over the FineWeb (Penedo et al., 2024) dataset. The raw perplexity samples were smoothed via a kernel density estimation method.

## E EXTENDED EXPERIMENTAL RESULTS

### E.1 EMPIRICAL VALIDATION

In this section, we provide additional results and commentary from the sweep detailed in Section D.3 that were not presented in the main portion of the paper. One thing to note is that most of the ESS results presented in Section 4 were computed using the entropy-ESS. However, we also computed ESS using the tolerance-based approach to affirm that both forms of ESS showcase similar trends. In particular, we examined tolerances of 1e-1, 1e-3 and 1e-5. Since we observe similar trends across tolerances, we provide plots for a tolerance of 1e-3 below and omit the others for the sake of brevity.

#### E.1.1 STATE COLLAPSE CONTINUED

Here, we continue our discussion on the state collapse phenomenon presented in Section 4.2. In particular, while we assert that state collapse is observable across all TSS in the high kv bucket for GLA/WLA, Figure 5 shows that accuracy differences between LA and GLA/WLA are only evident in the high TSS/high kv bucket of the task-model space. This is because state saturation is acting as a confounder, worsening performance in LA (see Figure 5 when TSS is 8). Therefore, although state collapse in GLA/WLA does not result in worse performance than LA in this specific task-model setting, it remains an issue even for models with smaller states when trained on sufficiently difficult tasks. This is the motivation behind the task-model setting explored in Section A.

#### E.1.2 ENTROPY-ESS MQAR RESULTS CONTINUED

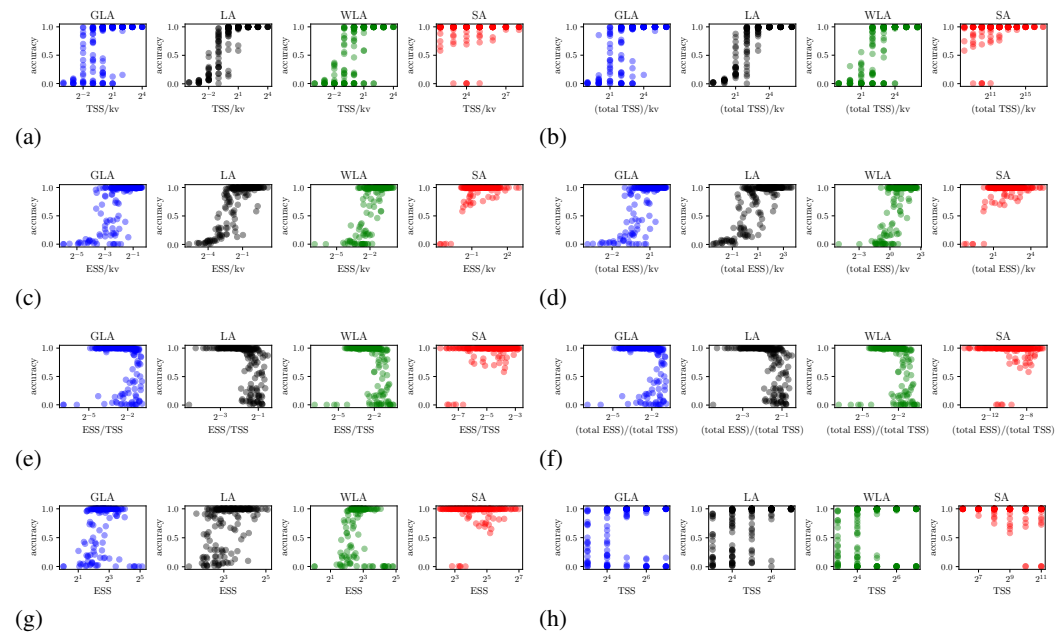


Figure 9: (a) TSS/kv vs accuracy across featureizers. This demonstrates that TSS/kv (i.e. memory capacity) is a worse proxy for model performance than ESS/kv as discussed in Section 4. (b) (total TSS)/kv vs accuracy across featureizers. This demonstrates that (total TSS)/kv is a worse proxy for model performance than (total ESS)/kv. (c) ESS/kv vs accuracy across featureizers. (d) (total ESS)/kv vs accuracy across featureizers. (e) ESS/TSS (i.e. state utilization) vs accuracy across featureizers. We note that models that saturate their state tend to perform worse on the task, which is evidence of the state saturation phenomenon discussed in Section 4.2. The models that do not saturate their state but still perform poorly are the models that undergo state collapse. (f) (total ESS)/(total TSS) vs accuracy across featureizers. (g) ESS vs accuracy across featureizers. Note that without normalizing by kv (i.e. the task memory), the correlation with accuracy breaks down substantially. (h) TSS vs accuracy across featureizers.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

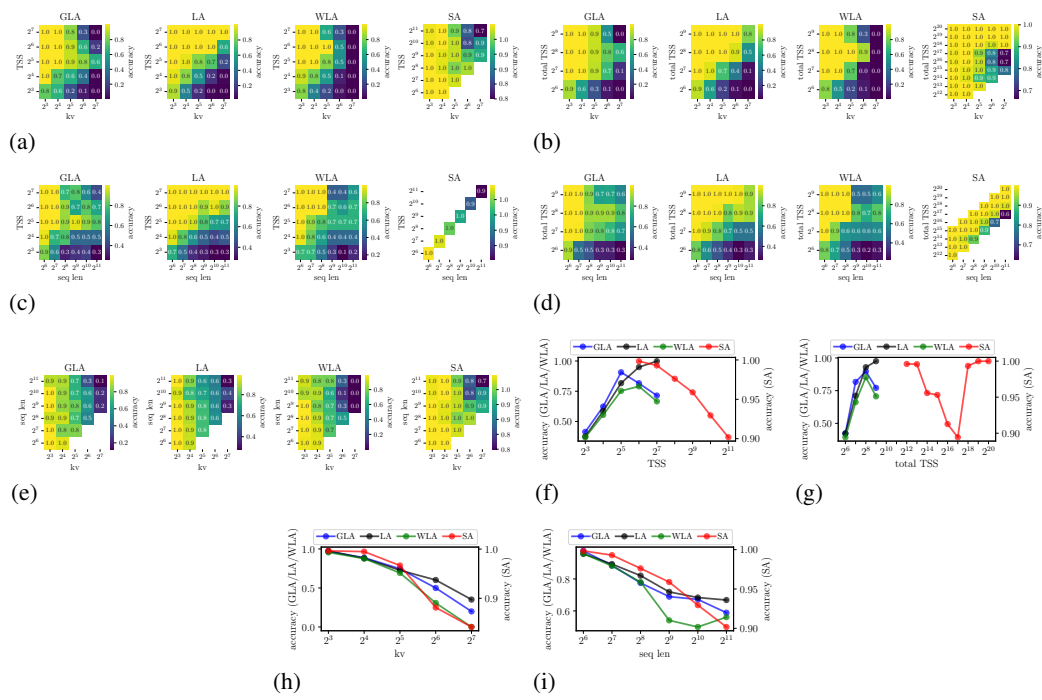


Figure 10: MQAR accuracies marginalized across different dimensions.

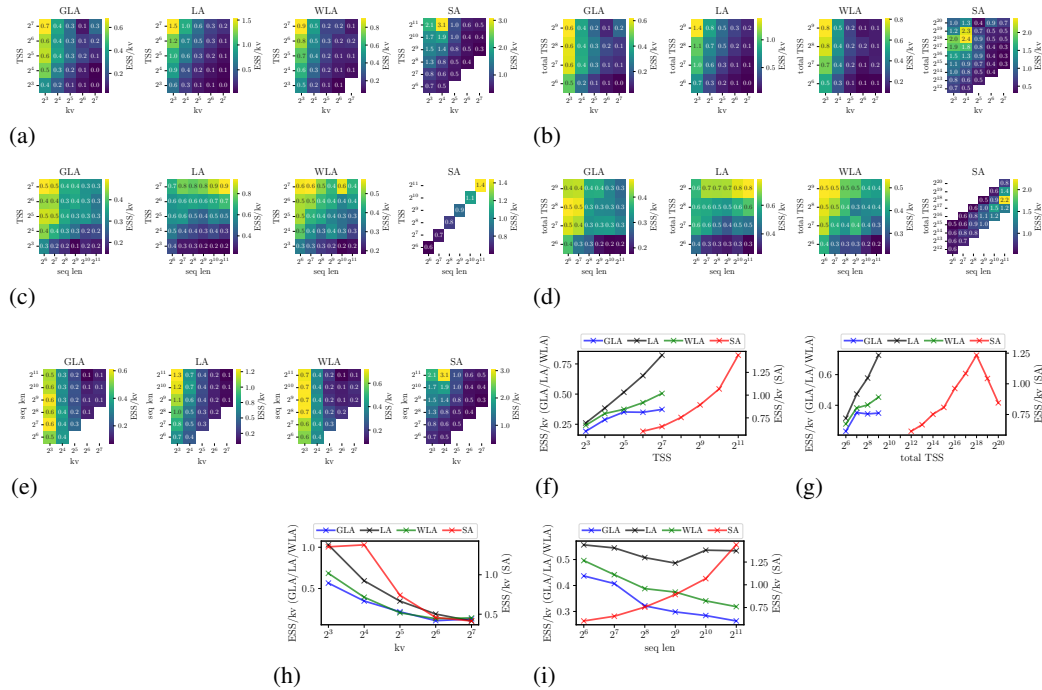


Figure 11: MQAR ESS/kv marginalized across different dimensions.

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

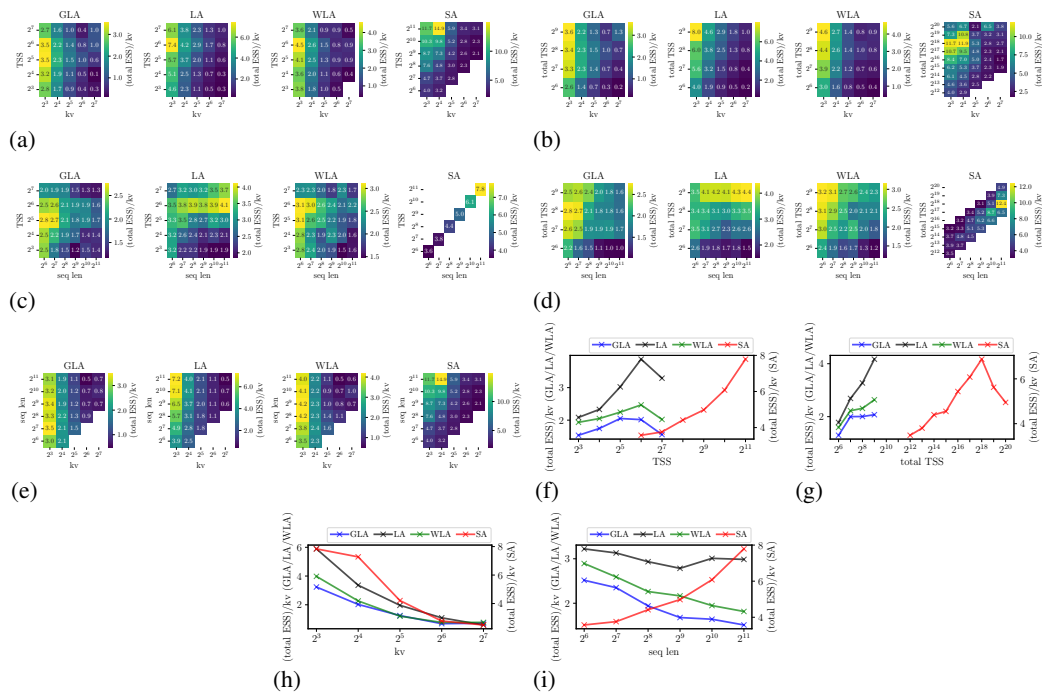


Figure 12: MQAR (total ESS)/kv marginalized across different dimensions.

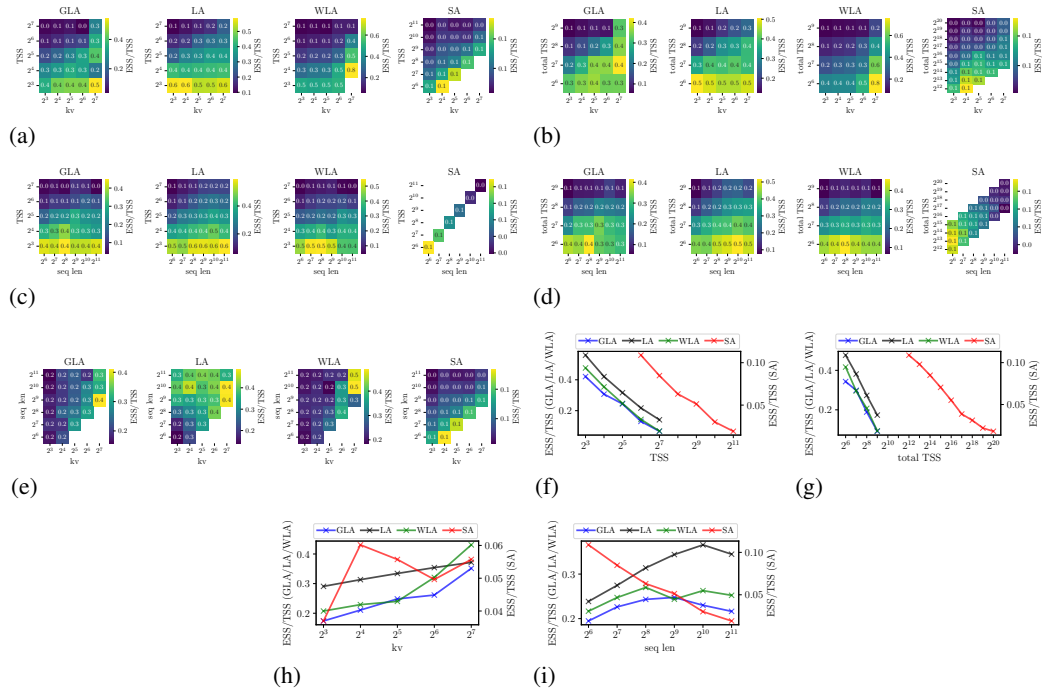


Figure 13: MQAR ESS/TSS marginalized across different dimensions.

2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159

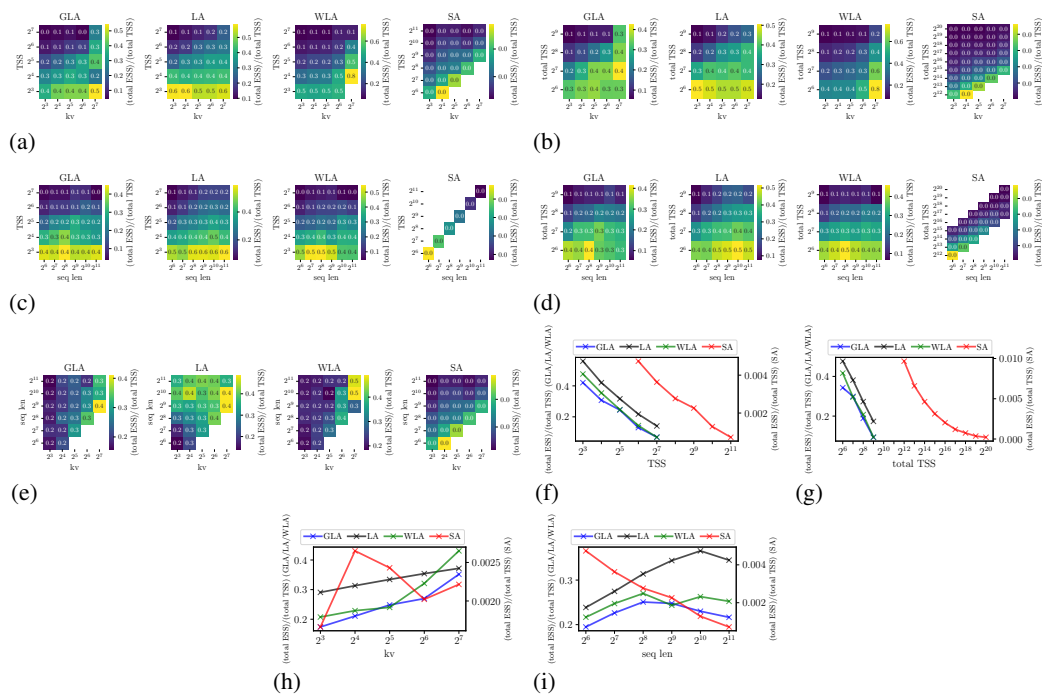


Figure 14: MQAR (total ESS)/(total TSS) marginalized across different dimensions.

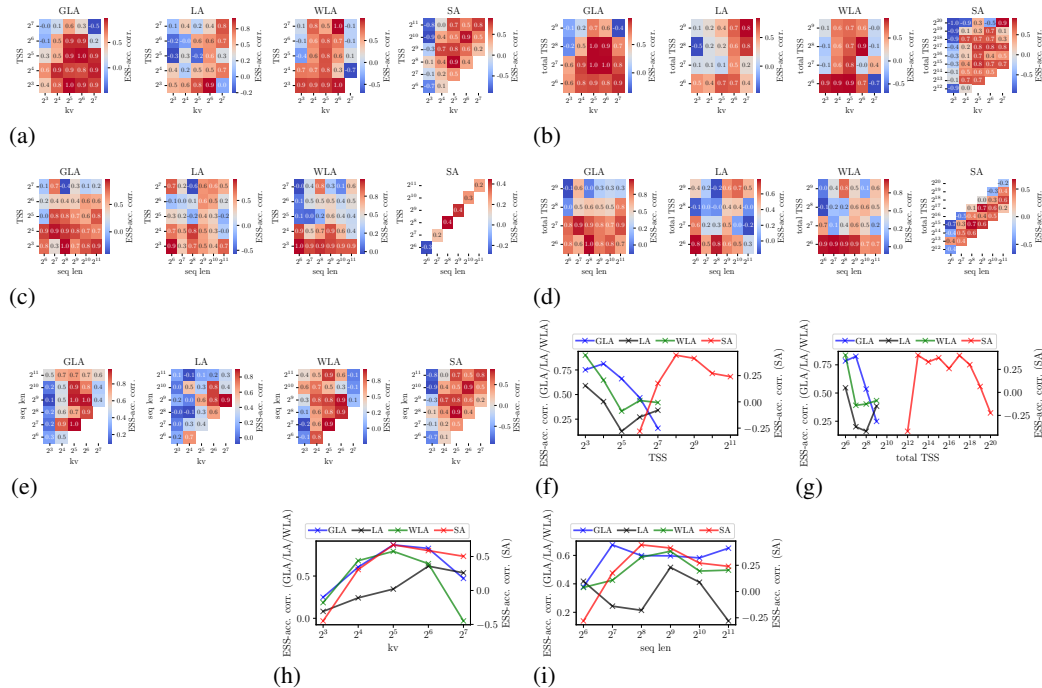


Figure 15: MQAR ESS-accuracy correlations computed over training marginalized across different dimensions.

### E.1.3 TOLERANCE-ESS MQAR RESULTS

Below are plots from the MQAR sweep using tolerance-ESS (tol=1e-3) instead of entropy-ESS. We note that all of the prevailing trends remain the same.



2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

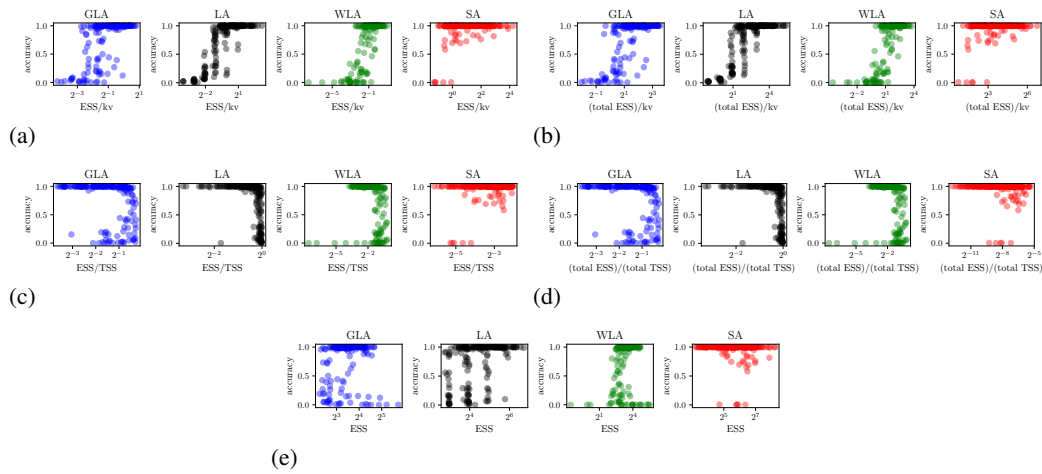


Figure 16: Accuracy vs various forms of tolerance-ESS across task-model space. Plots are entirely analogous to those shown in Figure 9.

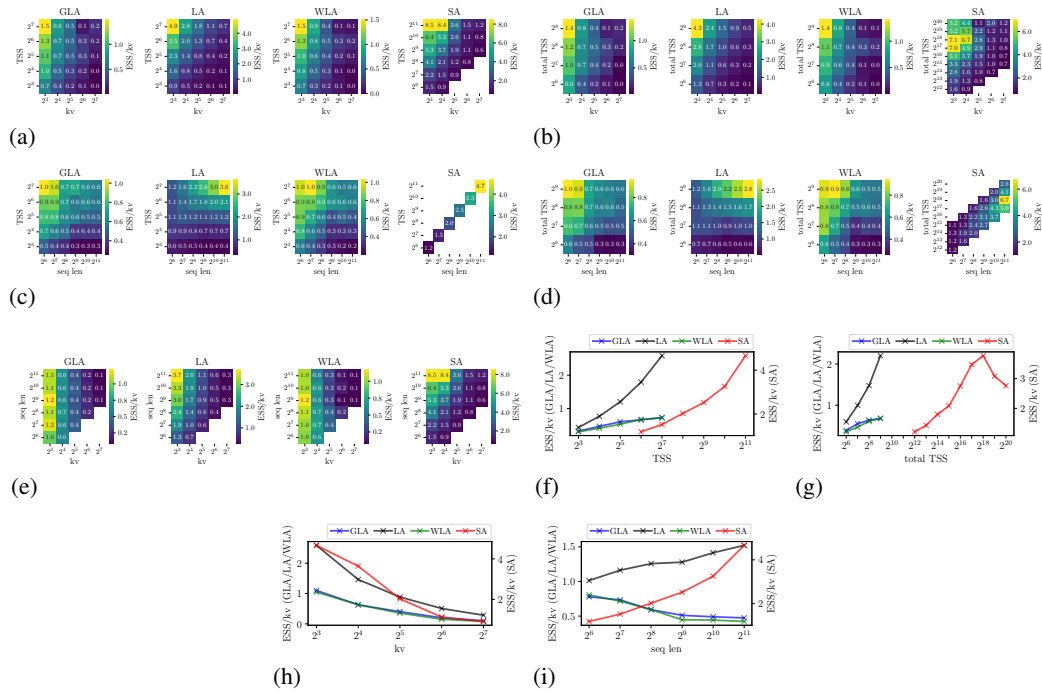


Figure 17: MQAR ESS/kv marginalized across different dimensions.

2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267

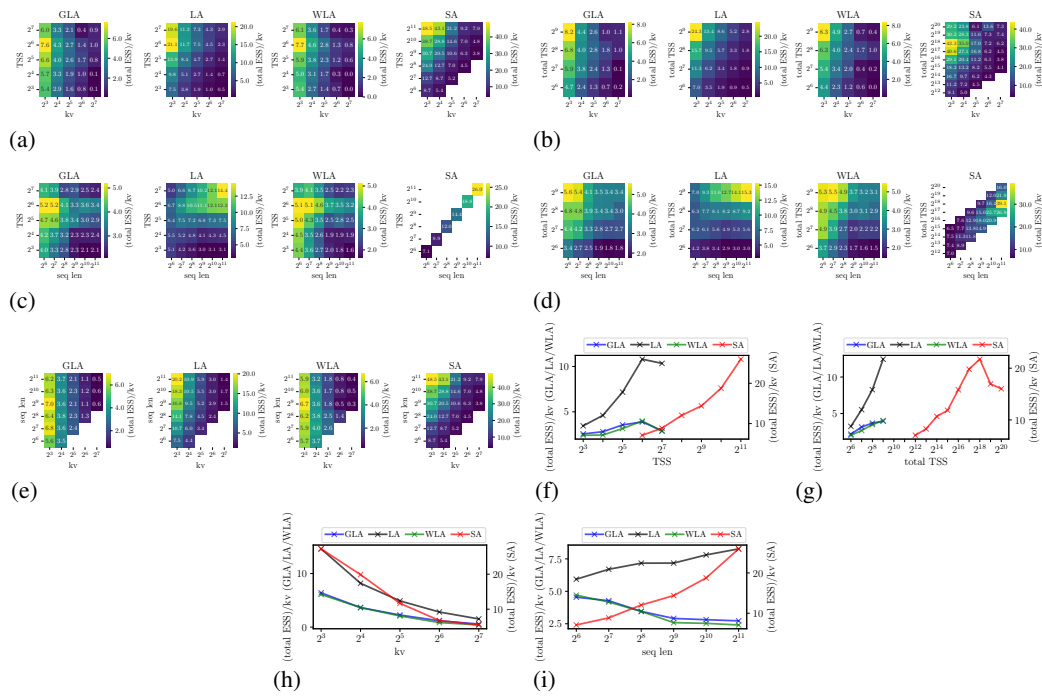


Figure 18: MQAR (total ESS)/kv marginalized across different dimensions.

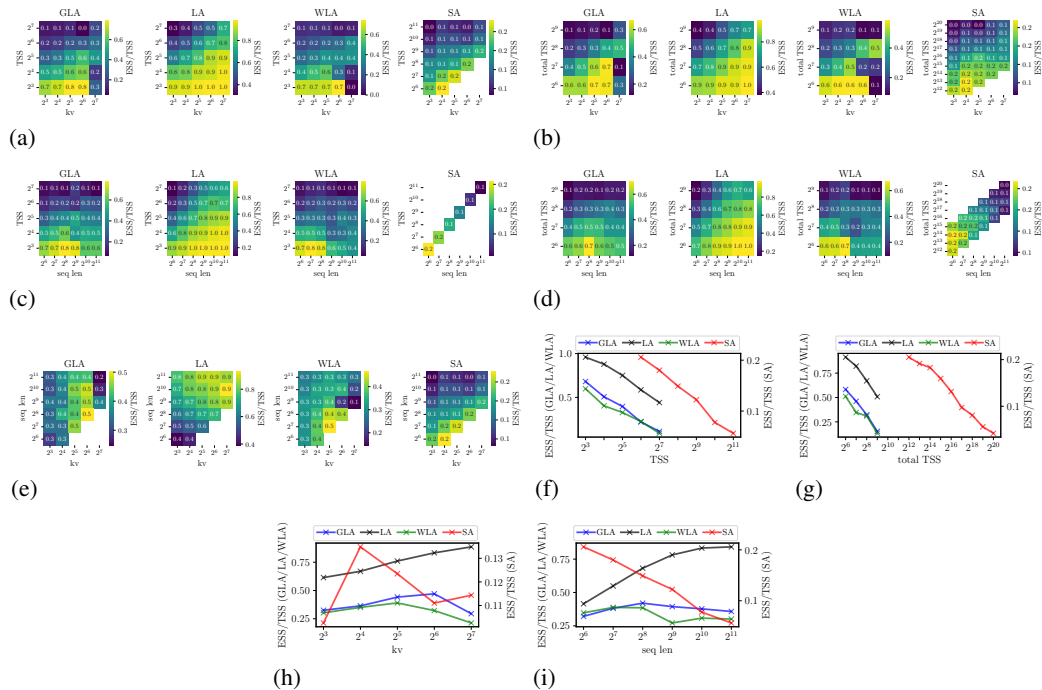


Figure 19: MQAR ESS/TSS marginalized across different dimensions.

2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321

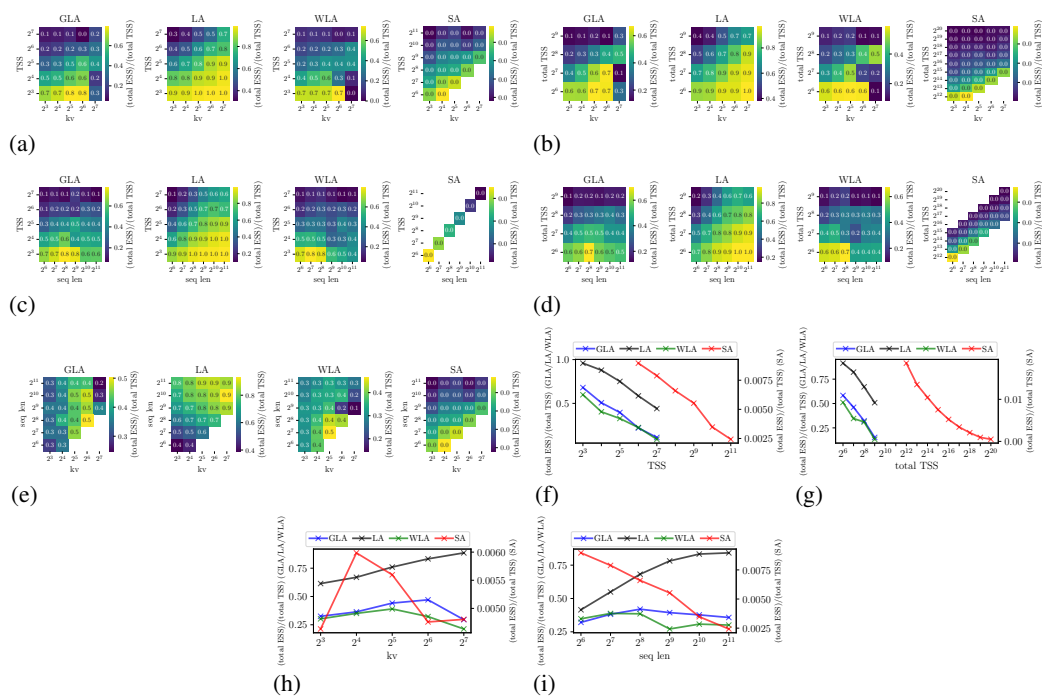


Figure 20: MQAR (total ESS)/(total TSS) marginalized across different dimensions.

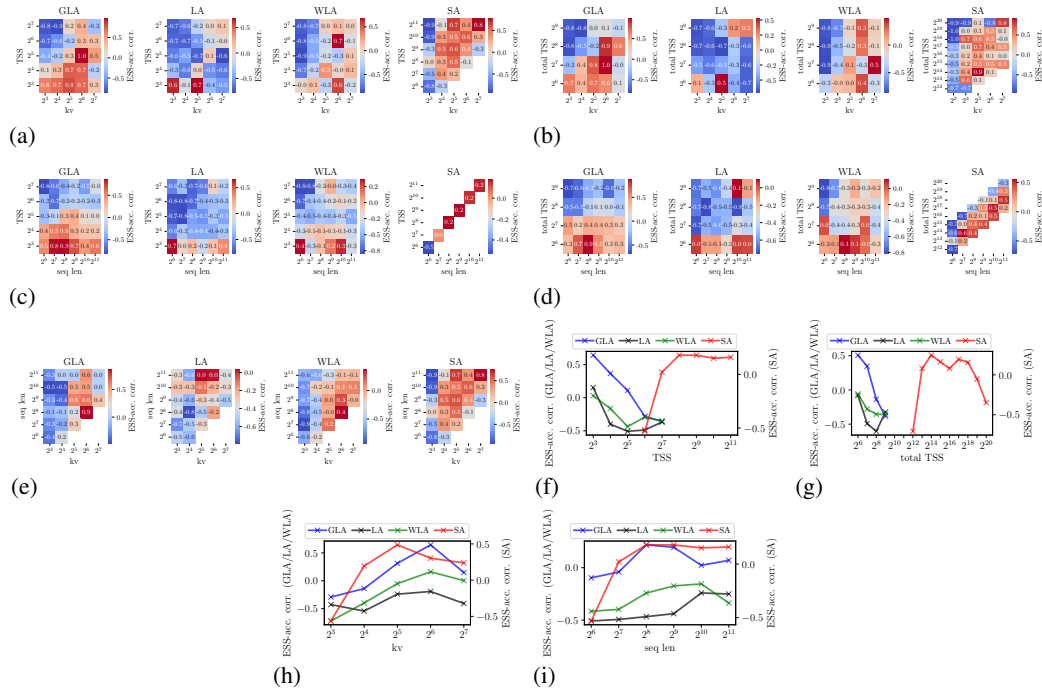


Figure 21: MQAR ESS-accuracy correlations computed over training marginalized across different dimensions.

### E.1.4 SELECTIVE COPYING AND COMPRESSION RESULTS

Below, we present results for the selective copying and compression tasks, analogous to the ones presented in Section 4 on MQAR.

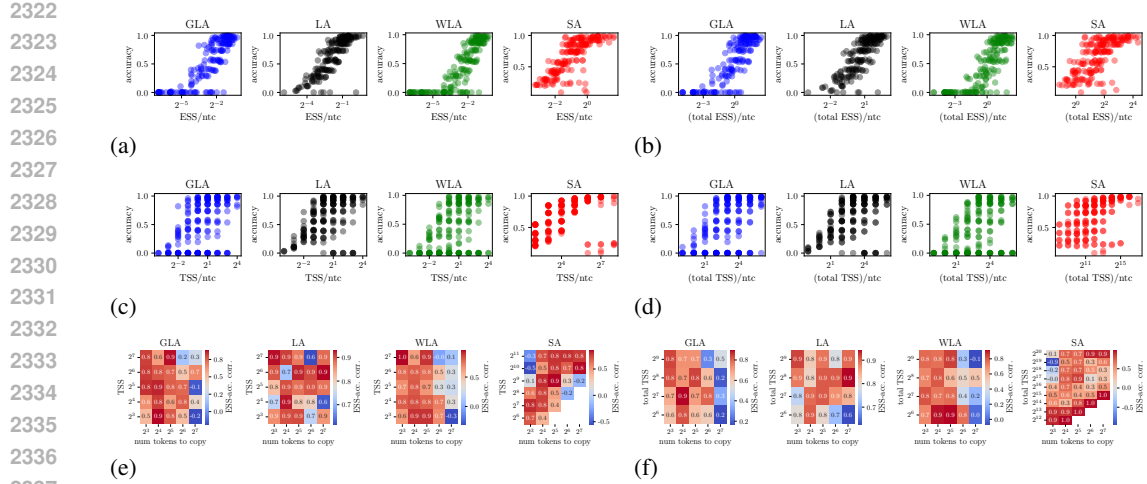


Figure 22: Selective copying results. Note that ESS here refers to entropy-ESS and we abbreviate num. tokens to copy as ntc in plots above. (a) ESS/ntc vs accuracy across featureizers. (b) (total ESS)/ntc vs accuracy across featureizers. (c) TSS/ntc vs accuracy across featureizers. (d) (total TSS)/ntc vs accuracy across featureizers. (e) ESS-accuracy correlation computed over the course of training in (TSS, kv) buckets. (f) ESS-accuracy correlation computed over the course of training in (total TSS, kv) buckets.

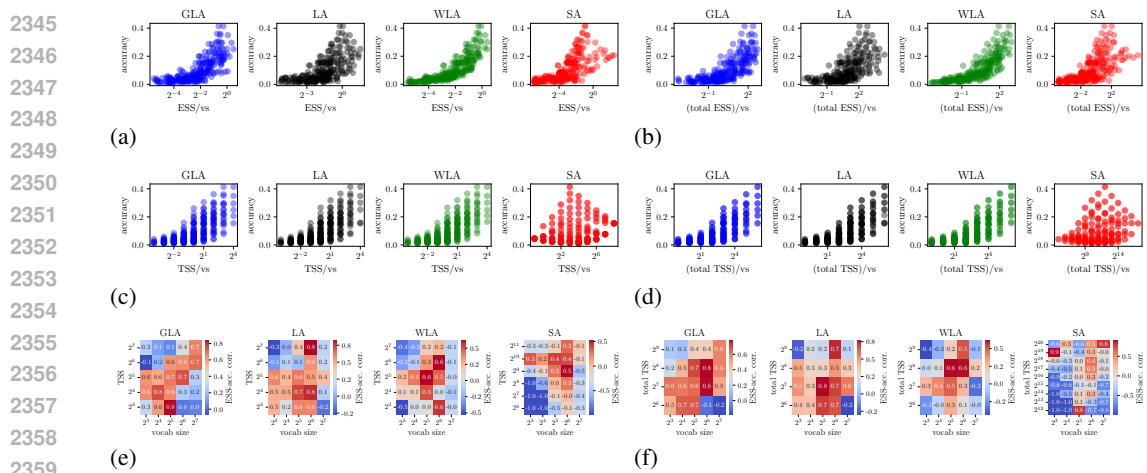


Figure 23: Compression results. Note that ESS here refers to entropy-ESS and we abbreviate vocab size as vs in plots above. (a) ESS/vs vs accuracy across featureizers. (b) (total ESS)/vs vs accuracy across featureizers. (c) TSS/vs vs accuracy across featureizers. (d) (total TSS)/vs vs accuracy across featureizers. (e) ESS-accuracy correlation computed over the course of training in (TSS, kv) buckets. (f) ESS-accuracy correlation computed over the course of training in (total TSS, kv) buckets.

We note that with respect to the cross task-model trends, we find that in both selective copying and compression, task-adjusted ESS is a better proxy for model performance than task-adjusted TSS (Figures 22a, 22c, 23a, 23c). This is substantial as it demonstrates the utility of the ESS metric beyond just MQAR.

Regarding within task-model trends, we observe similar patterns for selective copying as those seen in MQAR (Figure 22e), with one notable distinction. Namely, ESS and accuracy are positively correlated across a larger portion of the task-model space in selective copying than in MQAR. For compression, however, the within task-model trends look a bit different from what we observe in selective copying and MQAR (Figure 23e). One potential reason for this is that the compression task is significantly more difficult than the MQAR and selective copying tasks (as noted by the

lower accuracies in Figure 23a), leading to more instabilities over the course of training. But in any case, this does highlight the fact that the strength of ESS as a proxy for model performance changes as a function of the task. The precise nature of this relationship is something we hope to explore in future work.

### E.1.5 ESS TRAINING DYNAMICS IN MQAR

As mentioned in Section D.3, we observe a phase at the start of training in MQAR in which ESS tends to decrease. This is shown in Figure 24 in which we select an arbitrary task-model configuration from the sweep and plot its ESS and accuracy over the course of training on a per-featurizer basis.

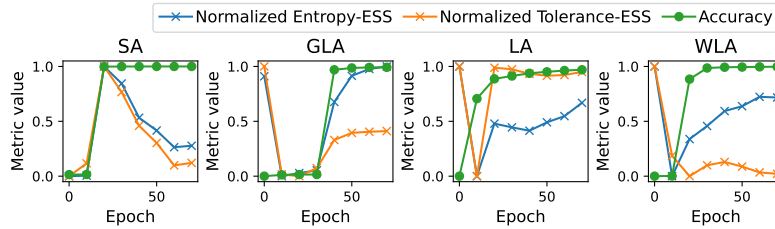
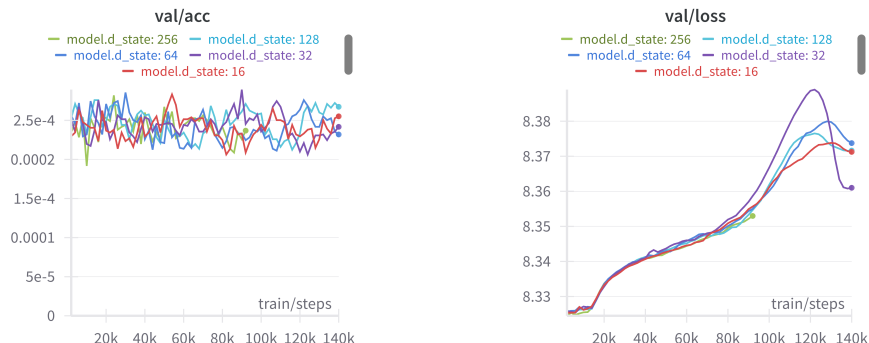


Figure 24: Training dynamics of ESS in select models ( $d_{\text{model}}=256$ ,  $\text{heads}=8$ ) trained on MQAR ( $\text{seqlen}=2048$ ,  $\text{kv}=64$ ). We min-max normalize the ESS curves over the course of training to emphasize the shape of the curve as opposed to its magnitude. Note that the tolerance-ESS shown here is computed using a tolerance of  $1e-3$ .

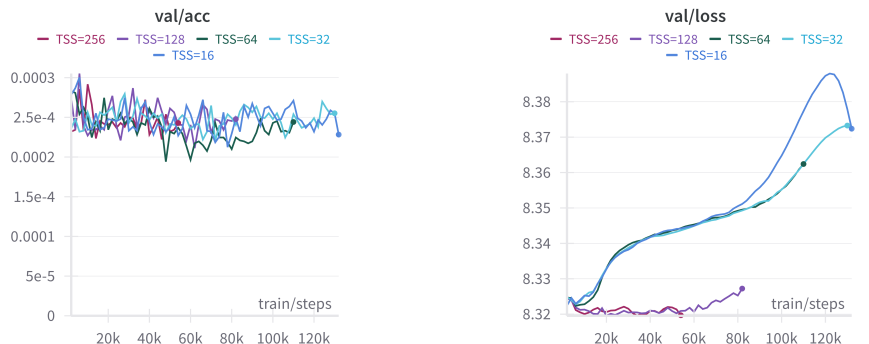
We find that at the start of training (i.e. in between epochs 0 and 10), even if the accuracy is not evolving, the ESS is. In particular, in the recurrent frameworks (GLA, LA and WLA), we note a sharp decrease in the ESS before it begins to rise later in training (and along with it the model accuracy). In contrast, in SA we observe the opposite: a sharp increase at the start of training followed by a steady decrease (even after it has solved the task). This points to a level of nuance in the training dynamics of MQAR ESS that we have yet to characterize and is something we hope to explore in future work.

E.2 INITIALIZATION-PHASE ANALYSIS

2430  
2431  
2432  
2433  
2434  
2435  
2436  
2437  
2438  
2439  
2440  
2441  
2442  
2443  
2444  
2445  
2446  
2447  
2448  
2449  
2450  
2451  
2452  
2453  
2454  
2455  
2456  
2457  
2458  
2459  
2460  
2461  
2462  
2463  
2464  
2465  
2466  
2467  
2468  
2469  
2470  
2471  
2472  
2473  
2474  
2475  
2476  
2477  
2478  
2479  
2480  
2481  
2482  
2483



(a) Validation accuracy of S6 (b) Validation loss of S6



(c) Validation accuracy of GLA-S6 (d) Validation loss of GLA-S6

Figure 25: Loss curves of S6 and GLA-S6 showing that the models are unable to improve beyond random guessing on MQAR, across various state-sizes.

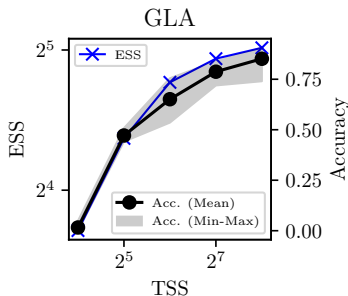


Figure 26: ESS and MQAR accuracy as a function of TSS on a custom task regime (sequence length = 1024, num. kv pairs = 256). This figure illustrates a strong correlation between MQAR accuracy, ESS and TSS.

E.3 MID-TRAINING ANALYSIS

First, we provide some additional commentary on the ESS-based regularization results discussed in Section A. Recall we showed that decaying the  $A$  matrices in GLA and WLA towards the identity matrix enables these models to outperform LA in the state collapse regime. Our intuition for this result is that by ameliorating state collapse, GLA and WLA can better leverage their increased expressivity, which stems from their learnable  $A$  matrices – a feature absent from LA.

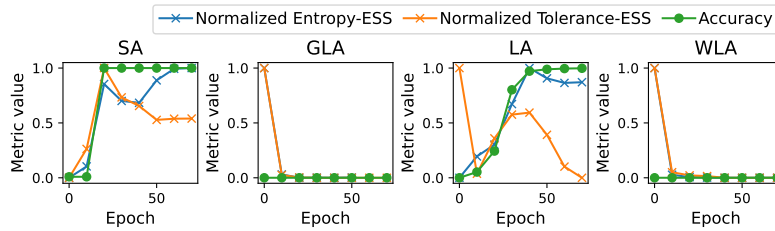


Figure 27: An example of the training dynamics of ESS in select models (dmodel=512, heads=4) trained on MQAR (seqlen=2048, kv=128) that undergo state collapse (i.e. GLA and WLA). We min-max normalize the ESS curves over the course of training to emphasize the shape of the curve as opposed to its magnitude. Note that the tolerance-ESS shown here is computed using a tolerance of  $1e-3$ .

Next, as mentioned in Section D.5, we provide some intuition behind the efficacy of regularizing only the second layer of the network as opposed to the first or both layers.

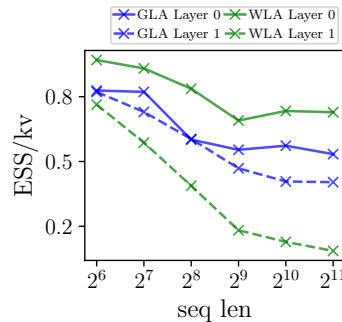


Figure 28: Per-layer ESS/kv as a function of MQAR sequence length for the GLA and WLA featureizers. ESS shown here is computed using a tolerance of  $1e-3$ . Layers are 0-indexed.

Using 0-indexing for the layers, Figure 28 shows that layer 1 realizes a lower ESS/kv than layer 0, particularly in the case of WLA. This suggests that layer 1 contributes disproportionately to the observed state collapse (Figure 27); consequently, it makes sense that layer 1 would need to be regularized more heavily. Now, this begs the question as to why only regularizing the second layer leads to better performance than regularizing both layers (results of which were not shown). We have two possible hypotheses for this outcome. First, introducing regularization terms for both layers may complicate optimization by creating potentially conflicting objectives. Second, excessive decay of the  $A$  matrices towards the identity matrix may cause the model to revert to the LA regime, which – as shown in Figure 8b – performs worse than GLA and WLA (when sufficiently regularized). Nonetheless, we hope to further explore this intuition and investigate other ESS-based forms of regularization in future work.

E.4 POST-TRAINING ANALYSIS

E.4.1 MODEL-ORDER REDUCTION

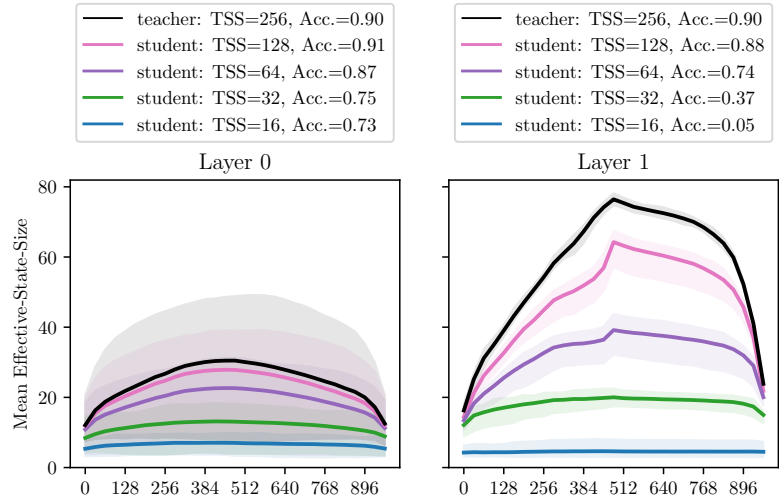


Figure 29: This figure compares MQAR accuracy and ESS across reduction scales for layers 0 and 1. The lower ESS in layer 0 of the teacher model leads to better downstream performance after distillation compared to distilling layer 1.

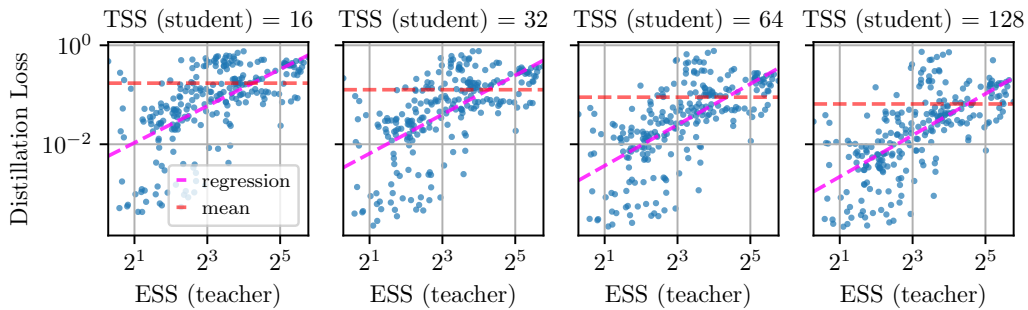


Figure 30: Correlation between ESS and distillation loss across multiple student TSSs (reduction ratios). The original teacher models have a TSS of 256.

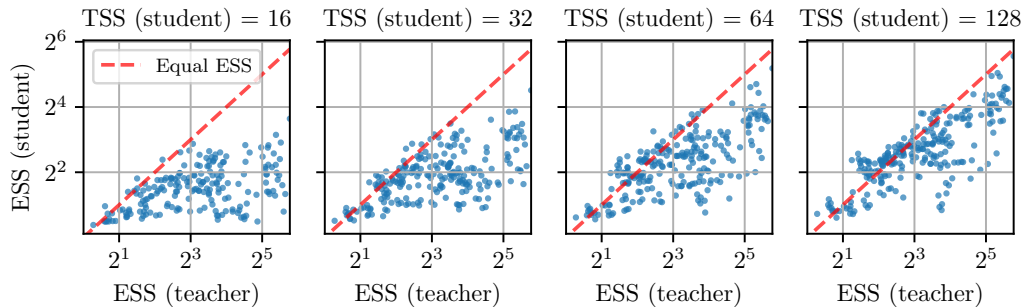


Figure 31: Teacher ESS vs distilled student ESS. As expected, we observe a clear trend: an increase in the student TSS results in the student’s ESS more closely matching the teacher’s ESS. Plots like these can help provide additional context during the distillation process.



E.4.2 HYBRIDIZATION

2592  
2593  
2594  
2595  
2596  
2597  
2598  
2599  
2600  
2601  
2602  
2603  
2604  
2605  
2606  
2607  
2608  
2609  
2610  
2611  
2612  
2613  
2614  
2615  
2616  
2617  
2618  
2619  
2620  
2621  
2622  
2623  
2624  
2625  
2626  
2627  
2628  
2629  
2630  
2631  
2632  
2633  
2634  
2635  
2636  
2637  
2638  
2639  
2640  
2641  
2642  
2643  
2644  
2645

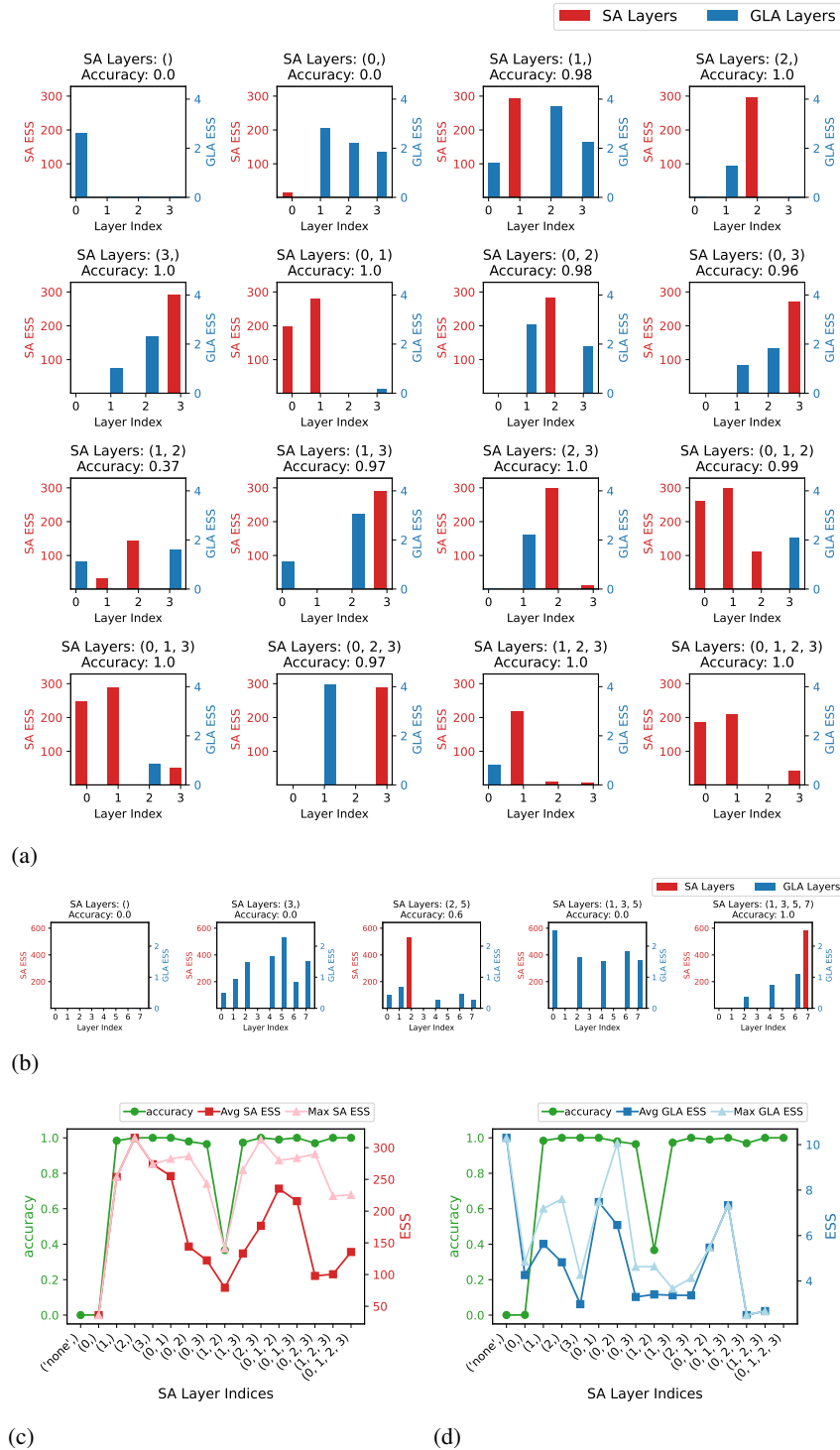


Figure 32: All results presented here are computed using tolerance-based ESS with a tolerance set at  $1e-1$ . Network layers are 0-indexed. (a) Per-layer ESS of all possible 4-layer GLA-SA hybrid networks. Experimental settings can be found in Section D.7. (b) Per-layer ESS of all possible 8-layer GLA-SA Jamba-inspired hybrid networks. Experimental settings can be found in Section D.7. (c) Model accuracy and max/average ESS of SA layers in the 4-layer GLA-SA hybrid networks. (d) Model accuracy and max/average ESS of GLA layers in the 4-layer GLA-SA hybrid networks.

2646 Recall in Section 5.2, we demonstrated an application of post-training ESS analysis through the lens  
2647 of model distillation. Here, we provide another example of post-training analysis that leverages ESS,  
2648 this time in order to gain intuition into learned network dynamics – hybridization. Hybridization is  
2649 the process of arranging different operators in a multi-layer sequence model (Lieber et al., 2024;  
2650 Glorioso et al., 2024; De et al., 2024). Specifically, we measure the per-layer ESS across various  
2651 hybrid networks and find that the precise ordering of layers significantly influences ESS dynamics,  
2652 offering intuition as to why certain hybrids outperform others.

2653 In this section, we present results from a post-training ESS analysis applied to GLA-SA hybrid  
2654 networks to demonstrate the ability of ESS to capture differences among hybrid networks with  
2655 varying topologies. In the first experimental setting, we train all possible 4-layer GLA-SA hybrid  
2656 networks and compute the per-layer ESS on each model. We use the tolerance-based ESS since we  
2657 want to analyze failure modes of learning in hybrid networks. In Figure 32a, we first note that in the  
2658 pure GLA model, many of the layers fail to learn expressive states (as evidenced by the tolerance-  
2659 ESS being 0), offering intuition as to why the model performs so poorly. Moving on to the hybrid  
2660 networks with a single attention layer, we note that all of them perform quite well, except for the  
2661 network that has attention in the first layer. Interestingly, we find that when attention is placed as  
2662 the first layer, it suffers from state collapse. At a higher level, this substantiates why many state-  
2663 of-the-art hybrid networks (such as Jamba) do not place attention as the first layer of the network.  
2664 However, such hybrids are typically constructed purely on the basis of performance: here, ESS is  
2665 able to provide a distinct perspective. Next, examining the hybrids with 2 SA layers, we find that  
2666 the only poor performing topology is with attention placed in the second and third layers. Again, we  
2667 find that the ESS of the attention layers is lower than what we observe in the hybrids that solve the  
2668 task, indicating its usefulness as a proxy for performance beyond the 2-layer non-hybrid networks  
we explored in Section 4.

2669 To clarify this, we examine the maximum/average ESS (computed across layers) of the SA and  
2670 GLA layers separately to understand how each relates to model performance. Notably, we find that  
2671 maximum ESS across attention layers best correlates with accuracy (Figure 32c). Interestingly, the  
2672 average SA layer ESS is a worse proxy for performance, potentially indicating that having a single  
2673 layer with high memory utilization in hybrid networks is more important than having many layers  
2674 with lower memory utilization. This offers support as to why hybrid networks like Jamba have a  
2675 1:7 ratio between attention and non-attention layers. Regarding the GLA layers, we find that despite  
2676 both the maximum and average SS varying across models, they do not correspond to changes in  
2677 accuracy. One possible explanation for this is that since the attention layers are responsible for  
2678 driving the total ESS of the network up due to their unbounded state size, the role of non-attention  
2679 layers in hybrid networks may not be captured entirely by the magnitude of their ESS. Nonetheless,  
2680 this is something we hope to explore in future work.

2681 In the second experimental setting, we move beyond 4-layer GLA-SA hybrids to 8-layer GLA-SA  
2682 hybrids. Here, instead of iterating over all possible topologies, we restrict the space of networks to  
2683 those constructed via the hybridization policy proposed by Jamba. The Jamba hybridization policy  
2684 takes in the number of layers as input and provides a particular hybrid topology as output (refer to  
2685 Lieber et al. (2024) for more details). Since most topologies explored in the 4-layer setting solved  
2686 the task, we both reduce the model dimension of the network and make the task more difficult to  
2687 see if we can observe performance differences across the architectures (model settings can be found  
2688 in Table 8). Unsurprisingly, we find that the pure GLA network is unable to solve the task and also  
2689 realizes a tolerance-based ESS of 0 in all layers (Figure 32b). However, more interesting is the fact  
2690 that while the 2 SA-layer Jamba hybrid partially learns the task, the 3 SA-layer does not. Examining  
2691 the ESS shows that the attention layers in the 3 SA-layer hybrid suffer from state collapse, which  
2692 we know is highly correlated with poor performance on MQAR. This points to a deficiency of  
2693 fixed-topology hybridization policies like Jamba which do not take into account factors like network  
2694 trainability which can significantly influence model performance. Furthermore, this suggests that  
2695 the ESS metric can be used to better inform the construction of hybrid networks. We hope to further  
2696 elucidate these per-layer ESS trends and leverage these insights to construct novel ESS-informed  
2697 hybridization policies in future work.

2698  
2699

E.5 STATE MODULATION OF LARGE LANGUAGE MODELS

State modulation patterns on various open-weight models are illustrated in Figures 33, 34, 35, 36, 37, and 38.

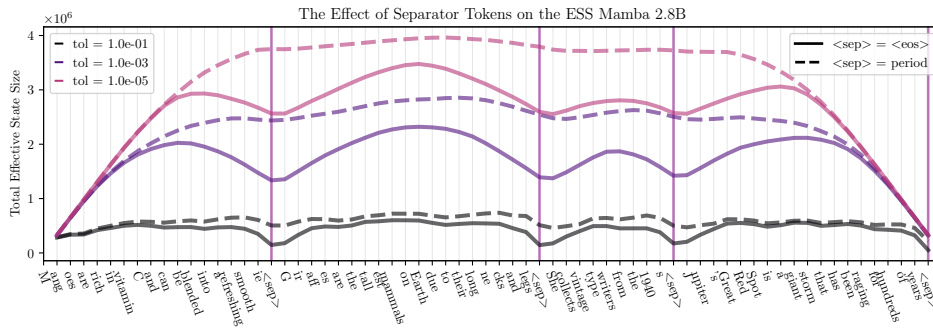


Figure 33

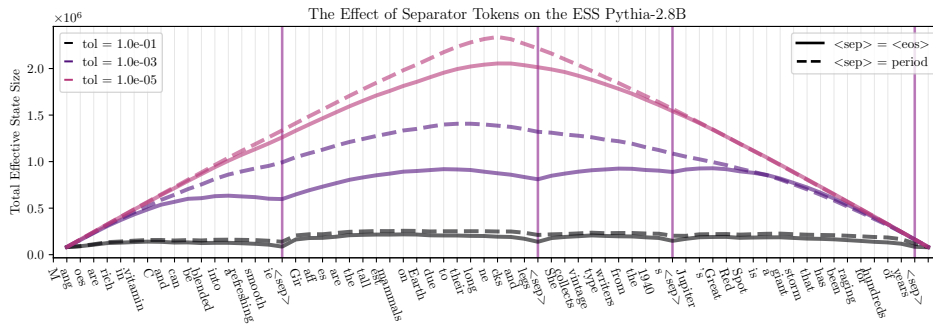


Figure 34

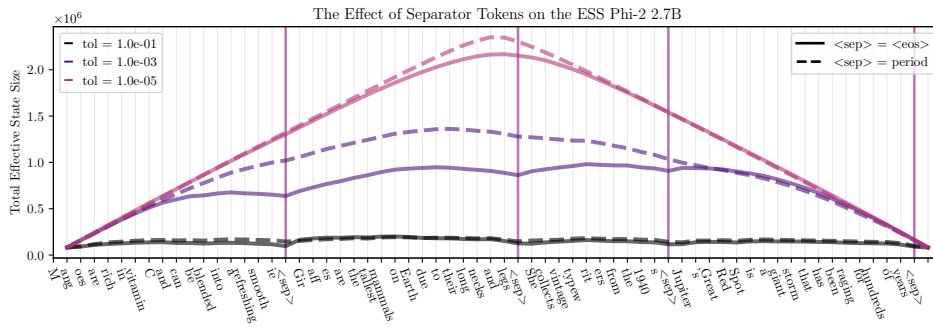


Figure 35

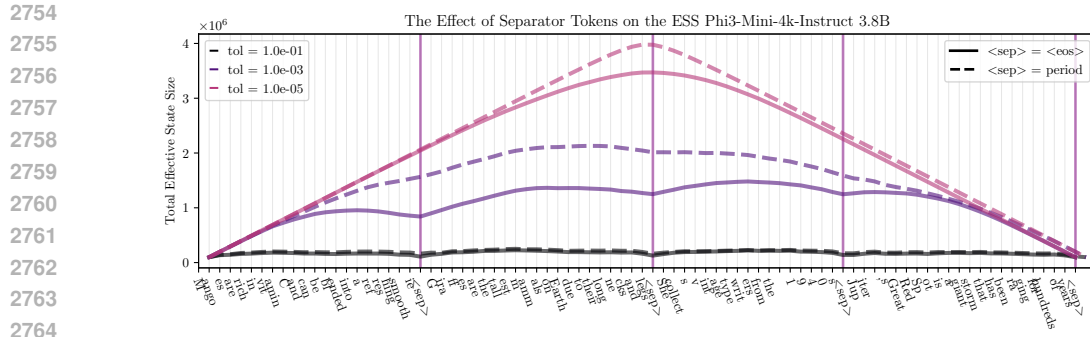


Figure 36

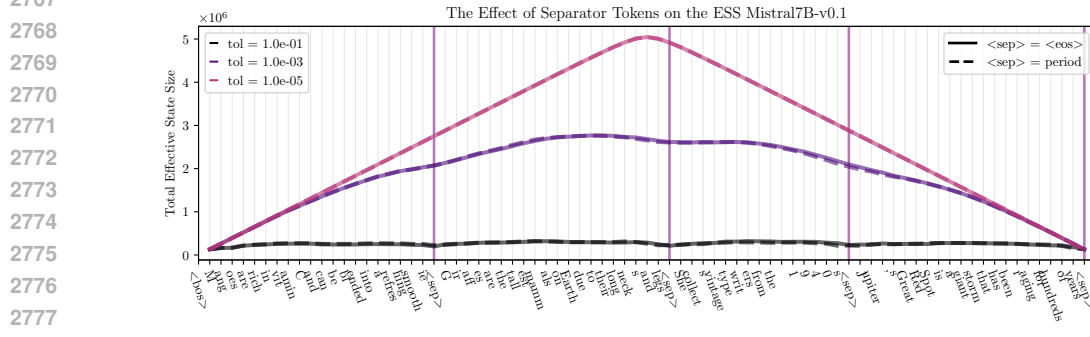


Figure 37

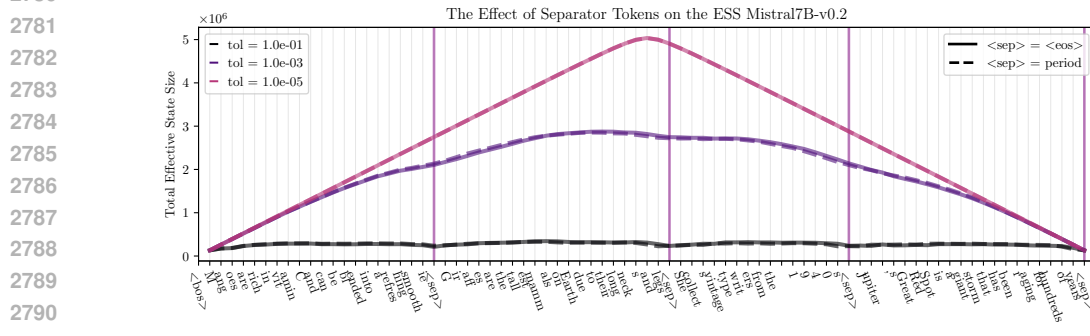


Figure 38

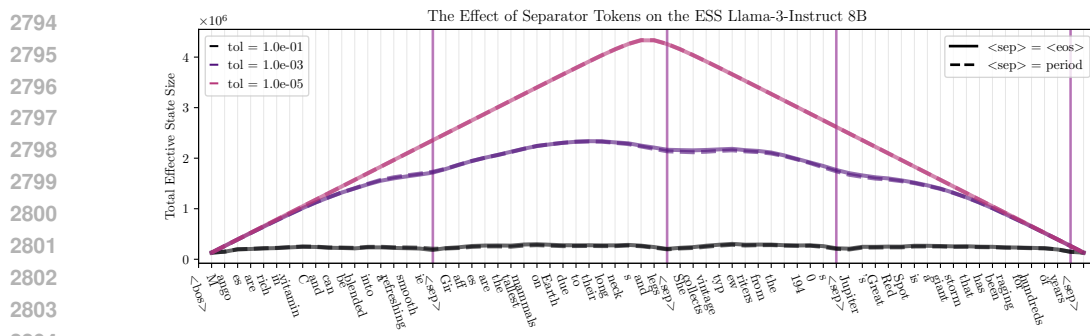


Figure 39

The figures above reveal significant cross-architectural differences in context processing. The attention-based model at a similar 7B scale (Figure 37, 38, and 39) shows minimal change in its ESS pattern when an EOS token is replaced with a period (“.”). In contrast, the limited cache-size state-space model (Falcon Mamba 7B, Figure 7a) exhibits a substantial reduction in state modulation under the same token substitution.

We attribute this difference to a phenomenon we term “preemptive state modulation” in limited state-size models, which stems from fundamental architectural differences. State-space models (SSMs) with limited cache must efficiently manage their finite memory capacity and learn to preemptively modulate state-size to optimize information retention, relying on explicit signals like EOS tokens to trigger context resets. In contrast, attention models with linearly increasing cache can store all past information without the need for selective forgetting, do not require preemptive state modulation, and show less sensitivity to explicit demarcation tokens. This distinction highlights the different strategies employed by various model architectures in managing context across diverse inputs, potentially influencing their performance on tasks requiring long-range recall or context separation.

However, a subset of attention models demonstrated varying state modulation patterns in response to different separator tokens, with this effect being more pronounced in smaller model sizes (see Figure 34, 35, and 36). This phenomenon, while not consistent across all attention architectures, merits deeper exploration.

Figure 40 illustrates the state modulation patterns at different tolerance levels for the four 1B language models (LA, WLA, GLA, SA), trained under identical conditions.

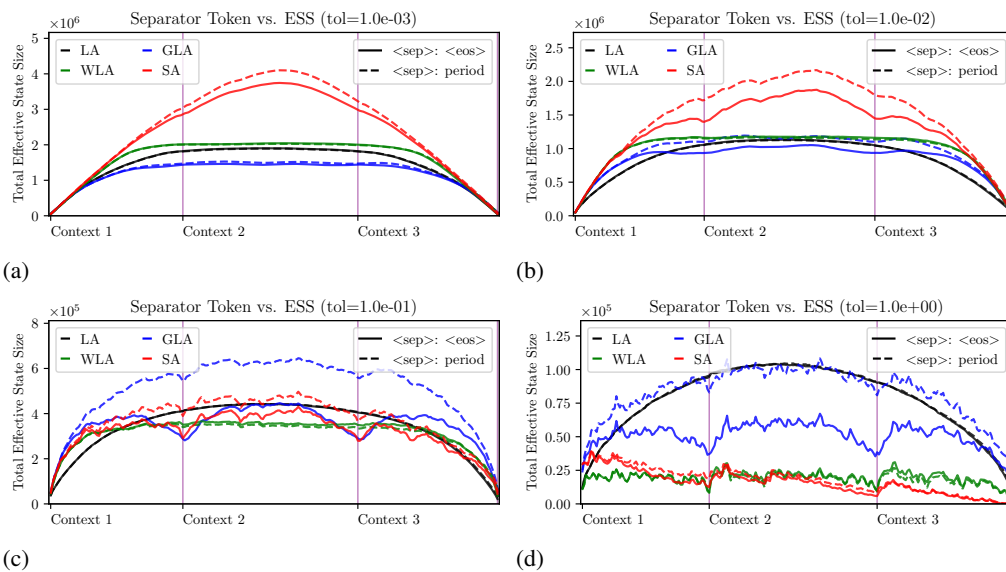


Figure 40: An illustration of the effect of different separator tokens over different layers across different tolerances. Softmax attention exhibits the most pronounced state modulation, beginning at a tolerance level of  $1e-2$ , followed by gated linear attention with significant modulation starting at a tolerance of  $1e-1$ . Weighted linear attention shows minimal modulation, only detectable at a tolerance of 1.0, while linear attention displays no discernible separator token-induced state modulation.

Notably, GLA exhibits a substantial variation in state modulation depending on the separator token, consistent with our earlier observations in Falcon Mamba regarding preemptive state modulation. In contrast, SA shows a smaller, yet non-trivial, effect. WLA and LA show no discernible differences across separator tokens, which may be attributed to their overall limited ability to modulate state size.

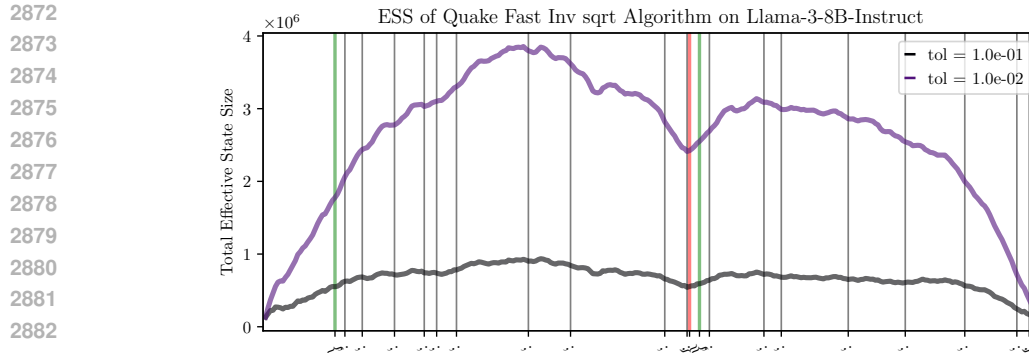
2862 E.6 MISCELLANEOUS

2863 E.6.1 EFFECTIVE STATE-SIZE ON C++ CODE

2864 Beyond sentence delimiters such as periods and end-of-speech tokens (discussed in Section 5.3), we  
2865 observe similar “dips” in effective state-size where there are scope delimiter tokens such as “}”.

2866 The following plots demonstrate the ESS pattern of Llama3-8B processing the C++ code of a fast  
2867 inverse square root algorithm and a Fibonacci sequence generator algorithm.

2868 **Quake fast inverse square-root algorithm:**



2884 Figure 41: Effective state-size over a quake fast inverse square root algorithm’s code.

```

2886
2887 1 #include <iostream>
2888 2 #include <cmath>
2889 3
2889 4 // Quake Fast Inverse Square Root function
2890 5 float quakeFastInvSqrt(float number) {
2891 6     long i;
2892 7     float x2, y;
2893 8     const float threehalfs = 1.5F;
2894 9
2894 10     x2 = number * 0.5F;
2895 11     y = number;
2896 12     i = *(long*)&y;           // Bit-level hacking: convert float to
2897 13     long                          // long
2898 14     i = 0x5f3759df - (i >> 1); // Initial magic number and bit shift
2899 15     y = *(float*)&i;         // Convert back from long to float
2900 16
2900 17     // Newton's method step for refining the result
2901 18     y = y * (threehalfs - (x2 * y * y)); // First iteration
2902 19
2903 20     return y;
2904 21 }
2905 22 int main() {
2906 23     float number;
2907 24
2907 25     // Input: Get the number from the user
2908 26     std::cout << "Enter a number: ";
2909 27     std::cin >> number;
2910 28
2911 29     // Output: Display the result using the Quake fast inverse sqrt
2912 30     float quake_result = quakeFastInvSqrt(number);
2913 31     std::cout << "Quake Fast Inverse Sqrt: " << quake_result << std::endl
2914 32     ;
2915 33
2915 34     // Compare with standard sqrt function
2915 35     float std_result = 1.0f / std::sqrt(number);

```

```

2916 35     std::cout << "Standard Inverse Sqrt: " << std_result << std::endl;
2917 36
2918 37     return 0;
2919 38 }

```

### Fibonacci sequence generating algorithm:

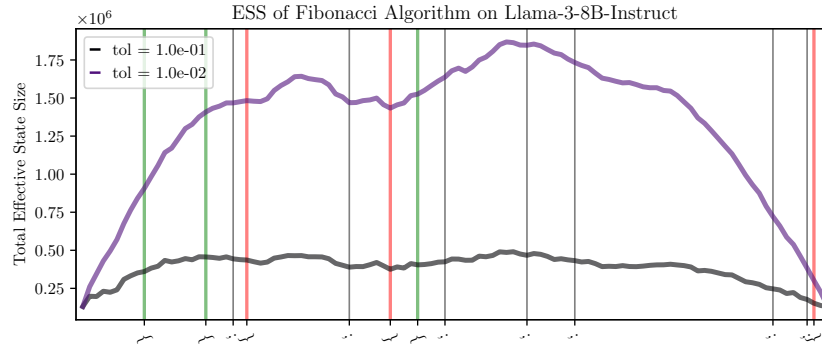


Figure 42: Effective state-size over a Fibonacci sequence generator algorithm’s code.

```

2939
2940 1 #include <iostream>
2941 2 int fibonacci(int n) {
2942 3     if (n <= 1) {
2943 4         return n;
2944 5     }
2945 6     return fibonacci(n - 1) + fibonacci(n - 2); // Recursive case
2946 7 }
2947 8 int main() {
2948 9     int n;
2949 10    std::cout << "Enter a positive integer: ";
2950 11    std::cin >> n;
2951 12    std::cout << "Fibonacci number at position " << n << " is: " <<
2952 13    fibonacci(n) << std::endl;
2953 14    return 0;

```

## E.6.2 HOW THE NUMBER OF PROMPTING SHOTS AFFECTS THE EFFECTIVE STATE-SIZE OF LANGUAGE MODELS

Here, we explore how varying the number of shots when prompting large language models affects their effective state-size patterns. We use Phi-2 as the candidate attention model and Mamba-2.8B as the state-space model. The task we tested this on is MMLU (elementary mathematics).

At the start of the Q&A section for the attention model, there is a noticeable difference in state size between 0-shot and 1-shot prompts. Beyond 1-shot, the difference in ESS appears minimal. For the state-space model, varying the number of shots has minimal impact on the effective state-size.

Although these sparse experimental results require further investigation, we note the stark difference in the effective state-size patterns between these two architectures, which provides additional insights into understanding the fundamental differences in the way prompts are processed across models.

2970  
 2971  
 2972  
 2973  
 2974  
 2975  
 2976  
 2977  
 2978  
 2979  
 2980  
 2981  
 2982  
 2983  
 2984  
 2985  
 2986  
 2987  
 2988  
 2989  
 2990  
 2991  
 2992  
 2993  
 2994  
 2995  
 2996  
 2997  
 2998  
 2999  
 3000  
 3001  
 3002  
 3003  
 3004  
 3005  
 3006  
 3007  
 3008  
 3009  
 3010  
 3011  
 3012  
 3013  
 3014  
 3015  
 3016  
 3017  
 3018  
 3019  
 3020  
 3021  
 3022  
 3023

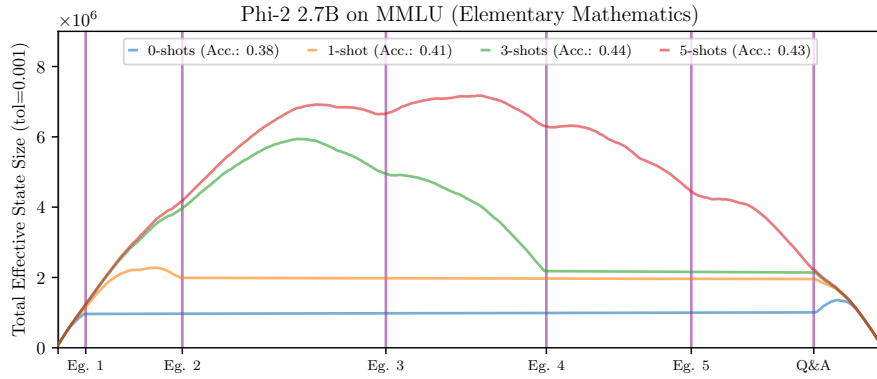


Figure 43: The variation in effective state-size with a varying number of shots (2.7B Attention).

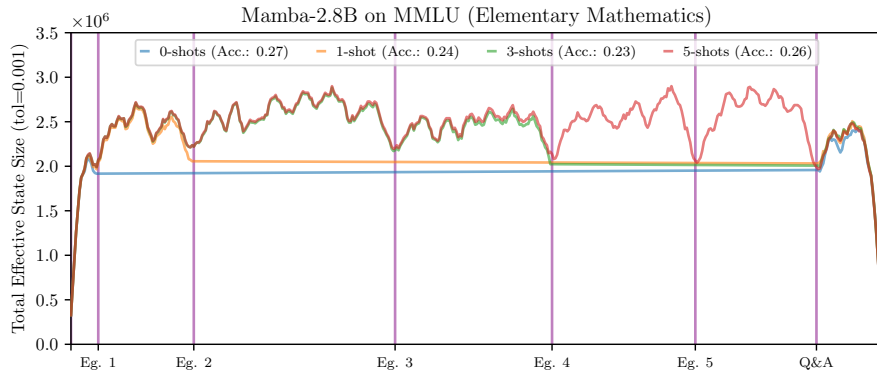


Figure 44: The variation in effective state-size with a varying number of shots (2.8B State-Space Model).