# 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.

### 028 029 030 031

032

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

## 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 trans former (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 040 expressed as either linear systems (y = Tu) or systems with *input-varying linear operators* 041  $(y = f_T(u)u)$ , the latter of which we abbreviate as LIV<sup>1</sup>. LIVs generalize the notion of adap-042 tive, or data-controlled operators to a broader class than previously described in Massaroli et al. 043 (2021); Poli et al. (2023). The input-varying linear operator framework decouples the input-varying 044 *featurization*  $u \mapsto T := f_T(u)$  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 046 uniformly formulated as linear systems, including models like convolutions [40; 47; 35; 56], linear 047 state-transition recurrences [15; 18; 27; 61; 29; 26; 75; 49; 16], and attention variants [68; 34; 65].

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) [45; 69; 1; 4; 74; 63]. However, these visualizations alone often fail to highlight critical properties that explain

 <sup>&</sup>lt;sup>1</sup>We abbreviate them as LIVs instead of IVLs to conform to the classical convention, where linear timevarying 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".



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

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>2</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.

078 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 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

105

065 066

067

068

069

071

079

081

082

084

085

090

091 092 093

094

 <sup>&</sup>lt;sup>2</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).

108 further be broken down into time-invariant (LTI) and time-varying (LTV) systems, depending on 109 whether or not the parameters of the realized recurrence change as a function of time (i.e. sequence 110 index). This taxonomy has also been used in prior work (Zancato et al., 2024) to motivate the con-111 struction of novel architectures. In modern deep learning, LTI systems (Lindquist & Picci, 1979; 112 Ljung, 1999; Ho & Kalman, 1966) 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 113 expressive LIV class (Gu & Dao, 2024; Poli et al., 2023; Vaswani et al., 2023) which parameterizes 114 the state-space dynamics as a function of the input. Since LIVs constitute the minimal superclass 115 containing most sequence models of practical interest, the goal of this work is to develop a quan-116 titative means of analyzing models under the LIV abstraction. To do so, we draw from minimal 117 realization theory (DeWilde & van der Veen, 1998; Akaike, 1974), which has previously been ap-118 plied within the scope of classical linear systems but has yet to be extended to modern LIV sequence 119 models.

120

121 Quantifying memory utilization in sequence models. As discussed in Section 1, the current 122 landscape of examining memory in sequence models primarily reduces to qualitative measures or 123 limited quantitative ones. Here, we expand on the latter, noting that commonly used measures such 124 as model size (Hoffmann et al., 2022) and cache/state size (which we define formally as theoretically 125 realizable state-size in Section 3) are imperfect means of measuring memory in sequence models. 126 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 127 model pipeline such as data, initialization, and optimization. To capture these aspects of network 128 training, we apply the notions of numerical and effective rank (Roy & Vetterli, 2007) to the minimal 129 realization problem, which gives rise to the ESS metric. We elaborate on this in Section 3. 130

We highlight the concept of semiseparable rank (Vandebril et al., 2005; Xia et al., 2010), which has been recently explored in Dao & Gu (2024). Similar to ESS, semiseparable rank examines the rank of specific submatrices of T. However, as semiseparable rank is primarily motivated by the design of efficient algorithms, we find it most closely aligns with the concept of theoretically realizable state size (TSS) in the context of our work. In this study, we provide complete derivations of ESS, clearly distinguish it from TSS (and by extension semiseparable rank), and demonstrate its effectiveness as a proxy for memory utilization in practical settings for modern LIV sequence models.

## 3 THEORY

139 140

150

138

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 (and other non-temporal dimensions) 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):

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

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:

$$s_{i+1} = A_i s_i + B_i u_i, \quad y_i = C_i s_i + D_i u_i,$$
 (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.

176 Refer to Section C.3 for the proof.

177 Note that since we use the flattened notation, where  $T \in \mathbb{R}^{d\ell \times d\ell}$ ,  $\operatorname{rank}(H_i)$  determines the minimal 178 state-size required for processing all d channels in a layer via a recurrence at the *i*th index in the 179 sequence. We formally refer to this metric as per-sequence index *effective state-size* (ESS) and 180 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 182 system is given by its minimal state-size which is directly interpretable as a measure of model 183 memory utilization. For LIVs, however, ESS is also a function of the input  $(\operatorname{rank}(f_T(u)_{i \leq i-1}))$ which means that the minimal realization process outlined in the proof of Theorem 3.1 is no longer 185 guaranteed to obtain recurrences that preserve causality (i.e. the minimally realized features  $A_i^*$ 186 depend on future inputs  $u_k$ , k > i<sup>3</sup>. Nevertheless, ESS lower bounds the state-size  $n_i$  (refer to 187 Section C.2.2), meaning that for any LIV, an equivalent recurrence must necessarily materialize a 188 state-size at least as large as its ESS. Therefore, we claim that ESS serves as a proxy for memory 189 utilization in not only linear systems, but in LIVs as well. We empirically validate the usefulness of 190 this interpretation of ESS for LIVs in Sections 4 and 5.

191

207

208

209

181

165

192 Memory capacity in LIVs. The memory capacity of LIVs is given by the state-size  $n_i$ . We 193 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 194 attention, we resort to the trivial realization as shown in Equation C.2.5, in which TSS (for a single 195 channel) is equal to the sequence index *i*. For models like SSMs and linear attention variants, TSS 196 is equal to the state-size defined by their recurrent formulation. We refer readers to Section D.2197 for detailed operator-specific derivations of TSS. Recall that ESS depends not only on the model's functional form, but also on the input data, optimization, and more generally anything that impacts 199 the realization of T; in contrast, TSS is limited in that it is determined solely by the model's structure. 200

201 3.1 COMPUTING EFFECTIVE STATE-SIZE

In practice, computing ESS requires a few additional considerations due to the numerical errors and approximations involved in computing the rank of matrices. We propose two approaches – both of which rely on singular values ( $\Sigma_i$ ) from taking the singular value decomposition (SVD) of  $H_i$ (which equals  $f_T(u)_{i::,i-1}$  for LIVs) – that provide complementary perspectives on the same metric.

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

tolerance-ESS $(H_i, \tau) \coloneqq |\{\sigma_i^m : \sigma_i^m > \tau\}|$ , where  $\sigma_i^m \in \Sigma_i$ ,  $U_i \operatorname{diag}(\Sigma_i) V_i = \operatorname{SVD}(H_i)$ . (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  $(||T_{ij} - T_{ij}^*||_2 \le \tau)$ .

<sup>&</sup>lt;sup>3</sup>An example of a causality preserving realization of LIVs is the trivial realization shown in Equation (C.2.5).

**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  $H_i$ :

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

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. An additional discussion comparing tolerance and entropy ESS, along with our code for computing them, can be found in Section D.1.

227

220

221

228 **Computational complexity of ESS.** Since SVD scales cubically with the size of a square matrix, 229 the time complexity of computing ESS for a single layer with d channels, processing an input of 230 sequence length  $\ell$ , is  $O((d\ell)^3)$ . Fortunately, most modern sequence models (outside of S5 (Smith et al., 2023)) process the d channels independently (i.e. they are SISO). This means that each layer's 231  $d\ell \times d\ell$  operator T can be decomposed into d independent  $\ell \times \ell$  operators, which reduces the 232 ESS computation to an SVD of d independent operators ( $\hat{T}^{\alpha} \coloneqq T^{\alpha\alpha}$ ;  $\alpha \in [d]$ ). Consequently, a 233 per-channel ESS can be computed and summed to obtain the equivalent per-sequence index ESS. 234 This reduces the computational cost to  $O(d\ell^3)$ . Note that in models like attention, where channels 235 share the same recurrence within a single head, the time complexity is further reduced by a factor 236 proportional to the head dimension. Additionally, in recurrent models with bounded TSS ( $n \ll \ell$ ), 237 a truncated SVD can be employed to decrease the time complexity to  $O(\ell^2 n d)$ . 238

239 Aggregation of ESS across model and data dimensions. Initially, we formulated ESS as a per-240 sequence index metric, which is general in that it applies to all models that realize an operator T. 241 Recall in the last section, we extended ESS along the channel dimension for SISO models. For 242 LIV SISO models in particular, we can further extend ESS along the batch size dimension, since 243 ESS is a function of the input. This means that for a multi-layer LIV SISO model with m layers, d244 channels, batch size b, and sequence length  $\ell$  (which describes the entire set of models we analyze 245 in this work), ESS  $\in \mathbb{R}^{m \times d \times b \times \ell}$ . Since ESS is a multidimensional tensor, there are various ways to aggregate it across the model and data dimensions. We define two particular modes of aggregation 246 termed *average* ESS and *total* ESS as follows: 247

average ESS = 
$$\frac{1}{mdb\ell} \sum_{\alpha \in [d]} \sum_{\beta \in [m]} \sum_{\gamma \in [b]} \sum_{i \in [\ell]} \text{ESS}_i^{\alpha\beta\gamma}$$
 total ESS = average ESS

\*d

Average ESS marginalizes across all of the ESS tensor dimensions by taking a mean, creating a perchannel measure of ESS. Since average ESS is the metric used throughout most of our experiments, we will refer to it as ESS unless otherwise specified. For models like softmax attention, where average TSS (which is computed analogously to average ESS) depends only on the sequence index *i* and thus remains constant as a function of model width (i.e. channel dimension), we instead capture a model-dependent statistic, by first summing the ESS across channels and then averaging over the remaining dimensions. This approach allows both ESS (and TSS) to vary as a function of model width, a metric that we refer to as total ESS. For more details, refer to Sections D.1 and D.2.

258 259 260

261 262

263

264

248

249 250

251

252

253

254

255

256

257

## 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.

265 266

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

284

285

286

287

288

289

290

291

292

293

295

296 297

298

299

300

301

302 303

304

305



Figure 2: Scatter plots of accuracy vs ESS/kv across featurizers. Within each featurizer plot, all 278 task-model configurations from the sweep corresponding to each featurizer are shown. 279

280 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 282 compression, which showcase analogous trends. 283

**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 306 case of MQAR, we do so by modulating the number of key-value (kv) 307 pairs the models are tasked to match, as well as the total sequence length 308 of the prompt. Within each model, we sweep across varying TSS. For 309 each task-model configuration, we compute the ESS and accuracy on a 310 validation set every 10 epochs. We will refer to the entire space of tasks 311 and models across which we sweep as the task-model space. Finally, we 312 split our profiling of ESS into two sections: cross task-model analysis 313 (Section 4.1) and within task-model analysis (Section 4.2). For more 314 details on the setup, refer to Section D.3. 315



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

316 4.1 CROSS TASK-MODEL ANALYSIS 317

318 Our first goal is to understand how ESS empirically captures memory utilization by studying its 319 correlation with post-training MQAR performance across the entire task-model space. To appropri-320 ately analyze ESS across tasks, we normalize it by the memory demands of MQAR, constructing an 321 adjusted form of ESS given by ESS/kv.

322 Finding 1: Measured over the entire task-model space, ESS/kv exhibits a significantly higher 323 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

329

330 331

332

333

334

335

336

337

338

339



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>4</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 356 has insufficient TSS to fully learn a task, resulting in its 357 ESS converging near its TSS. This is reflected in its ES-358 S/TSS (which we refer to as state utilization) residing near 359 1. We observe this in Figure 5 where we note that mod-360 els with a TSS of 8 perform worse as the task difficulty 361 scales due to a saturated state. State collapse, on the other 362 hand, refers to the scenario in which a model has sufficient TSS to learn (or partially learn) a task, but its ESS fails to 363 increase during training, resulting in a heavily underuti-364 lized state. With respect to state collapse, we observe the following: 366



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

**Finding 3:** For GLA and WLA, state collapse occurs in the high kv bucket of task-model space (i.e.  $kv = 2^7$ ) 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.

<sup>&</sup>lt;sup>4</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

402 Initialization in weight space plays a crucial role in machine learning, significantly impacting model 403 convergence and training stability (Glorot & Bengio, 2010). We extend this concept to the ini-404 tialization of recurrent models in state space, leaning on the intuition from Figure 2 that suggests 405 higher ESS can enhance performance. Namely, we illustrate how ESS at initialization can be used 406 to inform featurizer selection – the selection of the function that maps the input to the operator T = f(u) or equivalently the recurrent features  $(A_i(u_{ii}), B_i(u_{ii}), C_i(u_{ii}), D_i(u_{ii}))_{i \in [\ell]}$  - and ini-407 tialization schemes. In doing so, we uncover design flaws of a prominent model, S6 (Mamba) (Gu 408 & Dao, 2024). 409

410

387

388

389 390

391 392

393

394

395

396

397

398

399 400

401

411 **ESS-informed featurizer selection.** To study the relationship between memory capacity and 412 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 413 settings outlined in Poli et al. (2024) (see Tables 2, 3, and 4 for details), we observe that S6 is en-414 tirely unable to learn MQAR (accuracy  $\approx 0$ ) across multiple scales of TSS (16 - 256) as shown in 415 Figure 25. This aligns with results from Yang et al. (2024b), which independently demonstrate the 416 poor performance of S6 without the additional short convolutional layer on a different in-context 417 recall task. To investigate the cause, we look into how S6 is preconditioned to utilize its memory by 418 computing its ESS when processing a Gaussian noise input, prior to training. 419

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.

429

**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:

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

450

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

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

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

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

## 451 5.2 POST-TRAINING ANALYSIS

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

457 ESS-informed model-order reduction. Model-order reduction refers to the process of improv-458 ing model efficiency by reducing state-size while retaining performance. Previous works, such 459 as Massaroli et al. (2023), have explored the distillation of linear time-invariant (LTI) operators 460  $(T_{ij} = T_{i+k,j+k})$  into linear recurrences with small state-sizes using backpropagation. Other tech-461 niques for model-order reduction such as modal truncation and balanced truncation (Beliczynski 462 et al., 1992; Gawronski & Juang, 1990) are also applicable to LTIs. In this study, however, we 463 are concerned with improving the efficiency of general LIVs. Since ESS serves as a lower bound for the minimally realizable TSS (Section 3), we postulate that ESS can be used as a heuristic for 464 conducting model-order reduction. 465

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

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

These findings position ESS as a useful heuristic for predicting model compressibility, enabling efficient estimation of the potential for state-size reduction without extensive experimentation.

479 480 481

5.3 STATE MODULATION OF LARGE LANGUAGE MODELS

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



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).

We begin by evaluating the  $(\text{total ESS})_i$  of open-weight pre-trained models. Borrowing notation from Section 3.1,  $(\text{total ESS})_i$  is defined as follows:

$$(\text{total ESS})_i = \frac{1}{mb} \sum_{\alpha \in [d]} \sum_{\beta \in [m]} \sum_{\gamma \in [b]} \text{ESS}^{\alpha \beta \gamma}$$

Since we do not marginalize along the sequence dimension, we will examine this measure quali-504 tatively by observing how it changes as a function of the sequence index. Our analysis shown in 505 Figure 7a (and more broadly in Section E.5) reveals an intriguing phenomenon: ESS undergoes 506 a noticeable dip whenever an end-of-speech (EOS) token is encountered (refer to Section D.8 for 507 experimental details). This behavior aligns with our intuition regarding the role of EOS tokens and 508 provides a quantitative measure of how effectively a model can 'reset' or 'forget' past contexts when 509 transitioning between distinct segments of text. To investigate these effects in a more controlled environment, we trained four 1B parameter models (LA, WLA, GLA, and SA as described in Section 510 4) under identical conditions (see Table 9). 511

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).

514 SA exhibits the most pronounced state modulation, beginning at a tolerance level of 1e-2, while also 515 realizing the largest ESS. GLA follows, with modulation emerging at a tolerance of 1e-1. WLA 516 shows minimal modulation, only detectable at a tolerance of 1.0, while LA displays no discernible 517 state modulation in response to separator tokens, demonstrating a clear lack of ability to modulate 518 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 519 dataset (Penedo et al., 2024)) are similar across SA, WLA, and GLA, significant differences emerge 520 when considering the bigram recall perplexity metric introduced by Arora et al. (2023). 521

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 finding suggests that state modulation serves as a key mechanism enabling models to effectively manage complex context dependencies, directly impacting their performance on recall-heavy tasks.

526 527 528

529

495

496

497 498

499

500

501 502

## 6 CONCLUSION

530 In this work, we propose effective state-size (ESS), a measure of memory utilization in sequence 531 models derived using dynamical systems theory. We motivate this metric as a valuable tool for ana-532 lyzing memory utilization in LIVs by demonstrating its strong correlation with performance across 533 a wide range of synthetic tasks. In doing so, we find that ESS offers a versatile framework for under-534 standing 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 536 improve beyond the existing performance-efficiency trade-offs found in recurrent models. Finally, 537 we extend the ESS framework to language tasks, introducing the idea of state modulation – a concept that proves crucial for performance on bigram recall tasks. Overall, this work establishes ESS 538 as a foundational tool for understanding and improving sequence model performance, opening new avenues for optimizing memory utilization and, more generally, model efficiency.

# 540 REPRODUCIBILITY STATEMENT

To ensure reproducibility, we utilized open-source models and tasks, adhering to default task configurations unless otherwise specified. All crucial configurations are detailed in either the main text or the appendix. Additionally, our code for computing both the tolerance-ESS and entropy-ESS is provided in the appendix (Section D.1.3).

548 REFERENCES

546 547

554

555

556

558 559

561

562

577

578

579

580

584

585

- Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers, 2020. URL https://arxiv.org/abs/2005.00928. (pages 1, 20).
- H. Akaike. Stochastic theory of minimal realization. *IEEE Transactions on Automatic Control*, 19 (6):667–674, 1974. doi: 10.1109/TAC.1974.1100707. (page 3).
  - Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-context language learning: Architectures and algorithms, 2024. URL https://arxiv.org/abs/2401.12973. (page 20).
  - Ameen Ali, Itamar Zimerman, and Lior Wolf. The hidden attention of mamba models, 2024. URL https://arxiv.org/abs/2403.01590. (pages 1, 20).
  - Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.3, knowledge capacity scaling laws, 2024. URL https://arxiv.org/abs/2404.05405. (page 2).
- Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Alsubaie, Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwairesh, Areeb Alowisheq, M Saiful Bari, and Haidar Khan. When benchmarks are targets: Revealing the sensitivity of large language model leaderboards, 2024. URL https://arxiv.org/abs/2402.01781. (page 21).
- Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri Rudra, and Christopher Ré. Zoology: Measuring and improving recall in efficient language models, 2023. URL https://arxiv.org/abs/2312.04927. (pages 10, 10, 20, 21, 21, 30, 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 9, 20).
  - Jimmy Ba, Geoffrey Hinton, Volodymyr Mnih, Joel Z. Leibo, and Catalin Ionescu. Using fast weights to attend to the recent past, 2016. URL https://arxiv.org/abs/1610.06258. (page 21).
- Bartlomiej Beliczynski, Izzet Kale, and Gerald D Cain. Approximation of fir by iir digital filters:
   An algorithm based on balanced model reduction. *IEEE Transactions on Signal Processing*, 40 (3):532–542, 1992. (page 9).
  - Satwik Bhattamishra, Arkil Patel, Phil Blunsom, and Varun Kanade. Understanding in-context learning in transformers and llms by learning to learn discrete functions, 2023. URL https://arxiv.org/abs/2310.03016. (page 20).
- Srinadh Bhojanapalli, Chulhee Yun, Ankit Singh Rawat, Sashank J. Reddi, and Sanjiv Kumar. Low-rank bottleneck in multi-head attention models, 2020. URL https://arxiv.org/abs/2002.07028. (pages 2, 20).
- Aviv Bick, Kevin Y. Li, Eric P. Xing, J. Zico Kolter, and Albert Gu. Transformers to ssms: Distilling quadratic knowledge to subquadratic models, 2024. URL https://arxiv.org/abs/ 2408.10189. (page 9).

603

607

611

624

625

626

627

631

632

633

634

- 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).
- Chi-Tsong Chen. Linear System Theory and Design. Oxford University Press, Inc., USA, 3rd 598 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, 3, 6, 19, 20, 27, 29).
- Soham De, Samuel L. Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Al-604 bert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, Guillaume Des-605 jardins, Arnaud Doucet, David Budden, Yee Whye Teh, Razvan Pascanu, Nando De Freitas, and 606 Caglar Gulcehre. Griffin: Mixing gated linear recurrences with local attention for efficient language models, 2024. URL https://arxiv.org/abs/2402.19427. (page 50). 608
- P. DeWilde and A.J. van der Veen. Time-Varying Systems and Computations. Springer US, 609 1998. ISBN 9780792381891. URL https://books.google.co.jp/books?id= 610 n3bEniJ2Wx8C. (pages 1, 3, 4, 19, 24).
- 612 Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure 613 attention loses rank doubly exponentially with depth, 2023. URL https://arxiv.org/ 614 abs/2103.03404. (page 20). 615
- Emilien Dupont, Arnaud Doucet, and Yee Whye Teh. Augmented neural odes, 2019. URL https: 616 //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, 2021. https://transformer-circuits.pub/2021/framework/index.html. (page 20). 623
  - Daniel Y. Fu, Tri Dao, Khaled K. Saab, Armin W. Thomas, Atri Rudra, and Christopher Ré. Hungry hungry hippos: Towards language modeling with state space models, 2023. URL https:// arxiv.org/abs/2212.14052. (page 20).
- Wodek Gawronski and Jer-Nan Juang. Model reduction in limited time and frequency intervals. In-628 ternational Journal of Systems Science, 21(2):349–376, 1990. doi: 10.1080/00207729008910366. 629 URL https://doi.org/10.1080/00207729008910366. (page 9). 630
  - Paolo Glorioso, Quentin Anthony, Yury Tokpanov, James Whittington, Jonathan Pilault, Adam Ibrahim, and Beren Millidge. Zamba: A compact 7b ssm hybrid model, 2024. URL https: //arxiv.org/abs/2405.16712. (page 50).
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural 635 networks. In Yee Whye Teh and Mike Titterington (eds.), Proceedings of the Thirteenth Interna-636 tional Conference on Artificial Intelligence and Statistics, volume 9 of Proceedings of Machine Learning Research, pp. 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR. 638 URL https://proceedings.mlr.press/v9/glorot10a.html. (page 8). 639
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2024. 640 URL https://arxiv.org/abs/2312.00752. (pages 1, 3, 6, 8, 19, 20, 27, 28). 641
- 642 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured 643 state spaces, 2022a. URL https://arxiv.org/abs/2111.00396. (pages 1, 3, 19, 25, 644 27). 645
- Albert Gu, Ankit Gupta, Karan Goel, and Christopher Ré. On the parameterization and initialization 646 of diagonal state space models, 2022b. URL https://arxiv.org/abs/2206.11893. 647 (page 19).

661

681

682

683

684

685

686

687

688

- 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. (pages 1, 6, 20).
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja-652 cob Steinhardt. Measuring massive multitask language understanding, 2021. URL https: 653 //arxiv.org/abs/2009.03300. (page 21). 654
- 655 . L. Ho and R. E. Kalman. Editorial: Effective construction of linear state-variable models from 656 input/output functions. at - Automatisierungstechnik, 14(1-12):545–548, 1966. doi: doi:10.1524/ 657 auto.1966.14.112.545. URL https://doi.org/10.1524/auto.1966.14.112.545. (page 3). 658
- 659 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. Neural Comput., 9(8): 660 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). 662
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza 663 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hen-664 nigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, 665 Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 666 Training compute-optimal large language models, 2022. URL https://arxiv.org/abs/ 667 2203.15556. (page 3). 668
- 669 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are 670 rnns: Fast autoregressive transformers with linear attention, 2020. URL https://arxiv. 671 org/abs/2006.16236. (pages 1, 6, 9, 19, 20, 27).
- 672 Serkan Kiranyaz, Onur Avci, Osama Abdeljaber, Turker Ince, Moncef Gabbouj, and Daniel J. 673 Inman. 1d convolutional neural networks and applications: A survey, 2019. URL https: 674 //arxiv.org/abs/1905.03554. (page 1). 675
- 676 Opher Lieber, Barak Lenz, Hofit Bata, Gal Cohen, Jhonathan Osin, Itay Dalmedigos, Erez Safahi, Shaked Meirom, Yonatan Belinkov, Shai Shalev-Shwartz, Omri Abend, Raz Alon, Tomer Asida, 677 Amir Bergman, Roman Glozman, Michael Gokhman, Avashalom Manevich, Nir Ratner, Noam 678 Rozen, Erez Shwartz, Mor Zusman, and Yoav Shoham. Jamba: A hybrid transformer-mamba 679 language model, 2024. URL https://arxiv.org/abs/2403.19887. (pages 34, 50, 50). 680
  - Anders Lindquist and Giorgio Picci. On the stochastic realization problem. SIAM J. Control Optim., 17(3):365-389, May 1979. ISSN 0363-0129. doi: 10.1137/0317028. URL https://doi. org/10.1137/0317028. (page 3).
  - L. Ljung. System Identification: Theory for the User. Prentice Hall information and system sciences series. Prentice Hall PTR, 1999. ISBN 9780136566953. URL https://books.google. com/books?id=nHFoQgAACAAJ. (page 3).
  - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL https: //arxiv.org/abs/1711.05101. (pages 32, 33, 34, 35).
- 690 Paul A. Lynn and Wolfgang Fuerst. Introductory digital signal processing with computer applica-691 tions (revised ed.). John Wiley & Sons, Inc., USA, 1994. ISBN 0471943746. (page 1). 692
- Eric Martin and Chris Cundy. Parallelizing linear recurrent neural nets over sequence length, 2018. 693 URL https://arxiv.org/abs/1709.04057. (page 2). 694
- Stefano Massaroli, Michael Poli, Jinkyoo Park, Atsushi Yamashita, and Hajime Asama. Dissecting 696 neural odes, 2021. URL https://arxiv.org/abs/2002.08071. (page 1). 697
- Stefano Massaroli, Michael Poli, Daniel Y. Fu, Hermann Kumbong, Rom N. Parnichkun, Aman Timalsina, David W. Romero, Quinn McIntyre, Beidi Chen, Atri Rudra, Ce Zhang, Christopher 699 Re, Stefano Ermon, and Yoshua Bengio. Laughing hyena distillery: Extracting compact recur-700 rences from convolutions, 2023. URL https://arxiv.org/abs/2310.18780. (pages 9, 701 19).

702 Zeping Min and Zhong Li. On the efficiency of transformers: The effect of attention rank, 2024. 703 URL https://openreview.net/forum?id=U9sHVjidYH. (pages 2, 20). 704

Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, 705 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Gan-706 guli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, 707 Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam 708 McCandlish, and Chris Olah. In-context learning and induction heads. Transformer Circuits Thread, 2022a. https://transformer-circuits.pub/2022/in-context-learning-and-induction-710 heads/index.html. (pages 1, 20). 711

712 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, 713 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, 714 Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, 715 716 and Chris Olah. In-context learning and induction heads. Transformer Circuits Thread, 2022b. https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html. (page 717 20). 718

- 719 Alan V. Oppenheim, Alan S. Willsky, and S. Hamid Nawab. Signals & systems (2nd ed.). Prentice-720 Hall, Inc., USA, 1996. ISBN 0138147574. (page 1). 721
- Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pas-722 canu, and Soham De. Resurrecting recurrent neural networks for long sequences, 2023. URL 723 https://arxiv.org/abs/2303.06349. (pages 19, 27). 724
- 725 Rom N. Parnichkun, Stefano Massaroli, Alessandro Moro, Jimmy T. H. Smith, Ramin Hasani, Math-726 ias Lechner, Qi An, Christopher Ré, Hajime Asama, Stefano Ermon, Taiji Suzuki, Atsushi Ya-727 mashita, and Michael Poli. State-free inference of state-space models: The transfer function 728 approach, 2024. URL https://arxiv.org/abs/2405.06147. (pages 1, 19, 27). 729
- 730 Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural 731 networks, 2013. URL https://arxiv.org/abs/1211.5063. (page 2).
- 732 Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin 733 Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the 734 finest text data at scale, 2024. (pages 10, 36). 735
- 736 Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y. Fu, Tri Dao, Stephen Baccus, Yoshua 737 Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional language models, 2023. URL https://arxiv.org/abs/2302.10866. (pages 1, 3, 6, 738 19). 739
- 740 Michael Poli, Armin W Thomas, Eric Nguyen, Pragaash Ponnusamy, Björn Deiseroth, Kristian Kersting, Taiji Suzuki, Brian Hie, Stefano Ermon, Christopher Ré, Ce Zhang, and Stefano Massaroli. Mechanistic design and scaling of hybrid architectures, 2024. URL https: //arxiv.org/abs/2403.17844. (pages 5, 8, 20, 30). 744

741

742

- Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Gener-745 alization beyond overfitting on small algorithmic datasets, 2022. URL https://arxiv.org/ 746 abs/2201.02177. (page 20). 747
- 748 Hubert Ramsauer, Bernhard Schäfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Thomas 749 Adler, Lukas Gruber, Markus Holzleitner, Milena Pavlović, Geir Kjetil Sandve, Victor Greiff, 750 David Kreil, Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. 751 Hopfield networks is all you need, 2021. URL https://arxiv.org/abs/2008.02217. 752 (page 21). 753
- David W. Romero, Anna Kuzina, Erik J. Bekkers, Jakub M. Tomczak, and Mark Hoogendoorn. 754 Ckconv: Continuous kernel convolution for sequential data, 2022. URL https://arxiv. 755 org/abs/2102.02611. (page 1).

756 757 758	Olivier Roy and Martin Vetterli. The effective rank: A measure of effective dimensionality. In 2007 15th European Signal Processing Conference, pp. 606–610, 2007. (pages 3, 5).
759 760	David E. Rumelhart and James L. McClelland. <i>Learning Internal Representations by Error Propagation</i> , pp. 318–362. 1987. (page 2).
761 762	Noam Shazeer. Glu variants improve transformer, 2020. URL https://arxiv.org/abs/2002.05202. (page 8).
763 764 765	Huitao Shen. Mutual information scaling and expressive power of sequence models, 2019. URL https://arxiv.org/abs/1905.04271. (page 20).
766 767 768	Jimmy T. H. Smith, Andrew Warrington, and Scott W. Linderman. Simplified state space layers for sequence modeling, 2023. URL https://arxiv.org/abs/2208.04933. (pages 1, 5, 19, 25).
769 770 771 772	Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. URL https://arxiv.org/abs/2104.09864. (page 28, 28).
773 774	Mingjie Sun, Xinlei Chen, J. Zico Kolter, and Zhuang Liu. Massive activations in large language models, 2024. URL https://arxiv.org/abs/2402.17762. (pages 1, 20).
775 776 777	Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. Retentive network: A successor to transformer for large language models, 2023. URL https://arxiv.org/abs/2307.08621. (pages 6, 19).
778 779 780 781	Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan Salakhutdinov. Transformer dissection: A unified understanding of transformer's attention via the lens of kernel, 2019. URL https://arxiv.org/abs/1908.11775. (pages 1, 19).
782 783 784	Neehal Tumma, Mathias Lechner, Noel Loo, Ramin Hasani, and Daniela Rus. Leveraging low-rank and sparse recurrent connectivity for robust closed-loop control, 2023. URL https://arxiv.org/abs/2310.03915. (page 20).
785 786 787 788	Raf Vandebril, Marc Van Barel, and Nicola Mastronardi. A note on the representation and definition of semiseparable matrices. <i>Numerical Linear Algebra with Applications</i> , 12(8):839–858, 2005. doi: https://doi.org/10.1002/nla.455. URL https://onlinelibrary.wiley.com/doi/abs/10.1002/nla.455. (page 3).
789 790 791	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, 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).
792 793 794	Jesse Vig. A multiscale visualization of attention in the transformer model, 2019. URL https://arxiv.org/abs/1906.05714. (pages 1, 20).
795 796 797 798 799	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark, 2024. URL https://arxiv.org/abs/2406.01574. (page 21).
800 801 802	Jan C. Willems. <i>Models for Dynamics</i> , pp. 171–269. Vieweg+Teubner Verlag, Wiesbaden, 1989. ISBN 978-3-322-96657-5. doi: 10.1007/978-3-322-96657-5_5. URL https://doi.org/10.1007/978-3-322-96657-5_5. (page 2).
804 805 806	Xinyi Wu, Amir Ajorlou, Yifei Wang, Stefanie Jegelka, and Ali Jadbabaie. On the role of atten- tion masks and layernorm in transformers, 2024. URL https://arxiv.org/abs/2405. 18781. (pages 2, 20).
807 808 809	Jianlin Xia, Shivkumar Chandrasekaran, Ming Gu, and Xiaoye S. Li. Fast algorithms for hierarchi- cally semiseparable matrices. <i>Numerical Linear Algebra with Applications</i> , 17(6):953–976, 2010. doi: https://doi.org/10.1002/nla.691. URL https://onlinelibrary.wiley.com/doi/ abs/10.1002/nla.691. (page 3).

810 811 812	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming lan- guage models with attention sinks, 2024. URL https://arxiv.org/abs/2309.17453. (pages 1, 20).
013	Songlin Yang Railin Wang Yikang Shen Rameswar Panda and Yoon Kim. Gated linear atten-
814	tion transformers with hardware-efficient training 2024a URL https://arxiv.org/abs/
815 816	2312.06635. (pages 1, 6, 8, 19, 20, 33).
817	Songlin Vong Doilin Wong, Vy Thong, Vikang Chan, and Voon Vim Darollalizing linear trans
818	formers with the delta rule over sequence length 2024b LIRL https://arxiv.org/abs/
819	2406.06484. (page 8).
820	
821	Luca Zancato, Arjun Seshadri, Yonatan Dukler, Aditya Golatkar, Yantao Shen, Benjamin Bowman,
822	Matthew Trager, Alessandro Achille, and Stefano Soatto. B'mojo: Hybrid state space realizations
823	of foundation models with eidetic and fading memory, 2024. URL https://arxiv.org/
824	abs/2407.06324. (pages 1, 3).
825	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-
826	chine really finish your sentence?, 2019. URL https://arxiv.org/abs/1905.07830.
827	(page 21).
828	
829	
830	
831	
832	
833	
834	
835	
836	
837	
838	
839	
04U 9/1	
842	
843	
844	
845	
846	
847	
848	
849	
850	
851	
852	
853	
854	
855	
856	
857	
858	
859	
860	
001	
002	
000	

#### Supplementary Material CONTENTS Introduction **Related Work** Theory 3.1 **Empirical Validation of Effective State-Size** 4.14.2 Within Task-Model Analysis **Applications of Effective State-Size** 5.1 Initialization-Phase Analysis 5.2 Post-Training Analysis 5.3 State Modulation of Large Language Models Conclusion **Mid-Training Analysis** Α B **Extended Related Work C** Theoretical Background C.1 C.2 C.2.1 The Operator Realization of Linear Recurrences C.2.2 Factorizing The Operator Realization Submatrix $H_i$ .... C.2.3 The Trivial Recurrence Realization Existence of Infinite Recurrent Realizations (Proof of Theorem 3.1) . . . C.2.4 C.3 **D** Methods D.1.1 Computing ESS for SISO models Tolerance vs entropy ESS ..... D.1.2 D.1.3

918		D.2	Formulation of the Featurizers	7
919 920		D.3	Empirical Validation	С
921		D.4	ESS-Informed Featurizer Selection and Initialization Scheme	1
922		D.1	ESS Informed Peqularization 22	2
923		D.5		, -
924		D.6	ESS-Informed Model-Order Reduction	3
925		D.7	ESS Analysis for Hybrid Networks	1
927		D.8	State Modulation of Large Language Models	5
928				
929	Е	Exte	ended Experimental Results 37	7
930		E.1	Empirical Validation	7
932			E.1.1 State Collapse Continued	7
933			E.1.2 Entropy-ESS MOAR Results Continued	7
934			E 1.2 Talaranaa ESS MOAD Degulta	'n
935			E.1.5 Tolerance-ESS MQAR Results	J
936			E.1.4 Selective Copying and Compression Results	3
938			E.1.5 ESS Training Dynamics in MQAR	5
939		E.2	Initialization-Phase Analysis	5
940		E.3	Mid-Training Analysis	5
941 042		E.4	Post-Training Analysis 48	8
943		2	F 4.1 Model-Order Reduction 48	ŝ
944			E.4.2 Heteridiantian	, ,
945			E.4.2 Hybridization	1
946 947		E.5	State Modulation of Large Language Models	ł
948		E.6	Miscellaneous	1
949			E.6.1 Effective State-Size on C++ Code	1
950			E.6.2 How the Number of Prompting Shots Affects the Effective State-Size of	
951			Language Models	5
952				
953				
954				
955				
956				
957				
958				
959				
960				
961				
902				
903				
904				
303				
300				
068				
500				

## 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.

987 To motivate the idea of increasing ESS mid-training, we revisit the concept of state collapse – a 988 phenomenon that arises due to trainability issues in recurrent models (Figure 27), as discussed in 989 Section 4.2. Recall that state collapse describes a failure mode of learning in GLA and WLA which, 990 unlike LA, have learnable  $A_i$  matrices (where i denotes the index along the sequence dimension). To 991 see why this contributes to state collapse, we note that the values of the operator submatrices  $H_i$  are 992 disproportionately influenced by  $A_i$ , due to the presence of terms in the form of  $A_{i-1} \dots A_1$  for each 993 *i*. Hence, the closer  $A_i$  lies to the 0-matrix, the faster these terms decay, reducing the numerical rank 994 of  $H_i$ . We demonstrate this empirically in Figure 8a, which shows that for both GLA and WLA, 995 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. 996

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.
- 1005

1007 1008

972

973

## **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, 1015 can be formulated using a linear system representation (y = Tu) as detailed in the introductory and 1016 theoretical sections (Sections 1 and 3). For this work, we categorize these models into two types: 1017 linear systems and systems with input-varying linear operators. The key distinction between these 1018 two frameworks is that in the former, the operator T in input-invariant models is composed of fixed 1019 system parameters, whereas in the latter, the parameters are dynamically generated from the input. 1020 Linear systems encompass both linear time-varying (LTV) and linear time-invariant (LTI) systems. 1021 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 1023 minimal state-size of the model (DeWilde & van der Veen, 1998). Massaroli et al. (2023) explored 1024 methods to reduce the order of models by leveraging the Hankel matrix. Notably, the submatrix  $H_i$ 1025 (defined in Theory 3.2) exhibits a Hankel structure in LTI models and provides per-sequence-index information. In this work, however, we do not explore Hankel matrices further, as they are not easily generalizable to LTV systems.

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

1035

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 effective at capturing the inner workings of causal sequence models, which is the primary focus of 1042 this work. 1043

Intrinsic tools, conversely, directly visualize the model's internal mechanisms. A recently popular framework known as mechanistic interpretability provides one such example (Power et al., 2022).
 Mechanistic interpretability involves dissecting complex models to understand how specific components contribute to the model's overall behavior (Cammarata et al., 2020). Unlike our work, mechanistic interpretability does not target the operator view of the model but instead emphasizes 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 structures. However, raw visualizations are largely qualitative and oftentimes do not provide the 1056 quantitative metrics necessary for effectively evaluating a model's internal mechanisms – a gap we 1057 aim to address in this work. 1058

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.

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

1067

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 et al., 2024; Bhattamishra et al., 2023; Elhage et al., 2021; Olsson et al., 2022b). Training models 1071 on these tasks is efficient, making them well-suited for exploring a large space of tasks and models, 1072 as demonstrated in several prior works (Dupont et al., 2019; Arora et al., 2024; Fu et al., 2023). In 1073 this work, we investigate the effective state-size across a subset of the MAD tasks: multi-query as-1074 sociative recall (MQAR), selective copying, and compression, varying the difficulty of each to gain 1075 a nuanced understanding of how effective state-size evolves across these task landscapes. 1076

Among the synthetic tasks we examine, MQAR stands out in particular. Proposed by Arora et al.
(2023), MQAR was designed to bridge the gap between synthetic and real language tasks explained
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

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

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

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

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

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

The effective state-size analysis presented in this work reveals that strong recall performance as measured by bigram perplexity in language modeling tasks depends not only on memory capacity, but also on a model's ability to modulate its state-size within a given context.

- 1109 1110 1111 1112 1113 1113
- 1115
- 1116
- 1117
- 1118
- 1119
- 1120
- 1121
- 1122
- 1123 1124
- 1125
- 1126
- 1127 1128
- 1129
- 1130
- 1131
- 1132

С 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  $u_i$  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:  $s_0 = 0$  $s_1 = B_0 u_0$  $s_2 = B_1 u_1 + A_1 (B_0 u_0)$  $s_3 = B_2 u_2 + A_2 (B_1 u_1 + A_1 (B_0 u_0))$  $s_i = \left(\sum_{j=0}^{i-1} \left[\prod_{k=i-1}^{j+1} A_k\right] B_j u_j\right),\,$ (C.2.1)  $y_i = C_i \left( \sum_{j=1}^{i-1} \left[ \prod_{k=i-1}^{j+1} A_k \right] B_j u_j \right) + D_i u_i,$ (C.2.2) which corresponds to the operator:  $T_{ij} = \begin{cases} 0 & i < j \\ D_i & i = j \\ C_i A_{i-1 \cdot j+1} B_j & i > j \end{cases}$ (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  $(n_i \ge \operatorname{rank}(H_i))$ : 

$$T_{i:,:i-1} \equiv H_i = \begin{bmatrix} C_i & & \\ & \ddots & \\ & & C_{\ell-1} \end{bmatrix} \begin{bmatrix} A_{i-1;1} & \cdots & I \\ \vdots & \ddots & \vdots \\ A_{\ell-2;1} & \cdots & A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} B_0 & & \\ & \ddots & \\ & & B_{i-1} \end{bmatrix}$$

$$= \begin{bmatrix} C_i & & \\ & \ddots & \\ & & C_{\ell-1} \end{bmatrix} \begin{bmatrix} I \\ A_i \\ A_{i+1;i} \\ \vdots \\ A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{i-1;1} & A_{i-1;2} & \cdots & A_{i-1} & I \end{bmatrix} \begin{bmatrix} B_0 & & \\ & \ddots & \\ & & B_{i-1} \end{bmatrix}$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{i-1;1}B_0 & A_{i-1;2}B_1 & \cdots & A_{i-1}B_{i-2} & B_{i-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{i-1;1}B_0 & A_{i-1;2}B_1 & \cdots & A_{i-1}B_{i-2} & B_{i-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{i-1;1}B_0 & A_{i-1;2}B_1 & \cdots & A_{i-1}B_{i-2} & B_{i-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{i-1;1}B_0 & A_{i-1;2}B_1 & \cdots & A_{i-1}B_{i-2} & B_{i-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{i-1;1}B_0 & A_{i-1;2}B_1 & \cdots & A_{i-1}B_{i-2} & B_{i-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{i-1;1}B_0 & A_{i-1;2}B_1 & \cdots & A_{i-1}B_{i-2} & B_{i-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{i-1;1}B_0 & A_{i-1;2}B_1 & \cdots & A_{i-1}B_{i-2} & B_{i-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{i-1;1}B_0 & A_{i-1;2}B_1 & \cdots & A_{i-1}B_{i-2} & B_{i-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{i-1;1}B_0 & A_{i-1;2}B_1 & \cdots & A_{i-1}B_{i-2} & B_{i-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{\ell-1;1}B_0 & A_{\ell-1;2}B_1 & \cdots & A_{\ell-1}B_{\ell-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{i+1}A_i \\ C_{i+2}A_{i+1;i} \\ \vdots \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{\ell-1;1}B_0 & A_{\ell-1;2}B_1 & \cdots & A_{\ell-1}B_{\ell-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{\ell-1}A_{\ell-2;i} \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{\ell-1;1}B_0 & A_{\ell-1;2}B_1 & \cdots & A_{\ell-1}B_{\ell-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{\ell-1}A_{\ell-2;i} \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{\ell-1;1}B_0 & A_{\ell-2;i} \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_{\ell-1}A_{\ell-2;i} \\ C_{\ell-1}A_{\ell-2;i} \end{bmatrix} \begin{bmatrix} A_{\ell-1;1}B_1 & \cdots & A_{\ell-1;i}B_{\ell-1} \\ C_{\ell-1}A_{\ell-1} \end{bmatrix} \equiv \mathcal{O}_i C_i.$$

$$= \begin{bmatrix} C_i \\ C_$$

allows us to construct a more efficient realization of the recurrence:

- We can minimally factorize the causal projection as  $C_i = L_i R_i$ , where  $L_i \in \mathbb{R}^{n_i \times r}$ and  $R_i \in \mathbb{R}^{r \times di}$ , with  $r = \operatorname{rank}(\mathcal{C}_i)$ .
- 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:

$$s_{i+1} = \begin{bmatrix} I_{(di)} \\ 0_{(d)} \end{bmatrix} s_i + \begin{bmatrix} 0_{(di)} \\ I_{(d)} \end{bmatrix} u_i,$$
  
$$y_i = \begin{bmatrix} T_{i,0} & T_{i,1} & \cdots & T_{i,i-1} \end{bmatrix} s_i + T_{i,i} u_i.$$
 (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  $C_i = I_{di}$ .

1242 C.2.4 EXISTENCE OF INFINITE RECURRENT REALIZATIONS (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  $\mathcal{C}_i \in \mathbb{R}^{n_i \times di}$ , 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 C<sub>i</sub> = O<sub>i[:d-1]</sub> to obtain (C<sub>i</sub>)<sub>i∈[1,ℓ]</sub>, as given the assumptions above, the first set of rows of O<sub>i</sub> linearly projects s<sub>i</sub> onto y<sub>i</sub> − D<sub>i</sub>u<sub>i</sub>, which is identical to C<sub>i</sub> in Equation (1).
  - (b) Set B<sub>i-1</sub> = C<sub>i[:,-d:]</sub> to obtain (B<sub>i</sub>)<sub>i∈[ℓ-1]</sub>, for which the identity can be obtained by deconstructing the input-state projection matrix C<sub>i</sub> and equating its assumed state s<sub>i</sub> with Equation (1).

S

(c) Using the same state-dynamics equation, we could equate the assumed stateprojection matrices with each other, obtaining (A<sub>i</sub>)<sub>i∈[1,ℓ-1]</sub>:

$$s_{i+1} = A_i s_i + B_i u_i$$

$$C_{i+1} u_{:i} = A_i C_i u_{:i-1} + C_{i+1[:,-d:]} u_i$$

$$C_{i+1[:,:-d-1]} u_{:i-1} = A_i C_i u_{:i-1}$$

$$A_i = C_{i+1[:,:-d-1]} C_i^+.$$
(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.

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

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

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

1294 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 1295  $r = \text{rank of } H_i \in \mathbb{R}^{m \times n}$ . Theorem 3.2 demonstrates that these factors realizes a recurrence with minimal state-size.

# 1296 C.3 MINIMAL RECURRENT REALIZATION (PROOF OF THEOREM 3.2)

**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.

**1301** *Proof.* The proof of Theorem 3.1 shows that the operator submatrices  $H_i$  can be decomposed ar- **1302** bitrarily into two state-projection matrices,  $\mathcal{O}_i$  and  $\mathcal{C}_i$ , whose inner product dimension defines the **1303** state size of its recurrent realization at sequence index *i*. By the rank-nullity theorem, rank $(H_i)$  **1304** represents the minimum inner product dimension of any such state-projection matrices and thus **1305** corresponds to the minimally realizable state size of the operator *T* at sequence index *i*.

1306

1298

1299

1300

- <sup>1307</sup> D METHODS 1308
- 1309 D.1 Additional Details on Computing ESS 1310
- 1311 D.1.1 COMPUTING ESS FOR SISO MODELS

In this section, we provide additional details on the distinction between computing ESS in single-input single-output (SISO) models like S4 and multi-input multi-output (MIMO) models like S5.

Recall in Sections 3 and C.1, we introduced flattened notation, as it offers a general framework for formulating a wide range of operators and recurrences. Namely, a MIMO layer like S5 (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 a SISO layer like S4 (Gu et al., 2022a) can, 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; for S5, it is dense.

Note that since all of the models in our experiments are SISO, we can compute the ESS independently for each channel using the standard operator formulation  $T^{\alpha\alpha} \in \mathbb{R}^{\ell \times \ell}$ , where  $\alpha \in [d]$ . 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.

1327

D.1.2 TOLERANCE VS ENTROPY ESS

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. Moreover, in contrast to the tolerance-based metric which is discrete, entropy-ESS can assume continuous values ranging from 1 to  $|\Sigma_i|$ .

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. Finally, we note that for the analyses presented in Section 4, when computing ESS, we average across 8 samples (batch-size). For the rest of the ESS analyses, we average across 32 samples.

D.1.3 PyTorch Implementation

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

```
1346
1 import torch
1347
2
1348
3 def T2H_i(T, i, d=1):
1349
4
Extract H_i from T.
```

```
1350
     6
1351
     7
            Args:
1352 8
               - T: Flattened operator with shape [..., d*L, d*L].
1353 9
                - i: Index of H (H i) to retrieve.
               - d: Block size for multi-channel flattened operator
1354 <sup>10</sup>
            representation (default is 1).
1355
     11
1356 12
            Returns:
1357 13
               - H_i: Submatrix of the operator at index i.
1358 14
1359<sup>15</sup>
            return T[...,d*i:,:d*i]
1360<sup>16</sup>
    17 @torch.no_grad()
1361 <sub>18</sub> def T2Ss(T, d=1):
1362 19
            .....
1363 20
            Converts an operator into a list of singular values (Ss).
1364 <sup>21</sup>
1365<sup>22</sup>
            Args:
     23
                - T: Flattened operator with shape [..., d*L, d*L]
1366 24
                - d: Block size for multi-channel flattened operator
1367
           representation (default is 1).
1368 25
1369 <sup>26</sup>
            Returns:
1370<sup>27</sup>
              - Ss: A list of singular values for each sequence index in T.
            .....
     28
1371 29
            seqlen = T.size(-2)//d
1372 30
            Ss = []
            for i in range(1, seqlen):
1373 31
               H_i = T2H_i(T, i, d)
1374 <sup>32</sup>
1375 <sup>33</sup>
                _, S_i, _ = torch.svd(H_i)
               Ss.append(S_i)
     34
1376 35
            return Ss
1377 36
1378 37 @torch.no_grad()
1379 38 def Ss2ToleranceESS(Ss, tol=1e-4):
1380 <sup>39</sup>
            ....
            Computes the tolerance-ESS from the list of singular values.
     40
1381 41
1382 42
            Args:
               - Ss: List of singular values.
1383 43
1384 <sup>44</sup>
                - tol: Tolerance value.
1385 45
            Returns:
     46
1386 47
               - tolerance-ESS
1387 48
            ....
            ranks = []
1388 49
            for SV in Ss:
1389 50
1390 <sup>51</sup>
               rank = torch.sum(SV>=tol, dim = -1)
     52
                ranks.append(rank)
1391 53
            ranks = torch.stack(ranks, dim=-1)
1392 54
            return ranks
1393 55
1394 56 @torch.no_grad()
1395 57 def Ss2EntropyESS(Ss, clip=1e-12):
     58
1396 59
            Computes the entropy-ESS from the list of singular values.
1397 60
1398 61
            Args:
1399 <sup>62</sup>
                - Ss: List of singular values.
                - clip: clips probabilities below this value avoiding numerical
1400 <sup>63</sup>
            instabilities when the probabilities are too numerically close to 0.
1401 64
1402 65
            Returns:
1403 66

    entropy-ESS

            ....
    67
```

```
1404
             ranks = []
      68
1405
               for SV in Ss:
      69
1406 70
                     p = SV/SV.sum(dim=-1)[..., None]
1407 71
                      p = torch.clip(p, clip)
                    H = -torch.sum(p * torch.log(p), dim=-1)
1408 <sup>72</sup>
                     rank = torch.exp(H)
1409
      74
                     ranks.append(rank)
1410
      75
                ranks = torch.stack(ranks, dim=-1)
1411
      76
                return ranks
1412
          Example usage (Python-pseudocode):
1413
1414
      1 >>> out = model(u, output_attentions=True)
      2 >>> # T shape: [bs, layers, heads, len, len]
1415
       3 >>> T = out.attention_matrix
1416
       4 >>> Ss = T2Ss(T) # List of singular values
1417
       5 >>> # ESS shape [bs, layers, heads, len-1]
1418
       6 >>> ESS = Ss2ToleranceESS(Ss, tol=1e-3)
       7 >>> mean_ESS = torch.mean(ESS)
1419
1420
          We note that calculating the effective rank may cause numerical instability when p_i^m approaches 0
1421
          due to the logarithmic term. This is partially mitigated by clipping the normalized singular values,
1422
          as shown above.
1423
1424
          D.2 FORMULATION OF THE FEATURIZERS
1425
1426
          In this section, we first establish the equivalence between linear attention and state-space models,
1427
          then proceed with formulating the LIVs discussed in this paper.
1428
1429
          Linear attention and state-space model equivalence. We begin by demonstrating that linear
1430
          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.
1431
1432
          A single linear attention head with dimension d/h, typically formulated as
1433
                                                             y = qk^T v.
                                                                                                                     (D.2.1)
1434
          in which q, k, v \in \mathbb{R}^{\ell \times d/h} are input features. They can be reformulated as recurrences with matrix-
1435
          valued states s_i \in \mathbb{R}^{d/h \times d/h} as follows (Katharopoulos et al., 2020):
1436
                                                         s_i = s_{i-1} + k_i v_i^T
1437
                                                                                                                     (D.2.2)
1438
                                                             y_i = q_i^T s_i,
1439
          Without loss of generality, applying column-major flattening to the matrix-valued state and treating
1440
          v_i as the input u_i, the recurrence can be formulated like Equation (1), by setting A_i = I_{(d/h)^2} and:
1441
                                            0
1442
1443
         B_{i-1} = \begin{vmatrix} \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 \end{vmatrix}, \quad C_i = \begin{bmatrix} q_i^1 & \cdots & q_i^{d/h} & 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 & 0 & \cdots & q_i^1 & \cdots & q_i^{d/h} \end{bmatrix}.
1444
1445
1446
1447
1448
1449
1450
1451
1452
                     0
1453
                                                                                                                     (D.2.3)
1454
1455
          Notice that each individual channel forms a single-input-single-output (SISO) recurrence (like many
1456
          of the SSM architectures including S4, S6, Mamba2, and more [27; 26; 16; 48; 49]), as there is no
1457
          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:

$$s_{i+1} = A_i s_i + B_i f_u(u_i)$$
  

$$y_i = C_i^T s_i + D_i f_u(u_i).$$
(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):

$$A^{k} = I, \quad B^{k}_{i-1} = \operatorname{RoPE}(W^{k}_{B}u_{i}), \quad C^{k}_{i} = \operatorname{RoPE}(W^{k}_{C}u_{i}), \quad f_{u}(u_{i}) = W^{k}_{u}u_{i},$$
(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):

$$A_{i-1}^{k} = \text{diag}(\text{sigmoid}(W_{A_{2}}^{k}W_{A_{1}}u_{i})^{1/\beta}),$$
  

$$B_{i-1}^{k} = W_{B}^{k}u_{i}, \quad C_{i}^{k} = W_{C}^{k}u_{i}, \quad f_{u}(u_{i}) = W_{u}^{k}u_{i},$$
(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):

$$A^{k} = \text{diag}(\text{sigmoid}(A^{k})^{1/\beta}),$$
  

$$B^{k}_{i-1} = W^{k}_{B}u_{i}, \quad C^{k}_{i} = W^{k}_{C}u_{i}, \quad f_{u}(u_{i}) = W^{k}_{u}u_{i},$$
(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):

$$\hat{B}_i^k = \operatorname{RoPE}(W_B^k u_i), \quad \hat{C}_i^k = \operatorname{RoPE}(W_C^k u_i),$$
  

$$T^k = \operatorname{softmax}(\hat{C}^k (\hat{B}^k)^T), \quad f_u(u_i) = W_u^k u_i,$$
(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):

$$\Delta^{c} = \text{softplus}(W^{c}_{\Delta}u_{i} + b^{c}), \quad A^{c}_{i-1} = \text{diag}(\exp(-\hat{A}\Delta^{c})), \\ B^{c}_{i-1} = \Delta^{c}W_{B}u_{i}, \quad C_{i} = W_{C}u_{i},$$
(D.2.9)

where  $\hat{A} \in \mathbb{R}^n$  is initialized to  $\begin{bmatrix} 1 & 2 & \cdots & n \end{bmatrix}^T$ , 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*.

- 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:  $A_{i-1}^{h} = \operatorname{diag}(\exp(-[1/\alpha \quad 2/\alpha \quad \dots \quad n/\alpha]^{T} \odot \operatorname{softplus}(W_{A_{2}}^{h}W_{A_{1}}u_{i}))),$ (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  $(\begin{bmatrix} 1 & 2 & \cdots & n \end{bmatrix}^T)$  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

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

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

1577 1578

 $4 * \text{num kv pairs} \le \text{seq len}$ 

and the selective copying task requires that

1580 1581

1609

1610

2 \* num tokens to copy + 1 < seq len

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

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

Configuration	Value(s)
Tasks	MQAR, selective copying, compression
Num. key-value pairs	8, 16, 32, 64, 128
Num. tokens to copy	8, 16, 32, 64, 128
Vocab size (compression)	8, 16, 32, 64, 128
Vocab size (MQAR and selective copying)	8192
Sequence length	64, 128, 256, 512, 1024, 2048
Model (featurizer)	GLA, LA, WLA, SA
Model width	64, 128, 256, 512
Number of heads	4, 8
Optimizer	AdamW
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
C I	

Table 1: Set of hyperparameters for the task-model sweep.

1611 Regarding the post-hoc analysis performed on the sweep, we note the following:

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

1620		• Since we analyze the recurrent models across the average TSS dimension, we compute	
1621		average ESS in the plots presented in Section 4.1 in order to compare ESS and TSS as	
1622		proxies for performance. Similarly, since we analyze the SA models across the total	
1623		dimension, we compute total ESS for those plots. However, we note that plots for both the	
1624		average/total ESS and TSS are presented in Section E.1.	
1625		• When we marginalize across dimensions, we average across all models in that bucket of	
1626		task-model space. For example, in Figure 4, for each (TSS, kv) pair, we average over the	
1627		correlations of all models that correspond to that pair. Note, however, that we never average	
1628		across tasks (i.e. MQAR, selective copying, compression) or featurizers (i.e. GLA, LA,	
1629		WLA, SA).	
1630		• When we compute cross-model correlations (Figure 2) for SA, we filter out models which	
1631		have an accuracy $> 0.95$ . This is done in order to observe meaningful variation as a func-	
1632		tion of (total ESS)/kv and (total TSS)/kv since many of the SA models obtain an accuracy	
1633		of 1.	
1634		• When we compute within-model correlations (Figure 4) for MOAR we drop epoch 0 from	
1635		the computation since we observe a phase at the start of training in which ESS tends to	
1636		decrease, but accuracy does not change. We elaborate on this phenomenon in Section E.1	
1637		and hope to characterize it further in future work.	
1630		• Regarding the task-adjusted forms of FSS and TSS which in the case of MOAR are com-	
1640		puted by normalizing the raw ESS value by the number of ky-pairs in the task, we note	
1640		that this normalization factor is critical for observing the cross task-model correlations pre-	
1641		sented in Figure 2. In particular, in Figure 9, we find that correlations across the task-model	
1642		space break down when examining the unnormalized ESS. This points to the higher-level	
1643		notion that ESS is expected to scale with the memory demands of the task.	
1645		• We interpret the state utilization of a model, which is given by ESS/TSS, as a proxy for what	
1646		portion of the memory capacity of the network is realized in practice. By definition, state	
1647		utilization takes on values ranging continuously from 0 to 1. Recall that a state utilization	
1648		near 1 is indicative of state saturation.	
1649		• While for most of the ESS analysis conducted on the sweep we use the entropy-ESS, we	
1650		note that for the state utilization plot presented in Figure 5, we use the tolerance-ESS with	
1651		a tolerance level set at 1e-3. We do this because we find that entropy-ESS fails to capture	
1652		the state collapse phenomenon. This is because state collapse is primarily dictated by	
1653		the magnitude of the singular values, as opposed to the relative decay rate of the entire	
1654		spectrum. In particular, if all of the singular values are close to 0, the layer is likely failing	
1655		to learn an expressive state, resulting in poor performance. Due to the normalization applied to the gradient the entropy ESS matrix may respectively present this state as having a high	
1656		effective rank: however in practice we know that this is a misrepresentation of the true	
1657		dynamics. Tolerance-ESS in contrast appropriately captures the dynamics of the state	
1658		with respect to the norm of the operator. Because of this, whenever we analyze ESS as it	
1659		pertains to state collapse (e.g. Figure 8a), we present the tolerance-ESS instead.	
1660			
1661	D.4	ESS-INFORMED FEATURIZER SELECTION AND INITIALIZATION SCHEME	

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

Table 2: Default GLA hyperparameters.

Value Configuration 128 Model width State expansion (d\_state) 16

<sup>a</sup>For GLA-S6.

1666

1667

1668

1669 1670

1671

Table 3: Default S6 hyperparameters.

<sup>1672</sup>  ${}^{b}K$ -expansion is used to vary TSS in the featurizer 1673 experiments.

1674	Configuration   Val	ue
10/5	Sequence length 204	18
10/0	Num. KV Pairs 12	8
10//	KV Dist. Const.	1
10/0 1670	Optimizer Adam	nW <sup>a</sup>
10/9	Learning Rate 0.0	02
1080	Weight Decay 0.	1
1661	Batch Size 64	1
1682	Epochs 70	)
1683	Steps Per Epoch 200	)0
1684	Num. Training Samples 128	3k
1685	Num. Testing Samples 6.4	rk
1686	Vocabulary Size 819	92
1687		
1688	Table 4: Default MQAR task settings	employed
1689	neriments in Section 5.1	Lation ex-
1690	perments in section 3.1.	
1691	<sup>a</sup> Loshchilov & Hutter (2019)	
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		
1793		
170/		
1729		
1726		
1720		
1121		

## 1728 D.5 ESS-INFORMED REGULARIZATION 1729

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

1732	· · · · · · · · · · · · · · · · · · ·		
1733	Configuration	Value	
1734	Sequence length	4096	
1735	Num. KV Pairs	128	
1736	KV Dist. Const.	0.1	
1737	Optimizer	AdamW <sup>a</sup>	
1738	Learning Rate	0.002	
1739	Weight Decay	0.1	
1740	Batch Size	64	
1740	Epochs	70	
1741	Steps Per Epoch	2000	
1742	Num. Training Samples	128k	
1743	Num. Testing Samples	6.4K	
1744	Vocabulary Size	8192	
1740	Niodel Width	128	
1740	Nuill. lieaus	0	
1747	Table 5. MOAD took acting	and mode	1 by
1748	nerparameters, employed through	s and mode	mid
1749	training experiments in Section	A	inid-
1750			
1/51	<sup>a</sup> Loshchilov & Hutter (2019)		
1752			
1753	Regarding the regularization scheme itself, since we ex-	amine mod	els with two sequence mixing
1754	layers, we explore the following strategies: regularizing l	both layers,	only regularizing the first layer
1755	and only regularizing the second layer. Empirically, we t	find that only	y regularizing the second layer
1756	performs the best and is thus the result presented in Figure	re <mark>8b</mark> . We ela	aborate on why this is the most
1757	successful strategy in Section E.3.		
1758			
1759	D.6 ESS INFORMED MODEL OPDER PEDUCTION		
1760	D.0 ESS-INFORMED MODEL-ORDER REDUCTION		
1761	The teacher models used in the distillation experiments a	re 2-laver Gl	A models (Vang et al. 2024a)
1762	with dimension $-128$ and TSS $-256$ (num heads $-8$	and expand	k = 16) We checkpointed the
1763	models every 10 epochs while training on MOAR across	different ta	k = 10, we encerpointed the sk difficulties. The task ranges
1764	are given as follows:	, uniorent tu	sk diffedities. The task funges
1765	6		
1766	• Sequence length: [512, 1024, 2048]		
1767			
1768	• Number of Key-Value Pairs: [64, 128]		
1769			
1770	Other settings follow the defaults shown in Table 4. For	r each task o	difficulty pair, we repeated the
1771	training run with three different seeds. For each teacher n	nodel checkp	point, both layers were distilled
1772	independently with student models of different state-sizes	s (16, 32, 64	, and 128). Distillation settings
1773	are shown in Table 6.		
1774	The ESS metric in Figures 6d, 30, and 31 was compute	d by taking	the minimum across the batch
1775	and channels, evaluated at the mid-point of the sequence	$(\ell/2)$ . Using	the mid-point of the sequence
1776	as a summary statistic was done to reduce compute. The	e midpoint i	n particular was chosen as it is
1777	the point in the sequence at which $H_i$ has the greatest d	imensions, 1	retaining the largest amount of
1778	information from the original operator. Note that in pra	ctice, where	a large task space isn't being
1779	tested to validate our approach, a more thorough comp	outation of I	2SS across sequence length is
1780	teasible. Other marginalization approaches such as taking	g the maxim	um or average across the batch
1781	and channels also snow similar trends, but we found tal	king the min	infum to dest demonstrate the
	uona.		

1782		Configuration	Value	
1783		Ontimizor	AdamW	
1784		Batch Size	Adamw	
1785		Learning Rate	0.001	
1786		Weight Decay	0.001	
1787				
1788		Dropout (Operator)		
1789			0.2	
1790		Training Steps (Activation)	) 3200	
1791		Dropout (Activation)	0.2	
1792	Table & Distill	tion acttings used for the res	ulto mascanto	ad in Section 5.2
1793	Table 0. Distina	ation settings used for the les	suns presente	20 III Section 5.2.
1705				
1796				
1797				
1798	D.7 ESS ANALYSIS FOR	Hybrid Networks		
1799				
1800	In our ESS analysis applied	l to hybrid networks, we re	strict our sc	cope to GLA-SA hybrids. In
1801	particular, we explore the fol	lowing two settings:		-F
1802		0 0		
1803				
1804	• 8-layer hybrid netw	orks in which 4 layers are s	sequence miz	xers (i.e. one of GLA or SA)
1805	and 4 layers are cha	annel mixers (i.e. MLPs). W	e exhaust all	l possible hybrid networks (of
1806	which there are 16)	and perform post-training, p	er-layer ESS	analysis on the networks. We
1807	train these hybrid m	nodels on MQAR with task-n	nodel setting	s given below in Table 7.
1808				
1809				
1810	<ul> <li>16-layer hybrid network</li> </ul>	works in which 8 layers are	sequence mi	xers (i.e. one of GLA or SA)
1811	and 8 layers are cha	innel mixers (i.e. MLPs). He	ere, we explo	bre all combinations of hybrid
1812	training per layer	v the Jamba hybridization po	We train th	et al., 2024) and perform post-
1813	with task-model set	tings given below in Table 8		lese hybrid models on MQAK
1814	with task model set		•	
1815				
1816			X 7 1	
1817		Configuration	Value	
1818		Sequence length	2048	
1819		Num. KV Pairs	512	
1820		KV Dist. Const.	0.1	
1821		Upumizer Learning Rate	Adam W $^{4}$	
1822		Weight Decay	0.002	
1823		Batch Size	64	
1824		Epochs	70	
1025		Steps Per Epoch	2000	
1020		Num. Training Samples	128k	
1027		Num. Testing Samples	6.4k	
1020		Vocabulary Size	8192	
1029		Model width	64	
1831		num, neaus	4	
1832	Tak	ale 7. Default MOAD teals a	attings amou	oved
1833	1aC thr	ne 7. Default MQAK task so	sungs emplo	con-
1834	duc	ted in the first setting descri	bed above.	
1835				
-				

<sup>&</sup>lt;sup>*a*</sup>Loshchilov & Hutter (2019)

		X X 1		
	Configuration	Value		
	Sequence length	4006		
	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. Testing Samples	128K 6.4k		
	Vocabulary Size	8192		
	Model width	16		
	Num. heads	2		
	Table 8: Default MQAR task s	ettings empl	loyed	
	throughout the hybridization	experiments	con-	
	ducted in the second setting dea	scribed abov	e.	
	<sup>a</sup> Loshchilov & Hutter (2019)			
Results fo	r these experiments can be found in Section E.	4.2.		
D 0 0-		ODDIC		
D.8 ST.	TE MODULATION OF LARGE LANGUAGE M	ODELS		
D.8 ST.	The Modulation of Large Language M	ODELS	mly concepted conter	
D.8 ST. State mo	ATE MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follo	ODELS wing rando	mly generated senter	ices were
D.8 ST. State mo used to st guage mo	ATE MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follo udy the effects of separator tokens on state m dels	ODELS wing rando odulation in	mly generated senter open-weights pre-tra	ices were ined lan-
D.8 ST. State mo used to st guage mo	ATE MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow udy the effects of separator tokens on state m dels.	ODELS wing rando odulation in	mly generated senter open-weights pre-tra	ices were lined lan-
D.8 ST. State mo used to st guage mo <bos>M raffes ar</bos>	ATE MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follo udy the effects of separator tokens on state m dels.	ODELS wing rando odulation in ded into a re	mly generated senter open-weights pre-tra	ep> Gi-
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite</bos>	ATE MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state m dels. Mangoes are rich in vitamin C and can be blend to the tallest mammals on Earth due to their long rs from the 1940s<	ODELS wing rando odulation in ded into a re necks and le Spot is a sia	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect.</sep></s 	tices were ined lan- ep> Gi- s vintage n raging
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund</bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep></sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect ont storm that has bee</sep></s 	aces were ined lan- ep> Gi- s vintage n raging
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund</bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep></sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect ont storm that has bee</sep></s 	aces were ined lan- ep> Gi- s vintage n raging
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund</bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red a reds of years<sep></sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect. Int storm that has bee</sep></s 	ep> Gi- s vintage n raging
D.8 ST. State mo used to st guage mo <bos>M raffes ara typewrite for hund</bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state m dels. Mangoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red of reds of years<sep></sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect. int storm that has bee</sep></s 	ep> Gi- s vintage n raging
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund</bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state m dels. Mangoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red of reds of years<sep> dulation on custom-trained 1B models. For</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect. ant storm that has bee</sep></s 	ep> Gi- s vintage n raging
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund State mo we used le</bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state m dels. Mangoes are rich in vitamin C and can be blend to the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect. ont storm that has been n-trained 1B language</sep></s 	e models, equences.
D.8 ST. State mo used to st guage mo <bos>M raffes ard typewrite for hund State mo we used le A collecti</bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow dudy the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w on of randomly generated sentences is shown b</sep></sep>	ODELS wing rando: odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below:	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect ont storm that has been n-trained 1B language</sep></s 	e models, equences.
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund State mo we used le A collecti</bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow duly the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red of reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w on of randomly generated sentences is shown be the deep blue ocean teeping with an extraordir</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc velow:	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect. ont storm that has bee n-trained 1B language ernible with shorter so</sep></s 	e models, equences.
D.8 ST. State mo used to st guage mo <bos>M raffes ard typewrite for hund State mo we used le A collecti <bos>T plankton</bos></bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow duty the effects of separator tokens on state models. Mangoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red of reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w to no f randomly generated sentences is shown be the deep blue ocean, teeming with an extraording to the largest whales, stretches out infinitely to</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below: ary array of owards the h	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect. int storm that has bee n-trained 1B language ernible with shorter se</sep></s 	e models, equences.
D.8 ST. State mo used to st guage mo <bos>M raffes ara typewrite for hund State mo we used le A collecti <bos>T planktom expanse</bos></bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow duty the effects of separator tokens on state m dels. Mangoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red of reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w on of randomly generated sentences is shown be the deep blue ocean, teeming with an extraordir. to the largest whales, stretches out infinitely to that has captivated the imaginations of explore</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc pelow: eary array of owards the h rs, scientists	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect. int storm that has bee n-trained 1B language ernible with shorter se f marine life, from the porizon, a vast and my and poets for centur</sep></s 	e models, equences. smallest osterious cies, hid-
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund State mo we used le A collecti <bos>T plankton expanse ing withi</bos></bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow dudy the effects of separator tokens on state m dels. Mangoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w on of randomly generated sentences is shown be the deep blue ocean, teeming with an extraording to the largest whales, stretches out infinitely to that has captivated the imaginations of explore in its depths secrets yet to be discovered and sto</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc pelow: wary array of owards the h rs, scientists ries yet to be	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect. ant storm that has bee n-trained 1B language ernible with shorter se f marine life, from the porizon, a vast and my s, and poets for centure told<sep> In a buss</sep></sep></s 	e models, equences. swallest esterious ries, hid- tling city
D.8 ST. State mo used to st guage mo <bos>M raffes ard typewrite for hund State mo we used le A collecti <bos>T plankton expanse ing withi where sk</bos></bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state m dels. Mangoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w on of randomly generated sentences is shown be the deep blue ocean, teeming with an extraording to the largest whales, stretches out infinitely to that has captivated the imaginations of explore in its depths secrets yet to be discovered and sto pascapers tower over narrow streets filled with its second second se</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below: ary array of owards the h rs, scientists ries yet to be the constant	mly generated senter open-weights pre-tra efreshing smoothie <s gs<sep> She collect. ant storm that has bee entited 1B language entible with shorter se f marine life, from the porizon, a vast and my s, and poets for centure told<sep> In a bust hum of cars and the c</sep></sep></s 	e models, equences. swallest esterious ries, hid- tling city hatter of
D.8 ST. State mo used to st guage mo <bos>M raffes ard typewrite for hund State mo we used lo A collecti <bos>T plankton expanse ing withi where sk pedestrid</bos></bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow dudy the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w for of randomly generated sentences is shown to the deep blue ocean, teeming with an extraordir to the largest whales, stretches out infinitely to that has captivated the imaginations of explore is depths secrets yet to be discovered and sto were over narrow streets filled with the ns, a small café, nestled between two imposing</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below: eary array of owards the h rs, scientists ries yet to be the constant g buildings,	mly generated senter open-weights pre-tra- efreshing smoothie <s gs<sep> She collect. ant storm that has bee metrained 1B language ernible with shorter set f marine life, from the porizon, a vast and my s, and poets for centure told<sep> In a buss hum of cars and the co offers a quiet refuge j</sep></sep></s 	e models, equences. smallest osterious cies, hid- bling city hatter of for those
D.8 ST. State mo used to st guage mo <bos>M raffes ara typewrite for hund State mo we used le A collecti <bos>T planktom expanse ing withi where sk pedestria seeking o</bos></bos>	And the effects of separator tokens on state m dulation of open-weight models. The following the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend to the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w for of randomly generated sentences is shown be the deep blue ocean, teeming with an extraordin to the largest whales, stretches out infinitely to that has captivated the imaginations of explore in its depths secrets yet to be discovered and stor we can be to be a stretches to be imposing to moment of peace, with the comforting aroma of the models. The top of to</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc pelow: eary array of owards the h rs, scientists ries yet to be the constant g buildings, of freshly bre	mly generated senter open-weights pre-tra- efreshing smoothie <s gs<sep> She collect. ont storm that has bee n-trained 1B language ernible with shorter se f marine life, from the porizon, a vast and my s, and poets for centure told<sep> In a busis hum of cars and the co offers a quiet refuge j ewed coffee and the se</sep></sep></s 	e models, equences. smallest osterious cies, hid- tling city hatter of for those oft sound
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund State mo we used le A collecti <bos>T plankton expanse ing withi where sk pedestrice seeking c of jazz n</bos></bos>	ATE MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow duty the effects of separator tokens on state models. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w for of randomly generated sentences is shown be the deep blue ocean, teeming with an extraordir to the largest whales, stretches out infinitely to that has captivated the imaginations of explore in its depths secrets yet to be discovered and sto performed as small café, nestled between two imposing moment of peace, with the comforting aroma of usic playing in the background, creating a co</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below: wards the h rs, scientists ries yet to be the constant g buildings, of freshly bre zy ambiance	mly generated senter open-weights pre-tra- efreshing smoothie <s gs<sep> She collect. and storm that has bee entited 1B language entible with shorter se f marine life, from the porizon, a vast and my s, and poets for centure told<sep> In a busis hum of cars and the co offers a quiet refuge j weed coffee and the so that feels like a wor</sep></sep></s 	e models, equences. smallest vsterious ries, hid- tling city hatter of for those oft sound -ld away
D.8 ST. State mo used to st guage mo <bos>M raffes ard typewrite for hund State mo we used le A collecti <bos>T plankton expanse ing withi where sk pedestria seeking a of jazz n fron the</bos></bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state models. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w on of randomly generated sentences is shown be the deep blue ocean, teeming with an extraordin- to the largest whales, stretches out infinitely to that has captivated the imaginations of explore in its depths secrets yet to be discovered and story performed as a small café, nestled between two imposing moment of peace, with the comforting aroma of usic playing in the background, creating a co urban chaos outside<sep> The ancient oak</sep></sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc pelow: ary array of owards the h rs, scientists ries yet to be the constant g buildings, of freshly bre zy ambiance tree, with its	mly generated senter open-weights pre-tra- efreshing smoothie <s gs<sep> She collect. In storm that has bee n-trained 1B language ernible with shorter se f marine life, from the porizon, a vast and my s, and poets for centure told<sep> In a bus hum of cars and the co offers a quiet refuge j ewed coffee and the sa e that feels like a wor s gnarled branches su</sep></sep></s 	e models, equences. smallest osterious ries, hid- tling city hatter of for those oft sound cld away tretching
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund State mo we used le A collecti <bos>T plankton expanse ing withi where sk pedestrice seeking co of jazz m from the wide and</bos></bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The following the effects of separator tokens on state midels. Mangoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red with the tallest mammals on Earth due to their long rs from the 1940s<sep> Jupiter's Great Red with the tallest mammals on Earth due to their long rs from the 1940s<sep> Jupiter's Great Red with the largest sep&gt; Hulation on custom-trained 1B models. For onger sentences, as state modulation patterns with the deep blue ocean, teeming with an extraordim to the largest whales, stretches out infinitely to that has captivated the imaginations of explore in its depths secrets yet to be discovered and stop performed by the comforting aroma of usic playing in the background, creating a co- urban chaos outside<sep> The ancient oak is thick, sturdy trunk standing firm against the milies arow seasons abaves and counterparts to the milies arow seasons abaves and counterparts abaves and counterparts to the table the table to table to the table to the tabaves to the table t</sep></sep></sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below: ary array of owards the h rs, scientists ries yet to be the constant g buildings, of freshly bre zy ambiance tree, with it: the passage of owards the h	mly generated senter open-weights pre-tra- efreshing smoothie <s gs<sep> She collect. int storm that has bee in-trained 1B language ernible with shorter se f marine life, from the porizon, a vast and my s, and poets for centure told<sep> In a busi- hum of cars and the co offers a quiet refuge j ewed coffee and the so is gnarled branches su f time, has witnessed banacth its correspondences</sep></sep></s 	e models, equences. smallest esterious ries, hid- tling city hatter of for those off sound cld away retching genera- ognow
D.8 ST. State mo used to st guage mo <bos>M raffes and typewrite for hund State mo we used le A collecti <bos>T plankton expanse ing withi where sk pedestria seeking a of jazz n from the wide ana tions of f becomin</bos></bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w on of randomly generated sentences is shown be the deep blue ocean, teeming with an extraording to the largest whales, stretches out infinitely to that has captivated the imaginations of explore in its depths secrets yet to be discovered and sto percent of peace, with the comforting aroma of usic playing in the background, creating a con- usic playing in the background con- ting and con- ting and con- ting and con- ting and con- ting and con- ting and con- ting an</sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below: ary array of owards the h rs, scientists ries yet to be the constant g buildings, of freshly bre zy ambiance tree, with its be passage of ories unfold	mly generated senter open-weights pre-tra- efreshing smoothie <s gs<sep> She collect. ant storm that has bee entited 1B language entible with shorter set f marine life, from the worizon, a vast and my c, and poets for centure told<sep> In a busch hum of cars and the co offers a quiet refuge j gewed coffee and the set that feels like a wor s gnarled branches su f time, has witnessed beneath its expansive how set splace and the</sep></sep></s 	e models, equences. smallest esterious ries, hid- tling city hatter of for those oft sound 'ld away retching genera- canopy, sense of
D.8 ST. State mo used to st guage mo <bos>M raffes ara typewrite for hund State mo we used le A collecti <bos>T plankton expanse ing withi where sk pedestrid seeking a of jazz n from the wide and tions of f becomin, continuiti</bos></bos>	Atte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow duty the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w on of randomly generated sentences is shown to the deep blue ocean, teeming with an extraordir to the largest whales, stretches out infinitely to that has captivated the imaginations of explore is depths secrets yet to be discovered and sto vescrapers tower over narrow streets filled with the series of peace, with the comforting aroma of usic playing in the background, creating a co urban chaos outside<sep> The ancient oak its thick, sturdy trunk standing firm against the amilies grow, seasons change, and countless states ye in a rapidly changing world<sep></sep></sep></sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below: eary array of owards the h rs, scientists ries yet to be the constant g buildings, of freshly bre zy ambiance tree, with its to passage o pries unfold to those wh	mly generated senter open-weights pre-tra- efreshing smoothie <s gs<sep> She collect. ant storm that has bee m-trained 1B language ernible with shorter se f marine life, from the corizon, a vast and my s, and poets for centure told<sep> In a bus hum of cars and the co offers a quiet refuge j ewed coffee and the so that feels like a wor s gnarled branches su of time, has witnessed beneath its expansive ho seek solace and a</sep></sep></s 	e models, equences. smallest esterious cies, hid- bling city hatter of for those oft sound cld away tretching genera- canopy, sense of
D.8 ST. State mo used to st guage mo <bos>M raffes ara typewrite for hund State mo we used la A collecti <bos>T plankton expanse ing withi where sk pedestria seeking a of jazz m from the wide and tions of f becomin, continuiti</bos></bos>	And the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> Hulation on custom-trained 1B models. For onger sentences, as state modulation patterns w for on of randomly generated sentences is shown be the deep blue ocean, teeming with an extraordin to the largest whales, stretches out infinitely to that has captivated the imaginations of explore is depths secrets yet to be discovered and stor percent of peace, with the comforting aroma of usic playing in the background, creating a co urban chaos outside<sep> The ancient oak its thick, sturdy trunk standing firm against than anilies grow, seasons change, and countless state or a rapidly changing world<sep></sep></sep></sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below: eary array of owards the h rs, scientists ries yet to be the constant g buildings, of freshly bre zy ambiance tree, with its to passage o ories unfold to those wh	mly generated senter open-weights pre-tra- efreshing smoothie <s gs<sep> She collect. and storm that has bee entities that has bee in-trained 1B language entible with shorter set fmarine life, from the corizon, a vast and my s, and poets for centure told<sep> In a busis hum of cars and the co offers a quiet refuge y ewed coffee and the set that feels like a work s gnarled branches su of time, has witnessed beneath its expansive ho seek solace and a</sep></sep></s 	e models, equences. smallest prices, hid- tling city hatter of for those off sound cld away pretching genera- canopy, sense of
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund State mo we used le A collecti <bos>T plankton expanse ing withi where sk pedestrice seeking co of jazz m from the wide and tions of f becomin, continuit We note t</bos></bos>	And the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> And the tallest mammals on Earth due to their long rs from the 1940s<sep> Jupiter's Great Red reds of years<sep> And the tallest mammals on Earth due to their long rs from the 1940s<sep> Jupiter's Great Red reds of years<sep> And the tallest mammals on Earth due to their long rs from the 1940s<sep> Jupiter's Great Red reds of years<sep> And the tallest mammals on Earth due to their long rs from the 1940s<sep> Jupiter's Great Red reds of years<sep> And the tallest whales, stretches out infinitely to the deep blue ocean, teeming with an extraordin to the largest whales, stretches out infinitely to that has captivated the imaginations of explore is depths secrets yet to be discovered and sto vscrapers tower over narrow streets filled with the ns, a small café, nestled between two imposing moment of peace, with the comforting aroma of usic playing in the background, creating a co urban chaos outside<sep> The ancient oak its thick, sturdy trunk standing firm against the amilies grow, seasons change, and countless state g a silent guardian of the park, offering shade y in a rapidly changing world<sep> hat the specific sentences and their order are n</sep></sep></sep></sep></sep></sep></sep></sep></sep></sep></sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below: eary array of wards the h rs, scientists ries yet to be the constant g buildings, of freshly bre zy ambiance tree, with it: the passage of ories unfold to those whe	mly generated senter open-weights pre-tra- efreshing smoothie <s gs<sep> She collect. and storm that has bee entited 1B language entible with shorter se f marine life, from the corizon, a vast and my s, and poets for centur e told<sep> In a busi- hum of cars and the co offers a quiet refuge j ewed coffee and the se that feels like a wor s gnarled branches su f time, has witnessed beneath its expansive ho seek solace and a this analysis. Simila</sep></sep></s 	e models, equences. smallest esterious ries, hid- tling city hatter of for those oft sound cld away tretching genera- canopy, sense of r patterns
D.8 ST. State mo used to st guage mo <bos>M raffes are typewrite for hund State mo we used le A collecti <bos>T plankton expanse ing withi where sk pedestrice seeking co of jazz n from the wide and tions of f becomin, continuit We note th</bos></bos>	Arte MODULATION OF LARGE LANGUAGE M dulation of open-weight models. The follow addy the effects of separator tokens on state m dels. Angoes are rich in vitamin C and can be blend the tallest mammals on Earth due to their long rs from the 1940s <sep> Jupiter's Great Red reds of years<sep> dulation on custom-trained 1B models. For onger sentences, as state modulation patterns w on of randomly generated sentences is shown be the deep blue ocean, teeming with an extraordin- to the largest whales, stretches out infinitely to that has captivated the imaginations of explore in its depths secrets yet to be discovered and story performed of peace, with the comforting aroma of usic playing in the background, creating a co urban chaos outside<sep> The ancient oak is thick, sturdy trunk standing firm against than amilies grow, seasons change, and countless state g a silent guardian of the park, offering shaded y in a rapidly changing world<sep> hat the specific sentences and their order are n rged with various sentence arrangements, provi-</sep></sep></sep></sep>	ODELS wing rando odulation in ded into a re necks and le Spot is a gia r our custon ere less disc below: wards the h rs, scientists ries yet to be the constant g buildings, of freshly bre zy ambiance tree, with its passage of pries unfold to those wh ot crucial to ded the sent	mly generated senter open-weights pre-tra- efreshing smoothie <s gs<sep> She collect. and storm that has bee on-trained 1B language ernible with shorter se f marine life, from the borizon, a vast and my s, and poets for centure told<sep> In a bust hum of cars and the c offers a quiet refuge j ewed coffee and the se that feels like a wor s gnarled branches su f time, has witnessed beneath its expansive to seek solace and a this analysis. Simila ences are sufficiently</sep></sep></s 	e models, equences. swallest esterious ries, hid- tling city hatter of for those oft sound end away retching genera- canopy, sense of r patterns long.

Training settings are outlined in Table 9.

1890				
1891		Configuration	Value	
1892		Batch Size	16	
1893		Max Sequence Length	32k	
1894		Training Steps	160k	
1895		Optimizer	AdamW	
1896		Learning Rate	0.001	
1907		Weight Decay	0.1	
1007		Num. Layers	24	
1000		Dimension	2048	
1099				
1900		Table 9: 1B LLM se	ttings.	
1901				
1902	The perplexity scores shown in	Figure 7b were computed	l on 16k rar	domly sampled sequences over
1903	the FineWeb (Penedo et al., 202	24) dataset. The raw perp	lexity samp	bles were smoothed via a kernel
1904	density estimation method.			
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				

#### 1944 Ε **EXTENDED EXPERIMENTAL RESULTS** 1945

#### 1946 E.1 EMPIRICAL VALIDATION 1947

1948 In this section, we provide additional results and commentary from the sweep detailed in Section D.31949 that were not presented in the main portion of the paper. One thing to note is that most of the ESS 1950 results presented in Section 4 were computed using the entropy-ESS. However, we also computed 1951 ESS using the tolerance-based approach to affirm that both forms of ESS showcase similar trends. 1952 In particular, we examined tolerances of 1e-1, 1e-3 and 1e-5. Since we observe similar trends across 1953 tolerances, we provide plots for a tolerance of 1e-3 below and omit the others for the sake of brevity.

1954 1955

1956

1965

#### E.1.1 STATE COLLAPSE CONTINUED

1957 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 1958 GLA/WLA, Figure 5 shows that accuracy differences between LA and GLA/WLA are only evident 1959 in the high TSS/high kv bucket of the task-model space. This is because state saturation is acting as a 1960 confounder, worsening performance in LA (see Figure 5 when TSS is 8). Therefore, although state 1961 collapse in GLA/WLA does not result in worse performance than LA in this specific task-model 1962 setting, it remains an issue even for models with smaller states when trained on sufficiently difficult 1963 tasks. This is the motivation behind the task-model setting explored in Section A. 1964

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



1988 Figure 9: (a) TSS/kv vs accuracy across featurizers. This demonstrates that TSS/kv (i.e. memory 1989 capacity) is a worse proxy for model performance than ESS/kv as discussed in Section 4. (b) (total 1990 TSS)/kv vs accuracy across featurizers. This demonstrates that (total TSS)/kv is a worse proxy for 1991 model performance than (total ESS)/kv. (c) ESS/kv vs accuracy across featurizers. (d) (total ESS)/kv vs accuracy across featurizers. (e) ESS/TSS (i.e. state utilization) vs accuracy across featurizers. We 1992 note that models that saturate their state tend to perform worse on the task, which is evidence of the 1993 state saturation phenomenon discussed in Section 4.2. The models that do not saturate their state 1994 but still perform poorly are the models that undergo state collapse. (f) (total ESS)/(total TSS) vs 1995 accuracy across featurizers. (g) ESS vs accuracy across featurizers. Note that without normalizing 1996 by kv (i.e. the task memory), the correlation with accuracy breaks down substantially. (h) TSS vs 1997 accuracy across featurizers.







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

E.1.3 TOLERANCE-ESS MQAR RESULTS

2157 2158

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.









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

2320

2316

### 2319 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.

2340

2341

2342

2343

2365



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 featurizers. (b) (to-tal ESS)/ntc vs accuracy across featurizers. (c) TSS/ntc vs accuracy across featurizers. (d) (total TSS)/ntc vs accuracy across featurizers. (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 featurizers. (b) (total ESS)/vs vs accuracy across featurizers. (c) TSS/vs vs accuracy across featurizers. (d) (total TSS)/vs vs accuracy across featurizers. (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.

2381 E.1.5

## 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 (dmodel=256, heads=8) trained on MQAR (seqlen=2048, 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.





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.

2477 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 featurizers. 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.

# 2538 E.4 POST-TRAINING ANALYSIS





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.





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 8layer 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.

Recall in Section 5.2, we demonstrated an application of post-training ESS analysis through the lens of model distillation. Here, we provide another example of post-training analysis that leverages ESS, this time in order to gain intuition into learned network dynamics – hybridization. Hybridization is
the process of arranging different operators in a multi-layer sequence model (Lieber et al., 2024;
Glorioso et al., 2024; De et al., 2024). Specifically, we measure the per-layer ESS across various hybrid networks and find that the precise ordering of layers significantly influences ESS dynamics, 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 networks and compute the per-layer ESS on each model. We use the tolerance-based ESS since we 2656 want to analyze failure modes of learning in hybrid networks. In Figure 32a, we first note that in the 2657 pure GLA model, many of the layers fail to learn expressive states (as evidenced by the tolerance-2658 ESS being 0), offering intuition as to why the model performs so poorly. Moving on to the hybrid 2659 networks with a single attention layer, we note that all of them perform quite well, except for the 2660 network that has attention in the first layer. Interestingly, we find that when attention is placed as 2661 the first layer, it suffers from state collapse. At a higher level, this substantiates why many state-2662 of-the-art hybrid networks (such as Jamba) do not place attention as the first layer of the network. 2663 However, such hybrids are typically constructed purely on the basis of performance: here, ESS is able to provide a distinct perspective. Next, examining the hybrids with 2 SA layers, we find that 2665 the only poor performing topology is with attention placed in the second and third layers. Again, we 2666 find that the ESS of the attention layers is lower than what we observe in the hybrids that solve the 2667 task, indicating its usefulness as a proxy for performance beyond the 2-layer non-hybrid networks we explored in Section 4. 2668

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 1:7 ratio between attention and non-attention layers. Regarding the GLA layers, we find that despite 2675 both the maximum and average SS varying across models, they do not correspond to changes in 2676 accuracy. One possible explanation for this is that since the attention layers are responsible for 2677 driving the total ESS of the network up due to their unbounded state size, the role of non-attention 2678 layers in hybrid networks may not be captured entirely by the magnitude of their ESS. Nonetheless, 2679 this is something we hope to explore in future work. 2680

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

2697

2698







<sup>2807</sup> 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 2809 state-space model (Falcon Mamba 7B, Figure 7a) exhibits a substantial reduction in state modulation 2810 under the same token substitution. 2811

We attribute this difference to a phenomenon we term "preemptive state modulation" in limited state-2812 size models, which stems from fundamental architectural differences. State-space models (SSMs) 2813 with limited cache must efficiently manage their finite memory capacity and learn to preemptively 2814 modulate state-size to optimize information retention, relying on explicit signals like EOS tokens to 2815 trigger context resets. In contrast, attention models with linearly increasing cache can store all past 2816 information without the need for selective forgetting, do not require preemptive state modulation, 2817 and show less sensitivity to explicit demarcation tokens. This distinction highlights the different 2818 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. 2819

However, a subset of attention models demonstrated varying state modulation patterns in response 2821 to different separator tokens, with this effect being more pronounced in smaller model sizes (see 2822 Figure 34, 35, and 36). This phenomenon, while not consistent across all attention architectures, 2823 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 dif-2846 ferent tolerances. Softmax attention exhibits the most pronounced state modulation, beginning at a 2847 tolerance level of 1e-2, followed by gated linear attention with significant modulation starting at a 2848 tolerance of 1e-1. Weighted linear attention shows minimal modulation, only detectable at a toler-2849 ance of 1.0, while linear attention displays no discernible separator token-induced state modulation. 2850 Here ESS is summed across channels and layers. 2851

2852 Notably, GLA exhibits a substantial variation in state modulation depending on the separator token, 2853 consistent with our earlier observations in Falcon Mamba regarding preemptive state modulation. In 2854 contrast, SA shows a smaller, yet non-trivial, effect. WLA and LA show no discernible differences 2855 across separator tokens, which may be attributed to their overall limited ability to modulate state 2856 size.

2857

2825

2827

2829

2837

2841

2845

2858

# 2862 E.6 MISCELLANEOUS

# 2864 E.6.1 EFFECTIVE STATE-SIZE ON C++ CODE 2865

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

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

2870 2871 2872

2873

2874

2875

2876

2877 2878 2879

2880 2881

2882

2883 2884

2885

2886

#### Quake fast inverse square-root algorithm:



Figure 41: Effective state-size over a quake fast inverse square root algorithm's code. Here ESS is summed across channels and layers.

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

```
2916
     34
             float std_result = 1.0f / std::sqrt(number);
2917
     35
             std::cout << "Standard Inverse Sqrt: " << std_result << std::endl;</pre>
2918 36
2919 37
             return 0:
     38 }
2920
2921
2922
        Fibonacci sequence generating algorithm:
2923
2924
                                       ESS of Fibonacci Algorithm on Llama-3-8B-Instruct
                         \times 10^{6}
2925

    tol = 1.0e-01

                     1.75
2926
                          _
                           tol = 1.0e-02
2927
                   .50 Siz
                  State
State
2929
                   Effective
                    1.00
2930
                    0.75
2931
                   Total
                    0.50
2933
                     0.25
2934
2935
2936
        Figure 42: Effective state-size over a Fibonacci sequence generator algorithm's code. Here ESS is
2937
        summed across channels and layers.
2938
2939
2940
      1 #include <iostream>
2941
      2 int fibonacci(int n)
                                   {
2942
             if (n <= 1) {
      3
2943
                  return n;
      4
2944
      5
             }
             return fibonacci(n - 1) + fibonacci(n - 2); // Recursive case
2945
      6
      7 }
2946
      8 int main() {
2947
      9
             int n;
2948
             std::cout << "Enter a positive integer: ";</pre>
     10
2949 11
             std::cin >> n;
             std::cout << "Fibonacci number at position " << n << " is: " <<</pre>
2950 <sup>12</sup>
             fibonacci(n) << std::endl;</pre>
2951
     13
             return 0;
2952
     14 }
2953
2954
2955
2956
2957
                HOW THE NUMBER OF PROMPTING SHOTS AFFECTS THE EFFECTIVE STATE-SIZE OF
        E.6.2
2958
                LANGUAGE MODELS
2959
2960
        Here, we explore how varying the number of shots when prompting large language models affects
2961
        their effective state-size patterns. We use Phi-2 as the candidate attention model and Mamba-2.8B
2962
        as the state-space model. The task we tested this on is MMLU (elementary mathematics).
2963
        At the start of the Q&A section for the attention model, there is a noticeable difference in state size
2964
        between 0-shot and 1-shot prompts. Beyond 1-shot, the difference in ESS appears minimal. For the
2965
        state-space model, varying the number of shots has minimal impact on the effective state-size.
2966
2967
        Although these sparse experimental results require further investigation, we note the stark differ-
        ence in the effective state-size patterns between these two architectures, which provides additional
2968
        insights into understanding the fundamental differences in the way prompts are processed across
2969
        models.
```



Figure 43: The variation in effective state-size with a varying number of shots (2.7B Attention). Here ESS is summed across channels and layers.



Figure 44: The variation in effective state-size with a varying number of shots (2.8B State-Space Model). Here ESS is summed across channels and layers.