TAIPAN: EFFICIENT AND EXPRESSIVE STATE SPACE LANGUAGE MODELS WITH SELECTIVE ATTENTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Efficient long-context language modeling remains a significant challenge in Natural Language Processing (NLP). While Transformers dominate language tasks, they struggle with long sequences due to quadratic computational complexity in training and linearly scaling memory costs during inference. Recent State Space Models (SSMs) such as Mamba offer alternatives with constant memory usage, but they underperform in tasks requiring extensive in-context retrieval. We introduce Taipan, a novel hybrid architecture that combines Mamba-2 with Selective Attention Layers (SALs). These SALs identify tokens requiring long-range interactions, remove less important features, and then augment their representations using the attention module. This approach balances Mamba's efficiency with Transformer-like performance in memory-intensive tasks. By constraining the attention budget, Taipan extends accurate predictions to context lengths of up to 1 million tokens while preserving computational efficiency. Our experiments demonstrate Taipan's superior performance across various scales and tasks, offering a promising solution for efficient long-context language modeling.

025 026

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

027 028

029 Transformer-based architectures Vaswani (2017); Brown (2020) have revolutionized Natural Language Processing (NLP), delivering exceptional performance across diverse language modeling tasks Touvron et al. (2023). This success stems from their ability to capture complex word de-031 pendencies using the self-attention mechanism. In addition, Transformers are highly scalable and well-suited for parallel training on large datasets. However, despite their success, they still face no-033 table challenges when handling long-context sequences. Specifically, the self-attention mechanism 034 suffers from quadratic computational complexity, and the memory requirement grows linearly with 035 context length during inference, as the model must store key-value vectors for the entire context. These factors impose practical constraints on sequence length due to the high computational and 037 memory costs. 038

To this end, recent advancements in recurrent-based architectures, particularly State Space Models (SSMs) Gu et al. (2021b;a), have emerged as promising alternatives for efficient language modeling 040 Gu & Dao (2023); Dao & Gu (2024). SSMs offer constant memory usage during inference, and 041 architectures like Mamba-2 Dao & Gu (2024), a variant of SSMs, have demonstrated performance 042 comparable to Transformers in certain language tasks Waleffe et al. (2024). Some studies even 043 suggest that SSMs can outperform Transformers in areas like state tracking Merrill et al. (2024) due 044 to their Markovian nature. However, despite these advancements, SSM-based models still fall short 045 in scenarios requiring in-context retrieval or handling complex long-range dependencies Arora et al. (2024); Waleffe et al. (2024). 046

To address these challenges, we introduce Taipan, a hybrid architecture that combines the efficiency of Mamba with enhanced long-range dependency handling through Selective Attention Layers (SALs). While Mamba is highly efficient, it relies on the Markov assumption—where predictions are based solely on the last hidden state—which can lead to information loss for tokens that need interactions with distant tokens. To mitigate this, Taipan incorporates SALs that strategically select key tokens in the input sequence requiring long-range dependencies. These selected tokens first undergo feature refinement to remove unimportant information, and then are passed through an attention module to capture long-range dependencies. Less critical tokens bypass the attention step, as

16K 32K 64K Context Length

a) Perplexity Across Context Lengths



054

14

12

Perplexity

2K Δĸ

1K



068

069

071 072

Figure 1: Model Performance Comparison. a) Perplexity across different context lengths. Lower perplexity indicates better performance. b) Latency comparison of models at various generation lengths. Taipan exhibits significantly lower latency and superior scaling compared to other strong baselines for longer sequences.

1 M

500

400

100

1K2K 4k

(s) Laten

Models

lamba

Mamba

Taipan

Transformer

16k

Generation Length

b) Latency Across Generation Lengths

32K

Models

lamba

Mamba

Taipan

Transforme

OOM

128K

we hypothesize that their Markovian representations from Mamba contain sufficient information for 073 accurate prediction, obviating the need for additional attention-based augmentation. This selective 074 approach enables Taipan to balance Mamba's computational efficiency with enhanced long-range 075 modeling capabilities. 076

077 SALs play a crucial role in Taipan's design, both in enhancing performance and ensuring computational efficiency. By focusing the attention mechanism on a subset of important tokens, SALs reduce the computational costs that come from attention modules. This targeted approach enables Taipan 079 to excel in memory-intensive tasks while maintaining efficiency during both training and inference. Importantly, Taipan retains the linear memory usage characteristic of SSMs, offering a significant 081 advantage over traditional Transformer models in handling extremely long sequences.

083 We scale Taipan to 190M, 450M, and 1.3B parameters, pre-training on 100B tokens. Experimental 084 results demonstrate Taipan's superior performance across a wide range of tasks. In zero-shot language modeling evaluations, Taipan consistently outperforms both Transformer and Mamba base-085 lines, showcasing its strong general language understanding capabilities. More notably, in memory-086 intensive tasks such as long-context retrieval and structured information extraction, Taipan exhibits 087 significant improvements over Mamba-2, addressing a key limitation of existing recurrent-based 088 models. Furthermore, Taipan demonstrates remarkable extrapolation capabilities, maintaining high 089 performance on sequences up to 1 million tokens in context-length - while preserving efficient gen-090 eration capabilities. This combination of broad task proficiency, superior performance in memory-091 intensive scenarios, and exceptional long-context modeling positions Taipan as a versatile and pow-092 erful architecture for advanced language processing tasks.

093 094

2 BACKGROUND

095 096

This section briefly overviews the foundational architectures relevant to our work. We first review 098 Causal Self-Attention Vaswani (2017), the core mechanism of Transformer models. We then discuss Linear Attention Katharopoulos et al. (2020), an efficient variant that achieves linear complexity. Finally, we examine Mamba-2, a recent architecture that generalizes Linear Attention using structured 100 state-space models (SSMs) Dao & Gu (2024). We emphasize how each model balances computa-101 tional efficiency and recall accuracy, particularly in memory-intensive tasks . 102

103

104 2.1 CAUSAL SELF-ATTENTION

105

Causal Self-Attention is the key component in Transformer architectures that allows each token 106 in a sequence to attend to all other previous tokens (Vaswani, 2017). Given an input sequence 107 $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_L] \in \mathbb{R}^{L \times d}$, where L is the sequence length and d is the embedding dimension, self-attention firsts computes the query, key, and value vectors for each token via linear projections:

$$\mathbf{q}_i = \mathbf{W}_Q \mathbf{x}_i, \quad \mathbf{k}_i = \mathbf{W}_K \mathbf{x}_i, \quad \mathbf{v}_i = \mathbf{W}_V \mathbf{x}_i$$

where $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d}$ are learnable weight matrices.

Then, the attention output o_i for each token x_i will be calculated as a weighted sum of the value vectors over the distribution of similarity matrix between its query vector and previous key vectors:

$$\mathbf{o}_i = \sum_{t=1}^i \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_t / \sqrt{d})}{\sum_{j=1}^i \exp(\mathbf{q}_i^\top \mathbf{k}_j / \sqrt{d})} \mathbf{v}_t$$

117 118 119

> 123 124

> 125

130 131

132 133

138

141

142 143

115 116

110

The non-linear softmax distribution allows the models to capture intricate relationships between tokens, and concentrate on salient features Qin et al. (2022); Zhao et al. (2019). As such, selfattention can encode complex language patterns and long-range dependencies that are crucial for complex language understanding and generation tasks.

2.2 LINEAR ATTENTION

To address the quadratic complexity, recent work has shown that it is possible to achieve linear complexity with the attention mechanism by replacing the softmax attention with dot-product attention (Shen et al., 2021; Katharopoulos et al., 2020). Given a feature transformation $\phi(\mathbf{x})$, causal self-attention can be rewritten as:

$$\mathbf{o}_i = \sum_{t=1}^i \frac{\phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_t)}{\sum_{j=1}^i \phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_j)} \mathbf{v}_t$$

Then, using the associate property of matrix multiplication, this can be reformulated as:

$$\mathbf{o}_i = \frac{\phi(\mathbf{q}_i)^\top \sum_{t=1}^i \phi(\mathbf{k}_t) \mathbf{v}_t^\top}{\phi(\mathbf{q}_i)^\top \sum_{t=1}^i \phi(\mathbf{k}_t)}$$

Let $\mathbf{S}_i = \sum_{t=1}^i \phi(\mathbf{k}_t) \mathbf{v}_t^{\top}$ and $\mathbf{z}_i = \sum_{t=1}^i \phi(\mathbf{k}_t)$. We can then rewrite the equation in a recurrent form:

$$\mathbf{S}_i = \mathbf{S}_{i-1} + \phi(\mathbf{k}_i) \mathbf{v}_i^\top$$

$$\mathbf{o}_i = rac{\mathbf{S}_i \phi(\mathbf{q}_i)}{\mathbf{z}_i^ op \phi(\mathbf{q}_i)} pprox \mathbf{S}_i \phi(\mathbf{q}_i)$$

144 145

156

158

This formulation allows for efficient training and inference. Let $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{L \times d}$ be the query, key, and value matrices of the sequence input \mathbf{X} . During training, we can use the matrix multiplication form: $\mathbf{O} = (\mathbf{Q}\mathbf{K}^{\top} \odot \mathbf{M}_L)\mathbf{V}$, where \mathbf{M}_L is a causal mask. At inference time, we can use the recurrent form for efficient sequential processing.

However, despite its computational efficiency, linear attention has notable limitations compared to softmax attention. The dot-product approximation in linear attention lacks the nonlinear normalization of softmax, often resulting in a more uniform distribution of attention weights Han et al. (2023).
This uniformity can impair the model's ability to focus sharply on specific and relevant tokens. Consequently, linear attention models may underperform in tasks requiring precise in-context retrieval or focused attention on particular input segments Han et al. (2023).

157 2.3 Мамва-2

Mamba Gu & Dao (2023) is a variant of structured state space models (SSMs) that uses the selective data-dependent mechanism. Mamba-2 Dao & Gu (2024) builds on this foundation, revealing deep connections between SSMs and linear attention Katharopoulos et al. (2020) through the framework of structured state-space duality (SSD).



Figure 2: An overview of the Taipan architecture.

The core of Mamba-2 can be defined by using the recurrent form:

$$\mathbf{h}_t = \mathbf{A}_t \mathbf{h}_{t-1} + \mathbf{B}_t \mathbf{x}$$

 $\mathbf{o}_t = \mathbf{C}_t \mathbf{h}_t$

185 where A_t is further simplified to a scalar multiplied by the identity matrix. This formulation allows 186 Mamba-2 to be interpreted as a generalization of linear attention.

The key insight of Mamba-2 is that this recurrence can be equivalently expressed as a matrix multiplication:

$$\mathbf{O}_t = (\mathbf{L}_t \odot \mathbf{C}_t \mathbf{B}_t^\top) \mathbf{X}_t$$

where L is a 1-semiseparable matrix. This matrix form reveals the duality between the recurrent (linear-time) and attention-like (quadratic-time) computations. Also, the 1-semiseparable matrix L encodes the temporal dependencies, while CB^{T} represents content-based interactions similar to attention. This formulation generalizes linear attention, which can be seen as a special case where L is the all-ones lower triangular matrix.

While Mamba-2 is efficient, it shares the same limitations as Linear Attention in terms of precise memory recall Arora et al. (2024); Wen et al. (2024), leading to reduced performance in tasks that demand accurate retrieval of specific sections in the input sequence.

200 3 TAIPAN MODEL

effective sequence representation.

179

181

182 183

187

188

189

199

207

To address the limited modeling capabilities of Mamba-2 and Linear Attention while preserving
their computational efficiency, we introduce Taipan, a new architecture for sequence encoding in language modeling. In Taipan, we strategically incorporate Selective Attention Layers (SALs) within
the Mamba framework, as shown in Figure 2. SALs are inserted after every *K* Mamba-2 blocks, creating a hybrid structure that combines Mamba-2's efficiency with Transformer-style attention for

The core of SALs is a gating network that identifies important tokens for enhanced representation modeling. These tokens undergo two phases: (1) feature refinement to filter out irrelevant information and (2) representation augmentation via softmax attention. This allows Taipan to capture complex, non-Markovian dependencies when necessary.

Taipan processes input through Mamba-2 blocks, with SALs periodically refining key token representations. These enhanced representations are then passed into the subsequent Mamba-2 layers, influencing further processing. This hybrid structure balances Mamba-2's efficiency with the expressive power of SALs, enabling Taipan to excel in tasks requiring both speed and accurate information retrieval. The following sections detail each component's structure and function.



Figure 3: Attention mechanisms in Taipan's Selective Attention Layers. White areas indicate no attention. (a) Full Causal Attention (b) Sliding Window Attention (w = 4) (c) Selective Attention (C = 0.3, w = 5)

3.1 SELECTIVE ATTENTION LAYERS

Selective Attention Layers (SALs) are the key innovation in Taipan, designed to enhance the model's ability to focus on critical tokens while maintaining overall efficiency. These layers employ a lightweight gating network G_{θ} to dynamically determine which tokens should undergo softmax attention processing.

For each token hidden representation h_i in the input sequence, the gating network G computes a score vector:

245

250 251

264

265

226

227

228 229 230

231

 $\mathbf{s}_i = G_\theta(\mathbf{h}_i) \tag{1}$

where $G_{\theta} : \mathbb{R}^d \to \mathbb{R}^2$ is parameterized by θ . This score vector $\mathbf{s}_i = [s_{i,0}, s_{i,1}]$ serves two purposes: 1) it is used to generate a binary mask m_i for token selection, and 2) it guides feature refinement.

To maintain differentiability while allowing for discrete token selection, we employ the Straight-Through Gumbel-Softmax trick Jang et al. (2017). A binary mask m_i is generated from s_i to select tokens during the forward pass of the network:

$$m_i = \operatorname{argmax}(\operatorname{GumbelSoftmax}(\mathbf{s}_i, \tau)) \tag{2}$$

where τ is the temperature parameter. \mathbf{h}_i will only be selected for attention processing if $m_i = 1$.

For the backward pass, we instead use continuous Gumbel-Softmax approximation of m_i to achieve computation differentiability for the network:

$$\tilde{m}_{i} = \frac{\mathbb{I}[m_{i} = 0] \exp((\mathbf{s}_{i,0} + g_{0})/\tau) + \mathbb{I}[m_{i} = 1] \exp((\mathbf{s}_{i,1} + g_{1})/\tau)}{\exp((\mathbf{s}_{i,0} + g_{0})/\tau) + \exp((\mathbf{s}_{i,1} + g_{1})/\tau)}$$
(3)

where $\mathbb{I}[]$ is the indicator function, and g_0 and g_1 are i.i.d samples from the Gumbel(0, 1) distribution. In this way, we are able to train our entire model, including the gating network, in an end-to-end fashion for language modeling.

For the selected tokens (those with a mask value m_i of 1), we compute their attention-based representations:

$$\mathbf{p}_i = \text{Attention}(\mathbf{q}_i, \mathbf{K}, \mathbf{V}) \tag{4}$$

where \mathbf{q}_i is the query vector for the *i*-th selected token (denoted \mathbf{h}_i^s), and **K** and **V** are the key and value matrices for previous tokens.

In our model, the score vector \mathbf{s}_i is also used to refine the representations of selected tokens. We employ the softmax of \mathbf{s}_i to compute the mixing weights: $[1 - \alpha_i, \alpha_i] = \operatorname{softmax}(\mathbf{s}_i)$. The final output for a selected token \mathbf{h}_i^s is a weighted combination:

$$\mathbf{h}_{i}^{s} = (1 - \alpha_{i})\mathbf{h}_{i}^{s} + \alpha_{i}\mathbf{o}_{i}$$
⁽⁵⁾

As such, Taipan can adaptively preserve key information in \mathbf{h}_i^s while enriching the representation with the attention output \mathbf{o}_i . In other words, α_i acts as the *data-dependent factor*, filtering out unimportant features from the original representation while integrating richer information from the attention outputs. Here, it is important to note that unselected tokens (i.e., $m_i = 0$) skip the attention module and retain their original representations from Mamba-2. Finally, all token representations are passed through a residual SwiGLU Shazeer (2020) layer:

$$\mathbf{h} = \mathbf{h} + \mathrm{SwiGLU}(\mathbf{h}) \tag{6}$$

This final transformation ensures that all token representations undergo consistent non-linear processing before being passed to the next layer in the network, enhancing the model's ability to capture complex dependencies.

278 3.2 SLIDING WINDOW ATTENTION

279 To maintain linear time complexity while leveraging the benefits of attention, Taipan employs Slid-280 ing Window Attention (SWA) Beltagy et al. (2020). SWA's computational complexity scales linearly 281 with sequence length, allowing Taipan to handle theoretically unlimited context lengths during infer-282 ence. Importantly, the combination of Selective Attention and Sliding Window Attention in Taipan 283 leads to a significantly sparser attention weight map compared to full attention or standard windowed 284 attention (Figure 3), thus enhancing the computational efficiency of Selective Attention for process-285 ing long sequences for our model. In addition, the sparser attention map allows us to afford a longer 286 sliding window (i.e., w = 2048 in our work) to effectively capture longer-range dependencies for 287 input sequences. In this way, our designed Taipan architecture offers a mechanism to balance the efficient processing of long sequences with the ability to capture important long-range dependencies, 288 thereby addressing a key limitation of existing efficient attention mechanisms. Finally, removing 289 positional embeddings from the Attention Module improves extrapolation capabilities, suggesting 290 that the model can better generalize temporal relationships. We explore this impact of positional 291 embeddings in more detail in Section 5.2. 292

293 294

295

272

273

277

3.3 TRAINING AND INFERENCE

To better balance efficiency and expressiveness, we introduce an attention budget constraint. Given a predefined budget C, representing the desired fraction of tokens to receive attention, we incorporate a constraint loss into our training objective:

$$\mathcal{L}_{\text{constraint}} = \sum_{n=1}^{N} \left\| C - \frac{\sum_{i=1}^{L} m_i}{L} \right\|_2^2 \tag{7}$$

Here, N is the number of SALs, L is the sequence length, and $\sum_{i=1}^{L} m_i$ represents the number of tokens selected for attention processing. During training, we employ the Straight Through Gumbel Softmax estimator for \tilde{m}_i in the backward pass Jang et al. (2017); Bengio et al. (2013), ensuring differentiability while maintaining discrete token selection in the forward pass, thereby enabling end-to-end training of the entire model. As such, our overall training objective includes a standard cross-entropy loss \mathcal{L}_{CE} for language modeling and the budget constraint term: $\mathcal{L} = \mathcal{L}_{CE} + \lambda \mathcal{L}_{constraint}$, where λ is a hyperparameter.

During inference, Taipan processes input tokens sequentially through Mamba-2 blocks. At each Selective Attention Layer, the gating network G_{θ} computes a score vector $\mathbf{s}_i = G_{\theta}(\mathbf{h}_i)$ for each token representation \mathbf{h}_i . This score computes a binary mask m_i to determine if \mathbf{h}_i should be used for attention processing. Consequently, our selective attention approach maintains Mamba-2's efficiency for most tokens while applying targeted attention to critical elements, enabling effective long-range dependency modeling with minimal computational overhead.

316 317

318

4 Experiments

We conducted extensive experiments to evaluate Taipan's performance across various scales and tasks. Our evaluation strategy focuses on three main areas: (1) zero-shot evaluation on diverse benchmarks to demonstrate Taipan's general language understanding capabilities (Section 4.2), (2) in-context retrieval tasks to assess Taipan's ability to retrieve information from historical contexts (Section 4.3), and (3) extrapolation ability in long-context scenarios to evaluate performance on extremely long sequences (Section 4.4).

324 4.1 EXPERIMENTAL SETUP

328

330

331

332

333

334

335

336

337

338

371

We evaluate Taipan across three model sizes: 190M, 450M, and 1.3B parameters. To ensure a comprehensive and fair comparison, we benchmark Taipan against three strong baselines:

- **Transformer++** Touvron et al. (2023): An enhanced version of the LLaMA architecture Touvron et al. (2023), incorporating Rotary Positional Embeddings Su et al. (2024), SwiGLU Shazeer (2020), and RMSNorm Zhang & Sennrich (2019).
- Mamba-2 Dao & Gu (2024): A state-of-the-art linear RNN model based on State Space Models (SSMs). Each Mamba-2 block consists of a depthwise convolutional layer Poli et al. (2023); Gu & Dao (2023), an SSM layer Dao & Gu (2024), and MLP layers.
- Jamba Lieber et al. (2024): A hybrid model combining full Causal Self-Attention layers (with Rotary Position Embedding Su et al. (2024)) and Mamba-2 layers. Unlike Taipan, Jamba uses full Causal self-attention instead of selective attention, retains positional embeddings, and lacks a feature refinement mechanism.

Implementation Details We train all models from scratch in three configurations: 190M, 450M, and 1.3B parameters. The training process is consistent across configurations with the following hyperparameters: a batch size of 0.5M tokens per step, a cosine learning rate schedule with 2000 warm-up steps, and AdamW Loshchilov (2017) optimization with a peak learning rate of 5e - 4decaying to a final rate of 1e - 5. We apply a weight decay of 0.01 and use gradient clipping with a maximum value of 1.0. All models are trained with a fixed context length of 4096 tokens.

The training data size varies by model scale: the 190M model is trained on 27 billion tokens, while the 450M and 1.3B models are trained on 100 billion tokens. The dataset details can be found in Appendix A.

For Taipan-specific implementation, we use a hybrid ratio of 6:1, inserting a Selective Attention Layer (SAL) after every 6 Mamba-2 Blocks. The Mamba-2 blocks are kept identical to the original work Dao & Gu (2024). We set the attention capacity C = 0.15. The sliding window attention mechanism uses a window size (w) of 2048 tokens.

Params & Data	Model	Wino.	PIQA	Hella.	\mathbf{ARC}_E	\mathbf{ARC}_C	OB.	Truth.	RACE	BoolQ	Avg.
	Transformer++	47.1	60.9	27.9	42.2	20.5	18.9	42.9	25.4	57.2	38.1
190M 27B	Mamba	49.6	60.7	29.3	45.3	21.8	20.6	40.8	27.2	59.3	39.4
	Jamba	49.9	60.3	29.2	46.3	21.4	18.5	39.8	27.4	58.6	39.1
	Taipan	51.0	62.6	29.4	46.7	20.7	21.8	41.1	26.6	58.7	39.9
450M 100B	Transformer++	51.5	67.6	42.3	60.8	27.7	33.4	39.2	30.5	54.7	45.3
	Mamba	52.7	68.9	42.7	61.4	27.1	34.0	38.5	29.3	53.2	45.3
	Jamba	53.1	69.3	44.3	62.6	28.7	34.4	37.5	31.3	55.7	46.3
	Taipan	53.0	69.6	46.6	65.6	32.9	36.6	38.6	30.7	60.4	48.2
1.3B 100B	Transformer++	53.8	71.6	53.8	63.2	36.3	36.4	44.0	31.2	59.4	49.9
	Mamba	55.2	73.0	55.6	70.7	38.0	39.0	39.9	32.0	61.8	51.7
	Jamba	54.7	73.8	55.8	69.7	37.6	41.8	40.4	32.8	59.2	51.8
	Taipan	57.0	74.9	57.9	71.2	39.3	40.4	43.0	34.4	61.5	53.3

Table 1: Z	Zero shot	results of	f Taipan	against	baseline	models.
14010 11 2	1010 01100	1000100 01	- impour	againse	casenne	

4.2 LANGUAGE MODELING PERFORMANCE

We report the zero-shot performance of Taipan and baseline models on a diverse set of commonsense reasoning and question-answering tasks. These include Winograd (Wino.) Sakaguchi et al.
(2021), PIQA Bisk et al. (2020), HellaSwag (Hella.) Zellers et al. (2019), ARC-easy and ARCchallenge (ARCe & ARCc) Clark et al. (2018), OpenbookQA (OB.) Mihaylov et al. (2018), TruthfulQA (Truth.) Lin et al. (2021), RACE Lai et al. (2017), and BoolQ Clark et al. (2019). It is
worth noting that these tasks are brief and do not involve in-context learning, thus inadequately
demonstrating long-context modeling or in-context learning retrieval abilities.

Table 1 presents the zero-shot results for models of three sizes: 190M, 450M, and 1.3B parameters.
 The results are evaluated using the lm-evaluation-harness¹ Gao et al. (2024) framework.

As can be seen, Taipan consistently outperforms the baseline models across most tasks for all model sizes. Notably, the performance gap widens as the model size increases, with the 1.3B Taipan model showing significant improvements over other baselines. This suggests that Taipan's architecture effectively captures and utilizes linguistic patterns, even in tasks that do not fully showcase its longcontext modeling capabilities.

386 387

388 389

390

391

392

394

395

396 397

398

399

400

401

402

403

404

405

406

407

408

409

413 414 415

416

4.3 IN-CONTEXT RECALL-INTENSIVE PERFORMANCE

To evaluate Taipan's proficiency in precise in-context retrieval, we assessed all models on a set of recall-intensive tasks Arora et al. (2024). These tasks are designed to test a model's ability to extract and utilize information from longer contexts, a capability particularly relevant to Taipan's architecture. Our evaluation suite includes two types of tasks: structured information extraction and question answering. For structured information extraction, we used the SWDE and FDA tasks Arora et al. (2024), which involve extracting structured data from HTML and PDF documents, respectively. To assess question-answering capabilities, we employed SQuAD Rajpurkar et al. (2018), which requires models to ground their answers in provided documents.

Params	Model	SWDE	FDA	SQuAD	Avg.
450M	Transformer++	43.0	48.7	18.1	36.6
	Mamba	27.9	9.8	12.5	16.7
	Jamba	35.4	36.6	16.3	29.4
	Taipan	<u>41.4</u>	<u>39.6</u>	<u>17.8</u>	<u>32.9</u>
1.3B	Transformer++	64.2	64.5	41.2	56.6
	Mamba	48.6	32.3	31.2	37.4
	Jamba	56.4	49.7	33.4	46.5
	Taipan	<u>61.5</u>	<u>59.7</u>	<u>36.9</u>	52.7

while consuming fewer computational resources than Jamba, which utilizes full attention mechanisms. This efficiency is attributed to Taipan's architecture, which combines Mamba-like elements with selective attention mechanisms, allowing it to filter out less important features. We also notice that Transformers excel at memory-intensive tasks in this experiment; however, they are constrained by

Table 4 demonstrates Taipan's significant

performance advantages over both Mamba

and Jamba in in-context retrieval tasks.

Notably, Taipan achieves this superiority

Figure 4: Performance on in-context retrieval tasks.

linear memory scaling with sequence length, limiting their effectiveness and applicability for very
long sequences. In contrast, Taipan maintains constant memory usage, offering a more efficient
solution for processing long documents.

4.4 LONG-CONTEXT EXTRAPOLATION

Figure 1 illustrates Taipan's superior performance in handling extended sequences compared to 417 Transformer, Jamba, and Mamba models. In perplexity evaluations across context lengths from 418 1K to 1M tokens (Figure 1a), Taipan yields the lowest perplexity, particularly excelling beyond 419 the training context length. This performance contrasts sharply with other models: Transformers 420 struggle with longer contexts due to quadratic computational complexity and linear memory scaling 421 with sequence length, often leading to out-of-memory errors. Jamba, despite its hybrid nature, faces 422 similar challenges due to its use of full attention mechanisms. Both Transformer and Jamba mod-423 els exhibit limited extrapolation ability beyond their training context lengths. Mamba, while more 424 efficient than Transformers and Jamba, still shows performance degradation for very long sequences. 425

Latency comparisons (Figure 1b) further highlight Taipan's exceptional efficiency. It demonstrates the lowest latency among all models, with linear scaling across sequence lengths. This contrasts with the quadratic scaling of Transformers and higher latency growth of Jamba. Notably, Taipan consistently outperforms Mamba-2, primarily due to its selective attention mechanism.

- 430
- 431

¹https://github.com/EleutherAI/lm-evaluation-harness

432 5 ABLATION STUDY

434

435

436

437 438

439

We conducted a comprehensive ablation study to investigate the effect of the two key components in Taipan's architecture, i.e., the attention budget capacity C and the inclusion of Positional Embeddings in the SALs, on its performance and efficacy.

5.1 EFFECT OF ATTENTION BUDGET CAPACITY

Our first experiment aimed to determine the optimal value of Capacity C that would maintain computational efficiency while maximizing performance on downstream tasks. We trained multiple variants of Taipan, each with 1.3B parameters, using different Capacity C values: 0.10, 0.15, 0.20, and 0.25. Each variant was trained for 24,000 steps, allowing us to observe both the immediate impact of different C values and their effect on model performance over time.

We evaluated the performance of each variant at regular intervals on two representative tasks: SWDE
Arora et al. (2024) (for structured information extraction) and HellaSwag Zellers et al. (2019) (for
commonsense reasoning). These tasks were chosen to assess both the model's ability to handle
long-context retrieval and its general language understanding capabilities.



Figure 5: Effect of Attention Budget Capacity C on Taipan's Performance

As illustrated in Figure 5, Taipan achieves optimal performance with a Capacity C = 0.15. We observed that increasing C beyond 0.15 does not lead to significant improvements in results while increasing computational costs. Conversely, reducing C below 0.15 resulted in a noticeable drop in performance on tasks requiring precise in-context retrieval or complex long-range dependencies. These findings support our hypothesis that computational demands vary across tokens, with many adequately represented by Mamba's Markovian structure without requiring attention mechanisms. By selectively applying attention only to tokens that benefit from it, Taipan optimizes resource allocation, enabling high performance while improving computational efficiency.

473 474

475

464 465

5.2 IMPACT OF POSITIONAL EMBEDDINGS

Our second experiment investigated the impact of Positional Embeddings in Taipan's Attention
mechanism, focusing on the model's ability to handle and generalize to various context lengths.
We trained two variants of the 1.3B parameter Taipan model for 24,000 steps with a fixed context
length of 4096 tokens. One variant incorporates Rotary Positional Embeddings Su et al. (2024) in
the Selective Attention layers, while the other excludes them. Figure 6 illustrates the performance
of both variants in terms of perplexity across different context lengths.

The results reveal that Taipan without Positional Embeddings performs superiorly in generalizing context lengths beyond the training context. Both variants show comparable performance for sequences similar to or shorter than the training context length. However, as the sequence length increases, the performance gap between the two variants widens, with Taipan without Positional Embeddings maintaining lower perplexity scores. This suggests that the absence of Positional Embeddings enables more robust scaling to longer sequences. We attribute this improved generalization to the model's increased reliance on attention representation rather than positional biases.

RELATED WORK 6

486

487

488 489 490

491

493

492 Our approach builds on a foundation of relevant previous research. We will now discuss 494 key studies that inform our methodology. 495

State Space Models: SSMs have emerged 496 as a promising approach in attention-free 497 architectures for language processing tasks. 498 These models offer improved computational 499 and memory efficiency compared to traditional 500 attention-based models. The development of 501 SSMs has progressed through several key iter-502 ations: S4 Gu et al. (2021a) introduced the first 503 structured SSM, focusing on diagonal and di-504 agonal plus low-rank (DPLR) structures. Sub-505 sequent variants like DSS Gupta et al. (2022), S4D Gu et al. (2022), and S5 Smith et al. 506 (2023) improved on this foundation. Frame-507 works like GSS Mehta et al. (2023), H3 Fu 508 et al. (2023), and RetNet Sun et al. (2023) in-509 corporated SSMs into broader neural network 510 architectures, often combining them with gat-



Figure 6: Perplexity comparison of Taipan variants with and without Positional Embeddings across different context lengths. Lower perplexity indicates better performance.

511 ing mechanisms or efficient attention approximations. Recently, Mamba Gu & Dao (2023) intro-512 duced time-varying or selective SSMs, which addresses limitations of static dynamics in previous 513 SSMs by incorporating input-dependent state transitions, leading to improved performance in vari-514 ous tasks.

515 Hybrid Architecture: Several recent studies H3 Fu et al. (2023), Griffin De et al. (2024), Zamba 516 Glorioso et al. (2024), Jamba Lieber et al. (2024) suggest the potential of blending SSM and the 517 attention mechanism. These hybrid designs show promise in outperforming both traditional Trans-518 formers and pure SSM architectures, such as Mamba, particularly in scenarios requiring in-context 519 learning capabilities. 520

Long Context Models: Recent advancements in sequence modeling have pushed the boundaries 521 of context length, each with distinct approaches and challenges. Recurrent Memory Transformer 522 Bulatov et al. (2023) demonstrated 1M token processing, but primarily on synthetic memorization 523 tasks. LongNet Ding et al. (2023) proposed scalability to 1B tokens, yet practical evaluations were 524 limited to sequences under 100K tokens. Hyena/HyenaDNA Poli et al. (2023); Nguyen et al. (2023) 525 claimed 1M token context, but faced efficiency issues at longer lengths. Mamba Gu & Dao (2023) 526 showed consistent improvements up to 1M tokens in DNA modeling and competitive performance 527 across various language tasks.

528 529

7 CONCLUSION

530 531

532 Taipan presents a significant advancement in long-context language modeling by combining the effi-533 ciency of Mamba with strategically placed Selective Attention Layers. Our experiments demonstrate 534 Taipan's superior performance across various scales and tasks, particularly in scenarios requiring extensive in-context retrieval, while maintaining computational efficiency. A key insight is that not all 536 tokens require the same computational resources. Taipan's architecture leverages this observation 537 through its selective attention mechanism, which dynamically allocates computational resources based on token importance. This hybrid approach addresses limitations of both Transformers and 538 SSMs, offering a promising solution for efficient, large-scale language processing. Future work could explore further optimizations and applications of this architecture.

540 REFERENCES 541

547

551

552

553

581

583

- Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalsina, Silas Alberti, Dylan Zinsley, 542 James Zou, Atri Rudra, and Christopher Ré. Simple linear attention language models balance 543 the recall-throughput tradeoff. arXiv preprint arXiv:2402.18668, 2024. 544
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. 546 arXiv preprint arXiv:2004.05150, 2020.
- Loubna Ben Allal, Anton Lozhkov, Guilherme Penedo, Thomas Wolf, and Leandro von 548 Werra. Smollm-corpus, 2024. URL https://huggingface.co/datasets/ 549 HuggingFaceTB/smollm-corpus. 550
 - Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. CoRR, abs/1308.3432, 2013. URL http://arxiv.org/abs/1308.3432.
- 554 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical com-555 monsense in natural language. In Proceedings of the AAAI conference on artificial intelligence, 556 volume 34, pp. 7432-7439, 2020.
- Tom B Brown. Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 2020. 558
- 559 Aydar Bulatov, Yuri Kuratov, and Mikhail S. Burtsev. Scaling transformer to 1m tokens and beyond 560 with rmt. ArXiv, abs/2304.11062, 2023. URL https://api.semanticscholar.org/ 561 CorpusID:258291566. 562
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina 563 Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. arXiv preprint 564 arXiv:1905.10044, 2019. 565
- 566 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and 567 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. arXiv preprint arXiv:1803.05457, 2018. 568
- 569 Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through 570 structured state space duality. arXiv preprint arXiv:2405.21060, 2024. 571
- Soham De, Samuel L. Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Al-572 bert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, Guillaume Des-573 jardins, Arnaud Doucet, David Budden, Yee Whye Teh, Razvan Pascanu, Nando De Freitas, and 574 Caglar Gulcehre. Griffin: Mixing gated linear recurrences with local attention for efficient lan-575 guage models, 2024. URL https://arxiv.org/abs/2402.19427. 576
- 577 Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang, Shaohan Huang, Wenhui Wang, Nanning 578 Zheng, and Furu Wei. Longnet: Scaling transformers to 1, 000, 000, 000 tokens. CoRR, abs/2307.02486, 2023. doi: 10.48550/ARXIV.2307.02486. URL https://doi.org/10. 579 48550/arXiv.2307.02486. 580
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 582 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024. 584
- Daniel Y Fu, Tri Dao, Khaled Kamal Saab, Armin W Thomas, Atri Rudra, and Christopher Re. 585 Hungry hungry hippos: Towards language modeling with state space models. In The Eleventh 586 International Conference on Learning Representations, 2023. URL https://openreview. 587 net/forum?id=COZDy0WYGg. 588
- 589 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Fos-590 ter, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework 592 for few-shot language model evaluation, 07 2024. URL https://zenodo.org/records/ 12608602.

594 Paolo Glorioso, Quentin Anthony, Yury Tokpanov, James Whittington, Jonathan Pilault, Adam 595 Ibrahim, and Beren Millidge. Zamba: A compact 7b ssm hybrid model, 2024. URL https: 596 //arxiv.org/abs/2405.16712. 597 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. arXiv 598 preprint arXiv:2312.00752, 2023. 600 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured 601 state spaces. arXiv preprint arXiv:2111.00396, 2021a. 602 Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Com-603 bining recurrent, convolutional, and continuous-time models with linear state space layers. Ad-604 vances in neural information processing systems, 34:572–585, 2021b. 605 606 Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameteriza-607 tion and initialization of diagonal state space models. In S. Koyejo, S. Mo-608 hamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems, volume 35, pp. 35971-35983. Curran Associates, Inc., 609 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/ 610 file/e9a32fade47b906de908431991440f7c-Paper-Conference.pdf. 611 612 Diagonal state spaces are as effec-Ankit Gupta, Albert Gu, and Jonathan Berant. 613 In S. Koyejo, S. Mohamed, A. Agarwal, D. Beltive as structured state spaces. 614 grave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing 615 Systems, volume 35, pp. 22982–22994. Curran Associates, Inc., 2022. URL 616 https://proceedings.neurips.cc/paper_files/paper/2022/file/ 617 9156b0f6dfa9bbd18c79cc459ef5d61c-Paper-Conference.pdf. 618 Dongchen Han, Xuran Pan, Yizeng Han, Shiji Song, and Gao Huang. Flatten transformer: Vi-619 sion transformer using focused linear attention. In Proceedings of the IEEE/CVF International 620 Conference on Computer Vision (ICCV), pp. 5961–5971, October 2023. 621 622 Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In International Conference on Learning Representations, 2017. URL https://openreview. 623 net/forum?id=rkE3y85ee. 624 625 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are 626 RNNs: Fast autoregressive transformers with linear attention. In International conference on 627 machine learning, pp. 5156-5165. PMLR, 2020. 628 Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. RACE: Large-scale 629 ReAding comprehension dataset from examinations. In Proceedings of the 2017 Conference 630 on Empirical Methods in Natural Language Processing, pp. 785-794, Copenhagen, Denmark, 631 September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1082. URL 632 https://aclanthology.org/D17-1082. 633 634 Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, 635 Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann 636 Fleureau, Guillaume Lample, and Stanislas Polu. Numinamath. [https://huggingface.

co/AI-MO/NuminaMath-CoT](https://github.com/project-numina/

637

638

aimo-progress-prize/blob/main/report/numina_dataset.pdf),2024.

639 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao 640 Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, 641 Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João 642 Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Lo-643 gesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra 644 Murthy, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, 645 Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, 646 Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex 647 Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva

648 649 650	Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder: may the source be with you! 2023.
651 652 653 654	Opher Lieber, Barak Lenz, Hofit Bata, Gal Cohen, Jhonathan Osin, Itay Dalmedigos, Erez Safahi, Shaked Meirom, Yonatan Belinkov, Shai Shalev-Shwartz, et al. Jamba: A hybrid transformer- mamba language model. <i>arXiv preprint arXiv:2403.19887</i> , 2024.
655 656	Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. <i>arXiv preprint arXiv:2109.07958</i> , 2021.
657 658	I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
659 660 661	Harsh Mehta, Ankit Gupta, Ashok Cutkosky, and Behnam Neyshabur. Long range language model- ing via gated state spaces. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=5MkYIYCbva.
662 663 664	William Merrill, Jackson Petty, and Ashish Sabharwal. The illusion of state in state-space models. <i>arXiv preprint arXiv:2404.08819</i> , 2024.
665 666	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In <i>EMNLP</i> , 2018.
667 668	Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking the potential of slms in grade school math. <i>arXiv preprint arXiv:2402.14830</i> , 2024.
670 671 672 673	Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Callum Birch-Sykes, Michael Wornow, Aman Patel, Clayton Rabideau, Stefano Massaroli, Yoshua Bengio, Stefano Ermon, Stephen A. Baccus, and Chris Ré. Hyenadna: Long-range genomic sequence modeling at single nucleotide resolution, 2023. URL https://arxiv.org/abs/2306.15794.
674 675	Keiran Paster, Marco Dos Santos, Zhangir Azerbayev, and Jimmy Ba. Openwebmath: An open dataset of high-quality mathematical web text, 2023.
676 677 678 679 680	Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional language models. In <i>International Conference on Machine Learning</i> , pp. 28043–28078. PMLR, 2023.
681 682 683	Zhen Qin, Weixuan Sun, Hui Deng, Dongxu Li, Yunshen Wei, Baohong Lv, Junjie Yan, Ling- peng Kong, and Yiran Zhong. cosformer: Rethinking softmax in attention. <i>arXiv preprint</i> <i>arXiv:2202.08791</i> , 2022.
684 685	Pranav Rajpurkar, Jian Zhang, and Percy Liang. Know what you don't know: Unanswerable questions for squad. In <i>ACL 2018</i> , 2018.
687 688	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 64(9):99–106, 2021.
689	Noam Shazeer. Glu variants improve transformer. arXiv preprint arXiv:2002.05202, 2020.
690 691 692 693	Zhuoran Shen, Mingyuan Zhang, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Efficient attention: Attention with linear complexities. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pp. 3531–3539, 2021.
694 695 696	Jimmy T.H. Smith, Andrew Warrington, and Scott Linderman. Simplified state space layers for se- quence modeling. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=Ai8Hw3AXqks.
697 698 699	Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. <i>Neurocomputing</i> , 568:127063, 2024.
700 701	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.

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 705 706 707

708

709

710

- A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.
- Roger Waleffe, Wonmin Byeon, Duncan Riach, Brandon Norick, Vijay Korthikanti, Tri Dao, Albert Gu, Ali Hatamizadeh, Sudhakar Singh, Deepak Narayanan, et al. An empirical study of mambabased language models. *arXiv preprint arXiv:2406.07887*, 2024.
- Kaiyue Wen, Xingyu Dang, and Kaifeng Lyu. Rnns are not transformers (yet): The key bottleneck
 on in-context retrieval. *arXiv preprint arXiv:2402.18510*, 2024.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 4791–4800. Association for Computational Linguistics, 2019. URL https://aclanthology.org/P19-1472.
- Biao Zhang and Rico Sennrich. Root mean square layer normalization. Advances in Neural Information Processing Systems, 32, 2019.
- Guangxiang Zhao, Junyang Lin, Zhiyuan Zhang, Xuancheng Ren, Qi Su, and Xu Sun. Explicit sparse transformer: Concentrated attention through explicit selection. arXiv preprint arXiv:1912.11637, 2019.
- 728 729

730

A DATASETS

731 Our training data comprises a diverse set of datasets, carefully curated to ensure breadth and depth 732 across various domains. This diverse collection includes specialized mathematics datasets (Meta-733 MathQA Yu et al. (2023), NuminaMath-CoT LI et al. (2024), OpenWebMath Paster et al. (2023), 734 Orca-Math Mitra et al. (2024)), high-quality web data (Fineweb-Edu-dedup Ben Allal et al. (2024)), 735 synthetic data (Cosmopedia-v2 Ben Allal et al. (2024)), code data (Starcoderdata-python-edu Li 736 et al. (2023)), and general knowledge sources (Wikipedia). This comprehensive approach aims to 737 enable our model to handle a wide array of language modeling tasks. The inclusion of both domain-738 specific and broad-coverage datasets is designed to enhance the model's versatility and robustness 739 across language modeling tasks.

All datasets were tokenized using the LLama3's tokenizer Dubey et al. (2024), resulting in 300B tokens.

The training data size varies by model scale: the 190M model is trained on 27 billion tokens (exclusively from Cosmopedia-v2), while the 450M and 1.3B models are trained on 100 billion tokens sampled from the combination of datasets mentioned above. Below are detailed descriptions of each dataset used:

746 747 748

749

750

751

752

- 1. **MetaMathQA** Yu et al. (2023): A comprehensive mathematics dataset designed to enhance the model's mathematical reasoning and problem-solving abilities.
- 2. **NuminaMath-CoT** LI et al. (2024): A chain-of-thought mathematics dataset that promotes step-by-step reasoning in mathematical problem-solving.
- 3. **Cosmopedia-v2** Ben Allal et al. (2024): A large-scale synthetic dataset for pre-training, consisting of over 39 million textbooks, blog posts, and stories.
- Fineweb-Edu-dedup Ben Allal et al. (2024): A high-quality subset of the FineWeb-Edu dataset, containing 220 billion tokens of educational web pages. This dataset was filtered using an educational quality classifier to retain only the most valuable educational content.

5. OpenWebMath Paster et al. (2023): A diverse collection of mathematical content from over 130,000 different domains, including forums, educational pages, and blogs. It covers mathematics, physics, statistics, computer science, and related fields. 6. Starcoderdata-Edu Li et al. (2023): A subset of the Starcoder dataset, specifically filtered for high-quality educational content related to Python programming. This dataset aims to enhance the model's coding capabilities. 7. Orca-Math Mitra et al. (2024): A dataset focused on mathematical word problems, de-signed to improve the model's ability to interpret and solve practical mathematical scenar-ios. 8. Wikipedia: An English Wikipedia dataset providing a broad range of general knowledge across various topics.