
Retentive Network

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this work, we propose Retentive Network (RETNET) as a foundation architecture
2 for large language models, simultaneously achieving training parallelism, low-cost
3 inference, and good performance. We theoretically derive the connection between
4 recurrence and attention. Then we propose the retention mechanism for sequence
5 modeling, which supports three computation paradigms, i.e., parallel, recurrent,
6 and chunkwise recurrent. Specifically, the parallel representation allows for training
7 parallelism. The recurrent representation enables low-cost $O(1)$ inference, which
8 improves decoding throughput, latency, and GPU memory without sacrificing
9 performance. The chunkwise recurrent representation facilitates efficient long-
10 sequence modeling with linear complexity, where each chunk is encoded parallelly
11 while recurrently summarizing the chunks. Experimental results on language
12 modeling show that RETNET achieves favorable scaling results, parallel training,
13 low-cost deployment, and efficient inference.

14 1 Introduction

15 Transformer [51] has become the de facto architecture for large language models, which was initially
16 proposed to overcome the sequential training issue of recurrent models [25]. However, training
17 parallelism of Transformers is at the cost of inefficient inference, because of the $O(N)$ complexity per
18 step and memory-bound key-value cache [42], which renders Transformers unfriendly to deployment.
19 The growing sequence length increases GPU memory consumption as well as latency and reduces
20 inference speed. Numerous efforts have continued to develop the next-generation architecture, aiming
21 at retaining training parallelism and competitive performance as Transformers while having efficient
22 $O(1)$ inference. It is challenging to achieve the above goals simultaneously.

23 There have been three main strands of research. First, linearized attention [27, 37] approximates
24 standard attention scores $\exp(\mathbf{q} \cdot \mathbf{k})$ with kernels $\phi(\mathbf{q}) \cdot \phi(\mathbf{k})$, so that autoregressive inference can
25 be rewritten in a recurrent form. However, the modeling capability and performance are worse than
26 Transformers, which hinders the method’s popularity. The second strand returns to recurrent models
27 for efficient inference while sacrificing training parallelism. As a remedy, element-wise operators [36]
28 are used for acceleration, however, representation capacity and performance are harmed. The third
29 line explores replacing attention with other mechanisms, such as S4 [20], and its variants [11, 38].
30 None of the previous work can achieve strong performance and efficient inference at the same time
31 compared to Transformers.

32 In this work, we propose retentive networks (RetNet), achieving low-cost inference, efficient long-
33 sequence modeling, Transformer-comparable performance, and parallel model training simultane-
34 ously. Specifically, we introduce a multi-scale retention mechanism to substitute multi-head attention,
35 which has three computation paradigms, i.e., parallel, recurrent, and chunkwise recurrent repre-
36 sentations. First, the parallel representation empowers training parallelism to utilize GPU devices
37 fully. Second, the recurrent representation enables efficient $O(1)$ inference in terms of memory
38 and computation. The deployment cost and latency can be significantly reduced. Moreover, the
39 implementation is greatly simplified without key-value cache tricks. Third, the chunkwise recurrent

40 representation can perform efficient long-sequence modeling. We parallelly encode each local block
 41 for computation speed while recurrently encoding the global blocks to save GPU memory.

42 We compare RetNet with Transformer and its variants. Experimental results on language modeling
 43 show that RetNet is consistently competitive in terms of both scaling curves and in-context learning.
 44 Moreover, the inference cost of RetNet is length-invariant. For a 7B model and 8k sequence
 45 length, RetNet decodes $8.4\times$ faster and saves 70% of memory than Transformers with key-value
 46 caches. During training, RetNet also achieves $3\times$ acceleration than standard Transformer with
 47 highly-optimized FlashAttention-2 [10]. Besides, RetNet’s inference latency is insensitive to batch
 48 size, allowing enormous throughput. The intriguing properties make RetNet a potential candidate to
 49 replace Transformer for large language models.

50 2 Retentive Network

51 Retentive network (RetNet) is stacked with L identical blocks, which follows a similar layout (i.e.,
 52 residual connection, and pre-LayerNorm) as in Transformer [51]. Each RetNet block contains two
 53 modules: a multi-scale retention (MSR) module, and a feed-forward network (FFN) module. We
 54 introduce the MSR module in the following sections. Given an input sequence $x = x_1 \cdots x_{|x|}$,
 55 RetNet encodes the sequence in an autoregressive way. The input vectors $\{\mathbf{x}_i\}_{i=1}^{|x|}$ is first packed
 56 into $X^0 = [\mathbf{x}_1, \cdots, \mathbf{x}_{|x|}] \in \mathbb{R}^{|x| \times d_{\text{model}}}$, where d_{model} is hidden dimension. Then we compute
 57 contextualized vector representations $X^l = \text{RetNet}_l(X^{l-1}), l \in [1, L]$.

58 2.1 Retention

59 In this section, we introduce the retention mechanism that has a dual form of recurrence and
 60 parallelism. So we can train the models in a parallel way while recurrently conducting inference.

61 Consider a sequence modeling problem that maps $v(n) \mapsto o(n)$ through states \mathbf{s}_n . Let v_n, o_n denote
 62 $v(n), o(n)$ for simplicity. We formulate the mapping in a recurrent manner:

$$\begin{aligned} \mathbf{s}_n &= A\mathbf{s}_{n-1} + K_n^\top v_n, \quad A \in \mathbb{R}^{d \times d}, \quad K_n \in \mathbb{R}^{1 \times d} \\ o_n &= Q_n \mathbf{s}_n = \sum_{m=1}^n Q_n A^{n-m} K_m^\top v_m, \quad Q_n \in \mathbb{R}^{1 \times d} \end{aligned} \quad (1)$$

63 where we map v_n to the state vector \mathbf{s}_n , and then implement a linear transform to encode sequence
 64 information recurrently. Next, we make the projection Q_n, K_n content-aware:

$$Q = XW_Q, \quad K = XW_K \quad (2)$$

65 where $W_Q, W_K \in \mathbb{R}^{d \times d}$ are learnable matrices.

66 We diagonalize the matrix $A = \Lambda(\gamma e^{i\theta})\Lambda^{-1}$, where $\gamma, \theta \in \mathbb{R}^d$. Then we obtain $A^{n-m} =$
 67 $\Lambda(\gamma e^{i\theta})^{n-m}\Lambda^{-1}$. By absorbing Λ into W_Q and W_K , we can rewrite Equation (1) as:

$$\begin{aligned} o_n &= \sum_{m=1}^n Q_n (\gamma e^{i\theta})^{n-m} K_m^\top v_m \\ &= \sum_{m=1}^n (Q_n (\gamma e^{i\theta})^n) (K_m (\gamma e^{i\theta})^{-m})^\top v_m \end{aligned} \quad (3)$$

68 where $Q_n (\gamma e^{i\theta})^n, K_m (\gamma e^{i\theta})^{-m}$ is known as xPos [45], i.e., a relative position embedding proposed
 69 for Transformer. We further simplify γ as a scalar, Equation (3) becomes:

$$o_n = \sum_{m=1}^n \gamma^{n-m} (Q_n e^{in\theta}) (K_m e^{im\theta})^\dagger v_m \quad (4)$$

70 where \dagger is the conjugate transpose. The formulation is easily parallelizable within training instances.

71 In summary, we start with recurrent modeling as shown in Equation (1), and then derive its parallel
 72 formulation in Equation (4). We consider the original mapping $v(n) \mapsto o(n)$ as vectors and obtain
 73 the retention mechanism as follows.

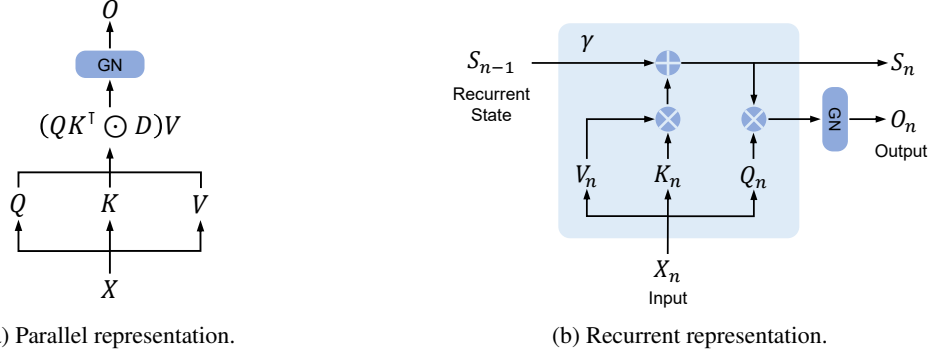


Figure 1: RetNet has three equivalent computation paradigms, i.e., parallel, recurrent, and chunkwise recurrent representations. Given the same input, three paradigms obtain the same output. “GN” is short for GroupNorm.

74 **The Parallel Representation of Retention** As shown in Figure 1a, the retention layer is defined as:

$$\begin{aligned}
 Q &= (XW_Q) \odot \Theta, \quad K = (XW_K) \odot \bar{\Theta}, \quad V = XW_V \\
 \Theta_n &= e^{in\theta}, \quad D_{nm} = \begin{cases} \gamma^{n-m}, & n \geq m \\ 0, & n < m \end{cases} \\
 \text{Retention}(X) &= (QK^\top \odot D)V
 \end{aligned} \tag{5}$$

75 where $D \in \mathbb{R}^{|x| \times |x|}$ combines causal masking and exponential decay along relative distance as one
 76 matrix, and $\bar{\Theta}$ is the complex conjugate of Θ . In practice, we map $Q, K \in \mathbb{R}^d \rightarrow \mathbb{C}^{d/2}$, add the
 77 complex position embedding Θ , then map them back to \mathbb{R}^d , following the implementation trick as in
 78 LLaMA [48, 44]. Similar to self-attention, the parallel representation enables us to train the models
 79 with GPUs efficiently.

80 **The Recurrent Representation of Retention** As shown in Figure 1b, the proposed mechanism can
 81 also be written as recurrent neural networks (RNNs), which is favorable for inference. For the n -th
 82 timestep, we recurrently obtain the output as:

$$\begin{aligned}
 S_n &= \gamma S_{n-1} + K_n^\top V_n \\
 \text{Retention}(X_n) &= Q_n S_n, \quad n = 1, \dots, |x|
 \end{aligned} \tag{6}$$

83 where Q, K, V, γ are the same as in Equation (5).

84 **The Chunkwise Recurrent Representation of Retention** A hybrid form of parallel representation
 85 and recurrent representation is available to accelerate training, especially for long sequences. We
 86 divide the input sequences into chunks. Within each chunk, we follow the parallel representation
 87 (Equation (5)) to conduct computation. In contrast, cross-chunk information is passed following the
 88 recurrent representation (Equation (6)). Specifically, let B denote the chunk length. We compute the
 89 retention output of the i -th chunk via:

$$\begin{aligned}
 Q_{[i]} &= Q_{Bi:B(i+1)}, \quad K_{[i]} = K_{Bi:B(i+1)}, \quad V_{[i]} = V_{Bi:B(i+1)} \\
 R_i &= K_{[i]}^\top (V_{[i]} \odot \zeta) + \gamma^B R_{i-1}, \quad \zeta_{ij} = \gamma^{B-i-1} \\
 \text{Retention}(X_{[i]}) &= \underbrace{(Q_{[i]} K_{[i]}^\top \odot D) V_{[i]}}_{\text{Inner-Chunk}} + \underbrace{(Q_{[i]} R_{i-1}) \odot \xi}_{\text{Cross-Chunk}}, \quad \xi_{ij} = \gamma^{i+1}
 \end{aligned} \tag{7}$$

90 where $[i]$ indicates the i -th chunk, i.e., $x_{[i]} = [x_{(i-1)B+1}, \dots, x_{iB}]$. The proof of the equivalence
 91 between recurrent representation and chunkwise recurrent representation is described in Appendix B.

92 2.2 Gated Multi-Scale Retention

93 We use $h = d_{\text{model}}/d$ retention heads in each layer, where d is the head dimension. The heads use
 94 different parameter matrices $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$. Moreover, multi-scale retention (MSR) assigns

```

def ParallelRetention(
q, k, v, # bsz * num_head * len * qkv_dim
decay_mask): # num_head * len * len
retention = q @ k.transpose(-1, -2)
retention = retention * decay_mask
output = retention @ v
output = group_norm(output)
return output

```

```

def RecurrentRetention(
q, k, v, # bsz * num_head * qkv_dim
past_kv, # bsz * num_head * qk_dim * v_dim
decay): # num_head * 1 * 1
current_kv = decay * past_kv + k.unsqueeze(-1) * v.
unsqueeze(-2)
output = torch.sum(q.unsqueeze(-1) * current_kv,
dim=-2)
output = group_norm(output)
return output, current_kv

```

```

def ChunkwiseRetention(
q, k, v, # bsz * num_head * chunk_size *
qkv_dim
past_kv, # bsz * num_head * qk_dim *
v_dim
decay_mask, # num_head * chunk_size *
chunk_size
chunk_decay, # num_head * 1 * 1
inner_decay): # num_head * chunk_size
retention = q @ k.transpose(-1, -2)
retention = retention * decay_mask
inner_retention = retention @ v
cross_retention = (q @ past_kv) *
inner_decay
retention = inner_retention +
cross_retention
output = group_norm(retention)
current_kv = chunk_decay * past_kv + k.
transpose(-1, -2) @ v
return output, current_kv

```

Figure 2: Pseudocode for the three computation paradigms of retention. Parallel implementation enables training parallelism to fully utilize GPUs. Recurrent paradigm enables low-cost inference. Chunkwise retention combines the above advantages (i.e., parallel within each chunk and recurrent across chunks), which has linear memory complexity for long sequences.

95 different γ for each head. For simplicity, we set γ identical among different layers and keep them
96 fixed. In addition, we add a swish gate [23, 40] to increase the non-linearity of retention layers.
97 Formally, given input X , we define the layer as:

$$\begin{aligned}
\gamma &= 1 - 2^{-5-\text{arange}(0,h)} \in \mathbb{R}^h \\
\text{head}_i &= \text{Retention}(X, \gamma_i) \\
Y &= \text{GroupNorm}_h(\text{Concat}(\text{head}_1, \dots, \text{head}_h)) \\
\text{MSR}(X) &= (\text{swish}(XW_G) \odot Y)W_O
\end{aligned} \tag{8}$$

98 where $W_G, W_O \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$ are learnable parameters, and GroupNorm [53] normalizes the
99 output of each head, following SubLN proposed in [43]. Notice that the heads use multiple γ scales,
100 which results in different variance statistics. So we normalize the head outputs separately.

101 The pseudocode of retention is summarized in Figure 2.

102 **Retention Score Normalization** We utilize the scale-invariant nature of GroupNorm to improve the
103 numerical precision of retention layers. Specifically, multiplying a scalar value within GroupNorm
104 does not affect outputs and backward gradients, i.e., $\text{GroupNorm}(\alpha * \text{head}_i) = \text{GroupNorm}(\text{head}_i)$.
105 We implement three normalization factors in Equation (5). First, we normalize QK^\top as QK^\top / \sqrt{a} .
106 Second, we replace D with $\tilde{D}_{nm} = D_{nm} / \sqrt{\sum_{i=1}^n D_{ni}}$. Third, let R denote the retention scores
107 $R = QK^\top \odot D$, we normalize it as $\tilde{R}_{nm} = R_{nm} / \max(\sum_{i=1}^n |R_{ni}|, 1)$. Then the retention output
108 becomes $\text{Retention}(X) = \tilde{R}V$. The above tricks do not affect the final results while stabilizing the
109 numerical flow of both forward and backward passes, because of the scale-invariant property.

110 2.3 Overall Architecture of Retention Networks

111 For an L -layer retention network, we stack multi-scale retention (MSR) and feed-forward network
112 (FFN) to build the model. Formally, the input sequence $\{x_i\}_{i=1}^{|x|}$ is transformed into vectors by a
113 word embedding layer. We use the packed embeddings $X^0 = [\mathbf{x}_1, \dots, \mathbf{x}_{|x|}] \in \mathbb{R}^{|x| \times d_{\text{model}}}$ as the
114 input and compute the model output X^L :

$$\begin{aligned}
Y^l &= \text{MSR}(\text{LN}(X^l)) + X^l \\
X^{l+1} &= \text{FFN}(\text{LN}(Y^l)) + Y^l
\end{aligned} \tag{9}$$

115 where $\text{LN}(\cdot)$ is LayerNorm [3]. The FFN part is computed as $\text{FFN}(X) = \text{gelu}(XW_1)W_2$, where
116 W_1, W_2 are parameter matrices.

117 **Training** We use the parallel (Equation (5)) and chunkwise recurrent (Equation (7)) representations
 118 during the training process. The parallelization within sequences or chunks efficiently utilizes
 119 GPUs to accelerate computation. More favorably, chunkwise recurrence is especially useful for
 120 long-sequence training, which is efficient in terms of both FLOPs and memory consumption.

121 **Inference** The recurrent representation (Equation (6)) is employed during inference, which nicely
 122 fits autoregressive decoding. The $O(1)$ complexity reduces memory and inference latency while
 123 achieving equivalent results.

124 3 Experiments

125 We perform language modeling experiments to evaluate RetNet. First, we present the scaling curves
 126 of Transformer and RetNet. Second, we follow the training settings of StableLM-4E1T [50] to
 127 compare with open-source Transformer models in downstream benchmarks. Moreover, for training
 128 and inference, we compare speed, memory consumption, and latency. The training corpus is a curated
 129 compilation of The Pile [16], C4 [14], and The Stack [29].

130 3.1 Comparison with Transformer Variants

131 We compare RetNet with various efficient Transformer variants, including RWKV [36], H3 [11],
 132 Hyena [38], and Mamba [19]. We use LLaMA [48] architecture, including RMSNorm [59] and
 133 SwiGLU [40, 7] module, as the Transformer backbone, which shows better performance and stability.
 134 Consequently, other variants follow these settings. Specifically, Mamba does not have FFN layers so
 135 we only implement RMSNorm. For RetNet, the FFN intermediate dimension is $\frac{5}{3}d$ and the value
 136 dimensions in W_G, W_V, W_O are also $\frac{5}{3}d$, where the overall parameters are still $12d^2$. All models
 137 have 400M parameters with 24 layers and a hidden dimension of 1024. For H3, we set the head
 138 dimension to 8. For RWKV, we use the TimeMix module to substitute self-attention layers while
 139 keeping FFN layers consistent with other models for fair comparisons. We train the models with 40k
 140 steps with a batch size of 0.25M tokens.

141 **Fine-Grained Language Modeling Evaluation** As shown in Table 1, we first report the language
 142 modeling perplexity of validation sets. Besides the overall validation set, following [2], we divide
 143 perplexity into “AR-Hit” and “First Occur”. Specifically, AR-Hit contains the predicted tokens that
 144 are previously seen bigrams in the previous context, which evaluates the associative recall ability.
 145 “First Occur” has the predicted tokens that can not be recalled from the context. Among various
 146 Transformer variants, RetNet outperforms previous methods on both “AR-Hit” and “First Occur”
 147 splits, which is important for real-world use cases.

148 **Knowledge-Intensive Tasks** We also evaluate Massive Multitask Language Understanding
 149 (MMLU; [24]) answer perplexity to evaluate models on knowledge-intensive tasks. We report
 150 the average perplexity of the correct answers, i.e., given input [Question, “Answer:”, Correct
 151 Answer], we calculate the perplexity of the “Correct Answer” part. RetNet achieves competitive
 152 results among the architectures.

	Language Modeling			MMLU				
	Valid. Set	AR-Hit	First-Occur	STEMs	Humanites	Social-Sci.	Others	Avg
Transformer [51]	3.320	1.118	3.826	0.584	0.229	0.279	0.402	0.356
<i>Transformer Variants</i>								
Hyena [38]	3.545	1.799	3.947	1.125	0.576	0.654	0.819	0.767
RWKV [36]	3.497	1.706	3.910	1.156	0.609	0.617	0.781	0.768
Mamba [19]	3.379	1.322	3.852	0.668	0.288	0.300	0.425	0.403
H3 [11]	3.563	1.722	3.986	1.169	0.532	0.637	0.792	0.752
RetNet	3.360	1.264	3.843	0.577	0.263	0.280	0.384	0.362

Table 1: Perplexity results on language modeling and MMLU [24] answers. We use the augmented Transformer architecture proposed in LLaMA [48] for reference. For language modeling, we report perplexity on both the overall validation set and fine-grained diagnosis sets [2], i.e., “AR-Hit” evaluates the associative recall capability, and “First-Occur” indicates the regular language modeling performance. Besides, we evaluate the answer perplexity of MMLU subsets.

153 3.2 Language Modeling Evaluation with Various Model Sizes

154 We train language models with various sizes (i.e., 1.3B,
155 2.7B, and 6.7B) from scratch. The training batch size
156 is 4M tokens with 2048 maximal length. We train the
157 models with 25k steps. The detailed hyper-parameters are
158 described in Appendix E. We train the models with 512
159 AMD MI200 GPUs.

160 Figure 3 reports perplexity on the validation set for the
161 language models based on Transformer and RetNet. We
162 present the scaling curves with three model sizes, i.e.,
163 1.3B, 2.7B, and 6.7B. RetNet achieves comparable results
164 with Transformers. More importantly, the results indicate
165 that RetNet is favorable in terms of size scaling. In addition
166 to performance, RetNet training is quite stable in our
167 experiments. Experimental results show that RetNet is a
168 strong competitor to Transformer for large language models.
169 Empirically, we find that RetNet starts to outperform
170 Transformer when the model size is larger than 2B.

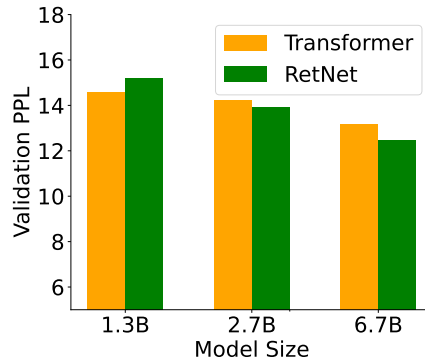


Figure 3: Validation perplexity (PPL) decreases along with scaling up the model size.

171 3.3 Long-Context Evaluation

172 We evaluate long-context modeling on the ZeroSCROLLS [41] benchmark. We train a hybrid model
173 of size 2.7B, RetNet+, which stacks the attention and retention layers. Specifically, we insert one
174 attention layer after every 3 retention layers. We follow most configurations of the 2.7B model as in
175 Section 3.2. We scale the number of training tokens to 420B tokens. The batch size is 4M tokens.
176 We first train the model with 4K length and then extend the sequence length to 16K for the last 50B
177 training tokens. The rotation base scaling [55] is used for length extension.

178 Figure 4 reports the answer perplexity given various lengths of input document. It shows that both
179 Transformer and RetNet+ perform better with longer input documents. The results indicate that the
180 language models successfully utilize the long-distance context. Notice that the 12K and 16K results
181 in Qasper are similar because the lengths of most documents are shorter than 16K. Moreover, RetNet+
182 obtains competitive results compared with Transformer for long-context modeling. Meanwhile,
183 retention has better training and inference efficiency.

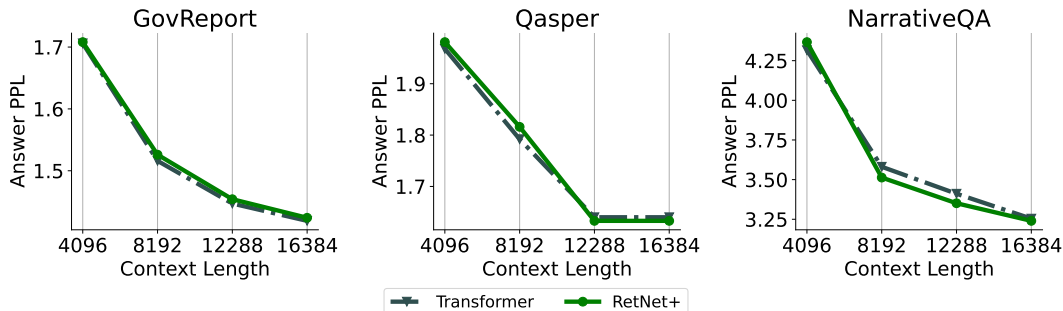
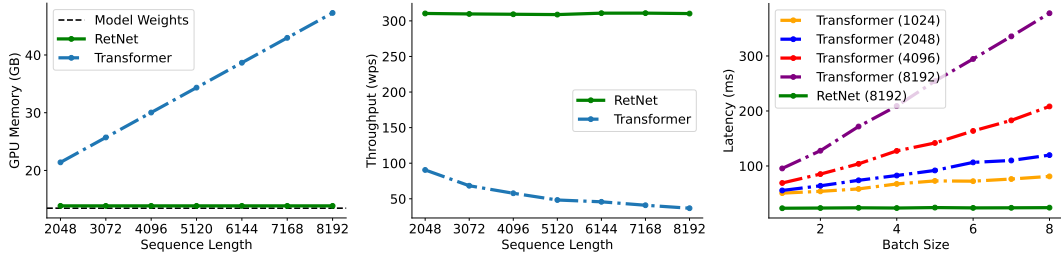


Figure 4: Answer perplexity decreases along with longer input documents. Transformer and RetNet+ obtain comparable performance for long-context modeling on the ZeroSCROLLS [41] benchmark.

184 3.4 Inference Cost

185 As shown in Figure 5, we compare memory cost, throughput, and latency of Transformer and RetNet
186 during inference. Transformers reuse KV caches of previously decoded tokens. RetNet uses the
187 recurrent representation as described in Equation (6). We evaluate the 6.7B model on the A100-80GB
188 GPU. Figure 5 shows that RetNet outperforms Transformer in terms of inference cost.

189 **Memory** As shown in Figure 5a, the memory cost of Transformer increases linearly due to KV
190 caches. In contrast, the memory consumption of RetNet remains consistent even for long sequences,



(a) GPU memory cost with varying sequence length. (b) Inference throughput with varying sequence length. (c) Inference latency with different batch sizes.

Figure 5: Inference cost of Transformer and RetNet with a model size of 6.7B. RetNet outperforms Transformers in terms of memory consumption, throughput, and latency.

191 requiring much less GPU memory to host RetNet. The additional memory consumption of RetNet is
 192 almost negligible (i.e., about 3%) while the model weights occupy 97%.

193 **Throughput** As presented in Figure 5b, the throughput of Transformer drops along with the
 194 decoding length increases. In comparison, RetNet has higher and length-invariant throughput during
 195 decoding, by utilizing the recurrent representation of retention.

196 **Latency** Latency is an important metric in deployment that greatly affects the user experience.
 197 We report the decoding latency in Figure 5c. Experimental results show that increasing batch size
 198 renders the Transformer’s latency larger. Moreover, the latency of Transformers grows faster with
 199 longer input. In order to make latency acceptable, we have to restrict the batch size, which harms the
 200 overall inference throughput of Transformers. By contrast, RetNet’s decoding latency outperforms
 201 Transformers and stays almost the same across different batch sizes and input lengths.

202 3.5 Training Throughput

203 Figure 6 compares the training throughput of Trans-
 204 former and RetNet, where the training sequence lengths
 205 range from 8192 to 65536. The model size is 3.5B,
 206 where the hidden dimension is 3072 and the layer size
 207 is 28. We use highly optimized FlashAttention-2 [10]
 208 for Transformers. In comparison, we implement chunk
 209 recurrent representation (Equation (7)) using Triton [46],
 210 where the computation is both memory-friendly and
 211 computationally efficient. The chunk size is set to 256.
 212 We evaluate the results with eight Nvidia H100-80GB
 213 GPUs because FlashAttention-2 is highly optimized for
 214 H100 cards.

215 Experimental results show that RetNet has higher training
 216 throughput than Transformers. The acceleration ratio increases as the sequence length is longer.
 217 When the training length is 64k, RetNet’s throughput is about 3 times than Transformer’s.

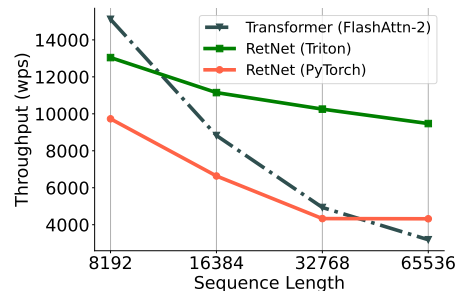


Figure 6: Training throughput (word per second; wps) of Transformer with FlashAttention-2 [10] and RetNet.

218 3.6 Zero-Shot and Few-Shot Evaluation on Downstream Tasks

219 We also compare the language models on a wide range of downstream tasks. We evaluate zero-shot
 220 and 4-shot learning with the 6.7B models. As shown in Table 2, the datasets include HellaSwag
 221 (HS; [57]), BoolQ [8], COPA [52], PIQA [6], Winograd, Winogrande [30], and StoryCloze (SC; [34]).
 222 The accuracy numbers are consistent with language modeling perplexity presented in Figure 3. RetNet
 223 achieves comparable performance with Transformer on zero-shot and in-context learning settings.

224 3.7 Ablation Studies

225 We ablate various design choices of RetNet and report the language modeling results in Table 3. The
 226 evaluation settings and metrics are the same as in Section 3.1.

	HS	BoolQ	COPA	PIQA	Winograd	Winogrande	SC	Avg
<i>Zero-Shot Performance</i>								
Transformer	55.9	62.0	69.0	74.6	69.5	56.5	75.0	66.07
RetNet	60.7	62.2	77.0	75.4	77.2	58.1	76.0	69.51
<i>Few-shot Performance (4-Shot)</i>								
Transformer	55.8	58.7	71.0	75.0	71.9	57.3	75.4	66.44
RetNet	60.5	60.1	78.0	76.0	77.9	59.9	75.9	69.76

Table 2: Zero-shot and few-shot learning performance. The language model size is 6.7B.

227 **Architecture** We ablate the swish gate and GroupNorm as described in Equation (8). Table 3
 228 shows that the above two components improve performance. First, the gating module is essential
 229 for enhancing non-linearity and improving model capability. Notice that we use the same parameter
 230 allocation as in Transformers after removing the gate. Second, group normalization in retention
 231 balances the variances of multi-head outputs, which improves training stability and language modeling
 232 results.

233 **Multi-Scale Decay** Equation (8) shows that we use different γ as the decay rates for the retention
 234 heads. In the ablation studies, we examine removing γ decay (i.e., “ $-\gamma$ decay”) and applying the
 235 same decay rate across heads (i.e., “ $-\text{multi-scale decay}$ ”). Specifically, ablating γ decay is equivalent
 236 to $\gamma = 1$. In the second setting, we set $\gamma = 1 - 2^{-6.5}$ for all heads. Table 3 indicates that both the
 237 decay mechanism and using multiple decay rates can improve the language modeling performance.

238 **Head Dimension** As indicated by the recurrent perspective of Equation (1), the head dimension
 239 implies the memory capacity of hidden states. In ablation, we reduce the default head dimension from
 240 256 to 64, i.e., 64 for queries and keys, and $\lfloor \frac{5}{3} \times 64 \rfloor \approx 108$ for values. We keep the hidden dimension
 241 d_{model} the same. Accordingly, we adjust the multi-scale decay as $\gamma = 1 - 2^{-5-\text{arange}(0,h)/4}$ to keep
 242 the same decay range. Table 3 shows that the larger head dimension achieves better performance.

	Language Modeling			MMLU				Avg
	Valid. Set	AR-Hit	First-Occur	STEMs	Humanites	Social-Sci.	Others	
RetNet	3.360	1.264	3.843	0.577	0.263	0.280	0.384	0.362
– swish gate	3.509	1.366	4.002	0.599	0.285	0.315	0.421	0.390
– GroupNorm	3.367	1.302	3.843	0.630	0.295	0.327	0.438	0.406
– γ decay	3.920	2.122	4.334	0.958	0.566	0.571	0.694	0.681
– multi-scale decay	3.524	1.768	3.928	0.921	0.433	0.471	0.590	0.582
Reduce head dim.	3.397	1.331	3.872	0.637	0.272	0.294	0.393	0.384

Table 3: Perplexity results on language modeling and MMLU [24] answers. For language modeling, we report perplexity on both the overall validation set and fine-grained diagnosis sets [2], i.e., “AR-Hit” evaluates the associative recall capability, and “First-Occur” indicates the regular language modeling performance. Besides, we evaluate the answer perplexity of the MMLU subsets.

243 3.8 Results on Vision Tasks

244 We also compare RetNet with vision Transformers [15, 47] in Table 4, where bidirectional en-
 245 coders are evaluated. Unlike causal language models, the vision encoders do not require recurrent
 246 representations. Specifically, we use retention as follows:

$$Q = (XW_Q) \odot \Theta, \quad K = (XW_K) \odot \bar{\Theta}, \quad V = XW_V$$

$$\text{Retention}(X) = (QK^T)V = Q(K^TV)$$

247 where multi-scale decay is removed in bidirectional computation. Notice that we can compute
 248 retention in different orders. Similar to linear attention [27], the $Q(K^TV)$ paradigm is an efficient
 249 operator in bidirectional settings, especially for high-resolution images.

250 We perform experiments on ImageNet-1K classification [13], COCO object detection [32], and
 251 ADE20K semantic segmentation [60]. We compare RetNet with DeiT [47] which is a well-tuned
 252 vision Transformer. Besides, we follow [21] and plug in a depth-wise convolution in experiments.
 253 We adopt the DeiT-M size, which has about 38M parameters. For ImageNet-1K image classification,

	ImageNet	COCO			ADE20K	
	Acc	AP ^b	AP ₅₀ ^b	AP ₇₅ ^b	mIoU	mAcc
DeiT [47]	80.76	0.458	0.678	0.502	43.52	55.08
RetNet	81.57	0.457	0.669	0.488	44.13	56.12

Table 4: Results on vision tasks, i.e., image classification (ImageNet), object detection (COCO), and semantic segmentation (ADE20K). RetNet achieves competitive performance with DeiT, which is a well-tuned vision Transformer.

we use AdamW [33] for 300 epochs, and 20 epochs of linear warm-up. The learning rate is 1×10^{-3} , the batch size is 1024, and the weight decay is 0.05. For COCO object detection, we use Mask R-CNN [22] as the task head, and the above models pre-trained on ImageNet as the backbone with 3x schedules. In ADE20K experiments, we use UperNet [54] as the segmentation head. The detailed configuration can be found in Appendix H.

Table 4 shows the results across various vision tasks. RetNet is competitive compared with DeiT. For classification and segmentation, RetNet is slightly better than DeiT, where RetNet achieves 0.81% accuracy improvement on ImageNet and 0.61% mIoU improvement on ADE20K. For object detection, the results are comparable.

4 Related Work

Numerous efforts are focused on reducing the quadratic complexity of attention mechanisms. Linear attention [27] uses various kernels $\phi^{(q_i)}\phi^{(k_j)}/\sum_{n=1}^{|x|}\phi^{(q_i)}\phi^{(k_n)}$ to replace the softmax function. In contrast, we reexamine sequence modeling from scratch, rather than aiming at approximating softmax. AFT [58] simplifies dot-product attention to element-wise and moves softmax to key vectors. RWKV [36] replaces AFT’s position embeddings with exponential decay and runs the models recurrently for training and inference. In comparison, retention preserves high-dimensional states to encode sequence information, which contributes to expressive ability and better performance. S4 [20] unifies convolution and recurrence format and achieves $O(N \log N)$ training complexity leveraging the FFT kernel. Unlike Equation (2), if Q_n and K_n are content-unaware, the formulation can be degenerated to S4 [20]. Hyena [38] generates the convolution kernels, achieving sub-quadratic training efficiency but keeping $O(N)$ complexity in single-step inference. Recently, most related work has focused on modifying γ in Equation (6) as a data-dependent variable, such as Mamba [19], GLA [56], Gateloop [28], and xLSTM [4]. Another strand explores hybrid architectures [31, 12] that interleave the above components with attention layers.

In addition, we discuss the training and inference efficiency of some related methods. Let D denote the hidden dimension, H the head dimension, and N the sequence length. For training, RWKV’s token-mixing complexity is $O(DN)$, and Mamba’s complexity is $O(DHN)$ with optimized CUDA kernels. Hyena’s is $O(DN \log N)$ with Fast Fourier Transform acceleration. In comparison, the chunk-wise recurrent representation is $O(DN(B + H))$, where B is the chunk size, and we usually set $H = 256$, $B \leq 512$. However, chunk-wise computation is highly parallelized, enabling efficient hardware usage. For large model size (i.e., larger D) or sequence length, the additional $b + h$ has negligible effects. For inference, among the efficient architectures compared, Hyena has the same complexity (i.e., $O(N)$ per step) as Transformer, while the others can perform $O(1)$ decoding.

5 Conclusion

We propose retentive networks (RetNet) for sequence modeling, which enables various representations, i.e., parallel, recurrent, and chunkwise recurrent. RetNet achieves significantly better inference efficiency (in terms of memory, speed, and latency), favorable training parallelization, and competitive performance compared with Transformers. The above advantages make RetNet an ideal successor to Transformers for large language models, especially considering the deployment benefits brought by the $O(1)$ inference complexity. In the future, we are interested in deploying RetNet on various edge devices, such as mobile phones.

295 **References**

- 296 [1] J. Ainslie, J. Lee-Thorp, M. de Jong, Y. Zemlyanskiy, F. Lebrón, and S. Sanghai. GQA: Training
297 generalized multi-query Transformer models from multi-head checkpoints. *arXiv preprint*
298 *arXiv:2305.13245*, 2023.
- 299 [2] S. Arora, S. Eyuboglu, A. Timalsina, I. Johnson, M. Poli, J. Zou, A. Rudra, and C. Ré. Zoology:
300 Measuring and improving recall in efficient language models. *arXiv preprint arXiv:2312.04927*,
301 2023.
- 302 [3] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*,
303 2016.
- 304 [4] M. Beck, K. Pöppel, M. Spanring, A. Auer, O. Prudnikova, M. Kopp, G. Klambauer, J. Brand-
305 stetter, and S. Hochreiter. xLSTM: Extended long short-term memory. *arXiv preprint*
306 *arXiv:2405.04517*, 2024.
- 307 [5] J. Berant, A. Chou, R. Frostig, and P. Liang. Semantic parsing on Freebase from question-
308 answer pairs. In *Proceedings of the 2013 Conference on Empirical Methods in Natural*
309 *Language Processing*, pages 1533–1544, Seattle, Washington, USA, Oct. 2013. Association for
310 Computational Linguistics.
- 311 [6] Y. Bisk, R. Zellers, R. L. Bras, J. Gao, and Y. Choi. Piqa: Reasoning about physical com-
312 monsense in natural language. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*,
313 2020.
- 314 [7] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W.
315 Chung, C. Sutton, S. Gehrmann, P. Schuh, K. Shi, S. Tsvyashchenko, J. Maynez, A. B. Rao,
316 P. Barnes, Y. Tay, N. M. Shazeer, V. Prabhakaran, E. Reif, N. Du, B. C. Hutchinson, R. Pope,
317 J. Bradbury, J. Austin, M. Isard, G. Gur-Ari, P. Yin, T. Duke, A. Levskaya, S. Ghemawat,
318 S. Dev, H. Michalewski, X. García, V. Misra, K. Robinson, L. Fedus, D. Zhou, D. Ippolito,
319 D. Luan, H. Lim, B. Zoph, A. Spiridonov, R. Sepassi, D. Dohan, S. Agrawal, M. Omernick,
320 A. M. Dai, T. S. Pillai, M. Pellat, A. Lewkowycz, E. O. Moreira, R. Child, O. Polozov, K. Lee,
321 Z. Zhou, X. Wang, B. Saeta, M. Díaz, O. Firat, M. Catasta, J. Wei, K. S. Meier-Hellstern,
322 D. Eck, J. Dean, S. Petrov, and N. Fiedel. PaLM: Scaling language modeling with pathways.
323 *ArXiv*, abs/2204.02311, 2022.
- 324 [8] C. Clark, K. Lee, M.-W. Chang, T. Kwiatkowski, M. Collins, and K. Toutanova. BoolQ:
325 Exploring the surprising difficulty of natural yes/no questions. In *Proceedings of the 2019*
326 *Conference of the North American Chapter of the Association for Computational Linguistics*,
327 pages 2924–2936, 2019.
- 328 [9] T. Computer. Redpajama-data: An open source recipe to reproduce llama training dataset, 2023.
329 URL <https://github.com/togethercomputer/RedPajama-Data>.
- 330 [10] T. Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. *arXiv*
331 *preprint arXiv:2307.08691*, 2023.
- 332 [11] T. Dao, D. Y. Fu, K. K. Saab, A. W. Thomas, A. Rudra, and C. Ré. Hungry hungry hippos:
333 Towards language modeling with state space models. *arXiv preprint arXiv:2212.14052*, 2022.
- 334 [12] S. De, S. L. Smith, A. Fernando, A. Botev, G. Cristian-Muraru, A. Gu, R. Haroun, L. Berrada,
335 Y. Chen, S. Srinivasan, G. Desjardins, A. Doucet, D. Budden, Y. W. Teh, R. Pascanu, N. D.
336 Freitas, and C. Gulcehre. Griffin: Mixing gated linear recurrences with local attention for
337 efficient language models. 2024.
- 338 [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical
339 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages
340 248–255. Ieee, 2009.
- 341 [14] J. Dodge, A. Marasović, G. Ilharco, D. Groeneveld, M. Mitchell, and M. Gardner. Documenting
342 large webtext corpora: A case study on the colossal clean crawled corpus. In *Conference on*
343 *Empirical Methods in Natural Language Processing*, 2021.

- 344 [15] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
345 M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for
346 image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 347 [16] L. Gao, S. Biderman, S. Black, L. Golding, T. Hoppe, C. Foster, J. Phang, H. He, A. Thite,
348 N. Nabeshima, et al. The Pile: An 800GB dataset of diverse text for language modeling. *arXiv*
349 *preprint arXiv:2101.00027*, 2020.
- 350 [17] L. Gao, J. Tow, B. Abbasi, S. Biderman, S. Black, A. DiPofi, C. Foster, L. Golding, J. Hsu,
351 A. Le Noac’h, H. Li, K. McDonell, N. Muennighoff, C. Ociepa, J. Phang, L. Reynolds,
352 H. Schoelkopf, A. Skowron, L. Sutawika, E. Tang, A. Thite, B. Wang, K. Wang, and A. Zou. A
353 framework for few-shot language model evaluation, 12 2023. URL [https://zenodo.org/
354 records/10256836](https://zenodo.org/records/10256836).
- 355 [18] X. Geng and H. Liu. Openllama: An open reproduction of llama, May 2023. URL [https://
356 //github.com/openlm-research/open_llama](https://github.com/openlm-research/open_llama).
- 357 [19] A. Gu and T. Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
358 *preprint arXiv:2312.00752*, 2023.
- 359 [20] A. Gu, K. Goel, and C. Ré. Efficiently modeling long sequences with structured state spaces.
360 *arXiv preprint arXiv:2111.00396*, 2021.
- 361 [21] D. Han, X. Pan, Y. Han, S. Song, and G. Huang. Flatten Transformer: Vision Transformer
362 using focused linear attention. In *Proceedings of the IEEE/CVF International Conference on*
363 *Computer Vision*, pages 5961–5971, 2023.
- 364 [22] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE*
365 *international conference on computer vision*, pages 2961–2969, 2017.
- 366 [23] D. Hendrycks and K. Gimpel. Gaussian error linear units (GELUs). *arXiv: Learning*, 2016.
- 367 [24] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring
368 massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- 369 [25] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9:1735–1780,
370 Nov. 1997.
- 371 [26] W. Hua, Z. Dai, H. Liu, and Q. Le. Transformer quality in linear time. In *International*
372 *Conference on Machine Learning*, pages 9099–9117. PMLR, 2022.
- 373 [27] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret. Transformers are rnns: Fast autoregressive
374 transformers with linear attention. In *International Conference on Machine Learning*, pages
375 5156–5165. PMLR, 2020.
- 376 [28] T. Katsch. Gateloop: Fully data-controlled linear recurrence for sequence modeling. *arXiv*
377 *preprint arXiv:2311.01927*, 2023.
- 378 [29] D. Kocetkov, R. Li, L. Ben Allal, J. Li, C. Mou, C. Muñoz Ferrandis, Y. Jernite, M. Mitchell,
379 S. Hughes, T. Wolf, D. Bahdanau, L. von Werra, and H. de Vries. The Stack: 3TB of permissively
380 licensed source code. *Preprint*, 2022.
- 381 [30] H. Levesque, E. Davis, and L. Morgenstern. The winograd schema challenge. In *Thirteenth*
382 *International Conference on the Principles of Knowledge Representation and Reasoning*, 2012.
- 383 [31] O. Lieber, B. Lenz, H. Bata, G. Cohen, J. Osin, I. Dalmedigos, E. Safahi, S. Meirum, Y. Belinkov,
384 S. Shalev-Shwartz, et al. Jamba: A hybrid Transformer-Mamba language model. *arXiv preprint*
385 *arXiv:2403.19887*, 2024.
- 386 [32] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick.
387 Microsoft COCO: Common objects in context. In *Computer Vision–ECCV 2014: 13th European*
388 *Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755.
389 Springer, 2014.

- 390 [33] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Conference*
391 *on Learning Representations*, 2019.
- 392 [34] N. Mostafazadeh, M. Roth, A. Louis, N. Chambers, and J. Allen. Lsdsem 2017 shared task: The
393 story cloze test. In *Proceedings of the 2nd Workshop on Linking Models of Lexical, Sentential*
394 *and Discourse-level Semantics*, pages 46–51, 2017.
- 395 [35] A. Orvieto, S. L. Smith, A. Gu, A. Fernando, C. Gulcehre, R. Pascanu, and S. De. Resurrecting
396 recurrent neural networks for long sequences. *ArXiv*, abs/2303.06349, 2023.
- 397 [36] B. Peng, E. Alcaide, Q. G. Anthony, A. Albalak, S. Arcadinho, H. Cao, X. Cheng, M. Chung,
398 M. Grella, G. Kranthikiran, X. He, H. Hou, et al. RWKV: Reinventing RNNs for the Transformer
399 era. *ArXiv*, abs/2305.13048, 2023.
- 400 [37] H. Peng, N. Pappas, D. Yogatama, R. Schwartz, N. A. Smith, and L. Kong. Random feature
401 attention. *arXiv preprint arXiv:2103.02143*, 2021.
- 402 [38] M. Poli, S. Massaroli, E. Nguyen, D. Y. Fu, T. Dao, S. Baccus, Y. Bengio, S. Ermon, and
403 C. Ré. Hyena hierarchy: Towards larger convolutional language models. *arXiv preprint*
404 *arXiv:2302.10866*, 2023.
- 405 [39] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang. SQuAD: 100,000+ questions for machine
406 comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in*
407 *Natural Language Processing*, pages 2383–2392, Austin, Texas, Nov. 2016. Association for
408 Computational Linguistics. doi: 10.18653/v1/D16-1264.
- 409 [40] P. Ramachandran, B. Zoph, and Q. V. Le. Swish: a self-gated activation function. *arXiv: Neural*
410 *and Evolutionary Computing*, 2017.
- 411 [41] U. Shaham, M. Ivgi, A. Efrat, J. Berant, and O. Levy. ZeroSCROLLS: A zero-shot benchmark
412 for long text understanding. *arXiv preprint arXiv:2305.14196*, 2023.
- 413 [42] N. M. Shazeer. Fast Transformer decoding: One write-head is all you need. *ArXiv*,
414 abs/1911.02150, 2019.
- 415 [43] M. Shoenberger, M. Patwary, R. Puri, P. LeGresley, J. Casper, and B. Catanzaro. Megatron-LM:
416 Training multi-billion parameter language models using model parallelism. *arXiv preprint*
417 *arXiv:1909.08053*, 2019.
- 418 [44] J. Su, Y. Lu, S. Pan, B. Wen, and Y. Liu. Roformer: Enhanced transformer with rotary position
419 embedding. *arXiv preprint arXiv:2104.09864*, 2021.
- 420 [45] Y. Sun, L. Dong, B. Patra, S. Ma, S. Huang, A. Benhaim, V. Chaudhary, X. Song, and F. Wei. A
421 length-extrapolatable transformer. In *Proceedings of the 61st Annual Meeting of the Association*
422 *for Computational Linguistics (Volume 1: Long Papers)*, pages 14590–14604, Toronto, Canada,
423 July 2023. Association for Computational Linguistics.
- 424 [46] P. Tillet and D. Cox. Triton: An intermediate language and compiler for tiled neural network
425 computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine*
426 *Learning and Programming Languages*, pages 10–19, 2019.
- 427 [47] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou. Training data-efficient
428 image transformers & distillation through attention. In *International conference on machine*
429 *learning*, pages 10347–10357. PMLR, 2021.
- 430 [48] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal,
431 E. Hambro, F. Azhar, et al. LLaMA: Open and efficient foundation language models. *arXiv*
432 *preprint arXiv:2302.13971*, 2023.
- 433 [49] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra,
434 P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv*
435 *preprint arXiv:2307.09288*, 2023.
- 436 [50] J. Tow, M. Bellagente, D. Mahan, and C. Riquelme. StableLM 3B 4E1T. [https://aka.ms/](https://aka.ms/StableLM-3B-4E1T)
437 [StableLM-3B-4E1T](https://aka.ms/StableLM-3B-4E1T), 2023.

- 438 [51] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and
439 I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*
440 *30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017,*
441 *Long Beach, CA, USA*, pages 6000–6010, 2017.
- 442 [52] A. Wang, Y. Pruksachatkun, N. Nangia, A. Singh, J. Michael, F. Hill, O. Levy, and S. R.
443 Bowman. SuperGLUE: A stickier benchmark for general-purpose language understanding
444 systems. *arXiv preprint arXiv:1905.00537*, 2019.
- 445 [53] Y. Wu and K. He. Group normalization. In *Proceedings of the European conference on computer*
446 *vision (ECCV)*, pages 3–19, 2018.
- 447 [54] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun. Unified perceptual parsing for scene understanding.
448 In *Proceedings of the European conference on computer vision (ECCV)*, pages 418–434, 2018.
- 449 [55] W. Xiong, J. Liu, I. Molybog, H. Zhang, P. Bhargava, R. Hou, L. Martin, R. Rungta, K. A.
450 Sankararaman, B. Oguz, et al. Effective long-context scaling of foundation models. *arXiv*
451 *preprint arXiv:2309.16039*, 2023.
- 452 [56] S. Yang, B. Wang, Y. Shen, R. Panda, and Y. Kim. Gated linear attention transformers with
453 hardware-efficient training. *arXiv preprint arXiv:2312.06635*, 2023.
- 454 [57] R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi. Hellaswag: Can a machine really
455 finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for*
456 *Computational Linguistics*, 2019.
- 457 [58] S. Zhai, W. Talbott, N. Srivastava, C. Huang, H. Goh, R. Zhang, and J. Susskind. An attention
458 free transformer. *arXiv preprint arXiv:2105.14103*, 2021.
- 459 [59] B. Zhang and R. Sennrich. Root mean square layer normalization. *Advances in Neural*
460 *Information Processing Systems*, 32, 2019.
- 461 [60] B. Zhou, H. Zhao, X. Puig, T. Xiao, S. Fidler, A. Barriuso, and A. Torralba. Semantic
462 understanding of scenes through the ADE20k dataset. *International Journal of Computer Vision*,
463 127:302–321, 2019.

464 A Scaling Up Number of Training Tokens

465 We scale up the number of training tokens to 350B for the 3B-size models. We compare with strong
 466 Transformer checkpoints including OpenLLaMA [18] and StableLM [50]. Moreover, we reproduce a
 467 Transformer language model (named Transformer_{Repro}) for apple-to-apple comparison.

468 Our model RetNet+ follows the same configuration as in Section 3.3, which is a hybrid model. The
 469 model’s hidden size is 3072, and the number of layers is 28. Without vocabulary embedding, the total
 470 number of parameters is 3.17B, which is between StableLM-3B-4E1T (2.7B) and OpenLLaMA-3B-
 471 v1 (3.19B). The batch size is 4M tokens. The training length is 4k. The learning rate is 3.2×10^{-4}
 472 with 1000 warm-up steps and linear learning rate decay. The training corpus includes The Pile [16]
 473 and RedPajama [9]. Transformer_{Repro} follows the exact same setting.

474 Table 5 reports accuracy numbers on the Harness-Eval benchmark [17]. We directly follow the evalua-
 475 tion protocol. The results show that RetNet+ achieves a performance comparable to Transformer_{Repro}
 476 on language tasks. Notice that OpenLLaMA-3B-v1 and StableLM-3B use different learning rate
 477 schedules. The results of these two models are used for reference purposes.

Model	ARC-C	ARC-C _{norm}	ARC-E	ARC-E _{norm}	Hellaswag	Hellaswag _{norm}
OpenLLaMA-3B-v1	0.303	0.323	0.641	0.599	0.449	0.608
StableLM-3B	—	—	0.649	0.610	—	—
Transformer _{Repro}	0.322	0.354	0.668	0.633	0.476	0.633
RetNet+	0.321	0.347	0.675	0.613	0.478	0.639
Model	OBQA	OBQA _{norm}	PIQA	PIQA _{norm}	Winogrande	Avg
OpenLLaMA-3B-v1	0.222	0.348	0.713	0.724	0.594	0.502
StableLM-3B	—	—	0.759	0.763	0.608	—
Transformer _{Repro}	0.258	0.358	0.746	0.755	0.612	0.529
RetNet+	0.258	0.362	0.750	0.763	0.614	0.529

Table 5: Accuracy on the Harness-Eval benchmark. All models are trained with 350B tokens with a batch size of 4M tokens. The results of OpenLLaMA-3B-v1 are taken from their official repository (<https://bit.ly/openllama-350b-results>), and StableLM-3B from their technical report (<https://bit.ly/StableLM-3B-4E1T>).

478 B Equivalence Between Chunk-wise Recurrent Representation and 479 Recurrent Representation

480 We illustrate the equivalence between the recurrent representation and the chunk-wise recurrent
 481 representation. Specifically, let B denote the chunk length. For the output O_n , n can be divided as
 482 $n = kB + r$ where B is the chunk size. Following Equation 6, we have:

$$\begin{aligned}
 O_n &= \sum_{m=1}^n \gamma^{n-m} Q_n K_m^T V_m \\
 &= (Q_n K_{kB+1:n}^T \odot \Gamma) V_{kB+1:n} + (Q_n \gamma^r) \sum_{c=0}^{k-1} \sum_{m=1}^B (K_{m+cB}^T V_{m+cB} \gamma^{B-m}) \gamma^{(k-1-c)B} \\
 &= (Q_n K_{kB+1:n}^T \odot \Gamma) V_{kB+1:n} + (Q_n \gamma^r) \sum_{c=1}^k (K_{[c]}^T (V_{[c]} \odot \zeta)) \gamma^{(k-c)B} \\
 &= (Q_n K_{kB+1:n}^T \odot \Gamma) V_{kB+1:n} + (Q_n \gamma^r) R_{i-1}
 \end{aligned} \tag{10}$$

483 where $\Gamma_i = \gamma^{n-i}$, $\zeta_{ij} = \gamma^{B-m}$, and $[i]$ indicates the i -th chunk, i.e., $x_{[i]} = [x_{(i-1)B+1}, \dots, x_{iB}]$.
 484 Then we write R_n as a recurrent function and compute the retention output of the i -th chunk via:

$$\begin{aligned}
 R_i &= K_{[i]}^T (V_{[i]} \odot \zeta) + \gamma^B R_{i-1} \\
 \zeta_{ij} &= \gamma^{B-i}, \quad \xi_{ij} = \gamma^i \\
 \text{Retention}(X_{[i]}) &= \underbrace{(Q_{[i]} K_{[i]}^T \odot D) V_{[i]}}_{\text{Inner-Chunk}} + \underbrace{(Q_{[i]} \odot \xi) R_{i-1}}_{\text{Cross-Chunk}}
 \end{aligned} \tag{11}$$

485 Finally, we show that the chunkwise recurrent representation is equivalent to the other representations.

486 C Results with Different Context Lengths

487 As shown in Table 6, we report the results of language modeling with different context lengths. In
 488 order to make the numbers comparable, we use 2048 text chunks as evaluation data and only compute
 489 the perplexity for the last 128 tokens. Experimental results show that RetNet performs comparably
 490 with Transformer in different context lengths.

Model	512	1024	2048
Transformer	13.55	12.56	12.35
RetNet	13.09	12.14	11.98

Table 6: Language modeling perplexity of RetNet and Transformer with different context length. The results show that RetNet has a consistent advantage across sequence length.

491 D Hyperparameters Used in Section 3.1

492 We use LLaMA [48] architecture, including RMSNorm [59] and SwiGLU [40, 7] module, as
 493 the Transformer backbone, which shows better performance and stability. The weights of word
 494 embedding and softmax projection are shared. Consequently, other variants follow these settings.
 495 For RetNet, the FFN intermediate dimension is $\frac{5}{3}d$ and the value dimensions in W_G, W_V, W_O are
 496 also $\frac{5}{3}d$, where the overall parameters are still $12d^2$.

497 For H3, we set the head dimension to 8. For RWKV, we use the TimeMix module to substitute
 498 self-attention layers while keeping FFN layers consistent with other models for fair comparisons.
 499 For Mamba, we follow all the details in the paper [19], where double-SSM layers are implemented
 500 instead of “SSM + SwiGLU”. In addition to RetNet and Mamba, the FFN intermediate dimension is
 501 all $\frac{8}{3}d$. All models have 400M parameters, 24 layers, and a hidden dimension of 1024. We train the
 502 models with 40k steps and a batch size of 0.25M tokens.

Params	Values
Layers	24
Hidden size	1024
Vocab size	100,288
Heads	24
Adam β	(0.9, 0.98)
LR	1.5×10^{-4}
Batch size	0.25M
Warmup steps	375
Weight decay	0.05
Dropout	0.0

Table 7: Hyperparamters used for the architecture comparison in Section 3.1.

503 E Hyperparameters Used in Section 3.2

504 We re-allocate the parameters in MSR and FFN for fair comparisons. Let d denote d_{model} for simplicity
 505 here. In Transformers, there are about $4d^2$ parameters in self-attention where $W_Q, W_K, W_V, W_O \in$
 506 $\mathbb{R}^{d \times d}$, and $8d^2$ parameters in FFN where the intermediate dimension is $4d$. In comparison, RetNet
 507 has $8d^2$ parameters in retention, where $W_Q, W_K \in \mathbb{R}^{d \times d}$, $W_G, W_V \in \mathbb{R}^{d \times 2d}$, $W_O \in \mathbb{R}^{2d \times d}$. Notice
 508 that the head dimension of V is twice Q, K , similar to GAU [26]. The widened dimension is
 509 projected back to d by W_O . In order to keep the parameter number the same as Transformer, the FFN
 510 intermediate dimension in RetNet is $2d$. Meanwhile, we set the head dimension to 256, i.e., 256 for

511 queries and keys, and 512 for values. For fair comparison, we keep γ identical among different model
 512 sizes, where $\gamma = 1 - e^{\text{linspace}(\log^{1/32}, \log^{1/512}, h)} \in \mathbb{R}^h$ instead of the default value in Equation (8).

Hyperparameters	1.3B	2.7B	6.7B
Layers	24	32	32
Hidden size	2048	2560	4096
FFN size	4096	5120	8192
Heads	8	10	16
Learning rate	6×10^{-4}	3×10^{-4}	3×10^{-4}
LR scheduler	Linear decay		
Warm-up steps	375		
Tokens per batch	4M		
Adam β	(0.9, 0.98)		
Training steps	25,000		
Gradient clipping	2.0		
Dropout	0.1		
Weight decay	0.05		

Table 8: Hyperparamters used for language modeling in Section 3.2.

513 F Results on Open-Ended Generation Tasks

514 Table 9 presents one-shot performance on two open-ended question-answering tasks, including
 515 SQUAD [39] and WebQS [5], with 6.7B models as follows. We report the recall metric in the table,
 516 i.e., whether the answers are contained in the generated response.

Dataset	SQUAD	WebQS
Transformer	67.7	36.4
RetNet	72.7	40.4

Table 9: Answer recall of RetNet and Transformer on open-ended question answering.

517 G Inference Cost of Grouped-Query Retention

518 We compare with grouped-query attention [1] and evaluate the method in the context of RetNet.
 519 Grouped-query attention makes a trade-off between performance and efficiency, which has been
 520 successfully verified in LLaMA2 34B/70B [49]. The method reduces the overhead of key/value cache
 521 during inference. Moreover, the performance of grouped-query attention is better than multi-query
 522 attention [42], overcoming the quality degradation brought by using one-head key value.

523 As shown in Table 10, we compare the inference cost with grouped-query attention and apply the
 524 method for RetNet. For the LLaMA2 70B model, the number of key/value heads is reduced by $8\times$,
 525 where the query head number is 64 while the key/value head number is 8. For RetNet-70B, the
 526 parameter allocation is identical to LLaMA [48], where the dimension is 8192, and the head number
 527 is 32 for RetNet. For RetNet-70B-GQ2, the key-value head number is 16, where grouped-query
 528 retention is applied. We run the inference with four A100 GPUs without quantization.

529 When the batch size is 256, LLaMA2 runs out of memory while RetNet without group query still
 530 has a high throughput. When equipped with grouped-query retention, RetNet-70B achieves 38%
 531 acceleration and saves 30% memory.

532 We evaluate LLaMA2 under 2k and 8k lengths separately. The batch size is reduced to 8 so that
 533 LLaMA2 can run without out of memory. Table 10 shows that the inference cost of Transformers
 534 increases with the sequence length. In contrast, RetNet is length-invariant. Moreover, RetNet-70B-
 535 GQ2 achieves better latency, throughput, and GPU memory than LLaMA2-70B-2k/8k equipped

536 with grouped-query attention. Notice that the evaluation metrics are averaged over positions of
 537 different sequence lengths for a fair comparison, rather than only considering the inference cost of
 538 the maximum length.

Model	Batch Size	Latency (ms)↓	Throughput (wps)↑	Memory (GB)↓
LLaMA2-70B-2k	256	—	—	OOM
LLaMA2-70B-8k	256	—	—	OOM
RetNet-70B	256	639.1	410.19	72.469
RetNet-70B-GQ2	256	461.8	567.66	52.726
LLaMA2-70B-2k	8	184.5	44.42	33.374
LLaMA2-70B-8k	8	277.7	29.50	37.386
RetNet-70B-GQ2	8	106.2	77.02	32.301

Table 10: Inference cost of RetNet and LLaMA2-70B with difference batch size and length. LLaMA2-70B is equipped with grouped-query attention, reducing key/value heads by $8\times$. “-GQ2” means grouped-query retention, which reduces half of key/value heads. “-2k” and “-8k” indicate sequence length for LLaMA2, while RetNet is length-invariant. RetNet is capable of large-batch inference and is favourable in terms of latency, throughput, and GPU memory.

539 **H Hyperparameters Used in Section 3.8**

Hyperparameters	DeiT	RetNet
Layers	12	12
Hidden size	512	512
Patch size	16	16
FFN size	2048	1024
Heads	8	2
Learning rate	1×10^{-3}	
LR scheduler	Cosine decay	
Batch size	1024	
Epochs	300	
Warmup epochs	5	
Smoothing	0.1	
Weight decay	0.05	
Drop path	0.3	

Table 11: Hyperparamters used for the ImageNet experiments in Section 3.8.

540 **NeurIPS Paper Checklist**

541 **1. Claims**

542 Question: Do the main claims made in the abstract and introduction accurately reflect the
543 paper's contributions and scope?

544 Answer: [\[Yes\]](#)

545 Justification: The abstract and introduction is carefully written.

546 Guidelines:

- 547 • The answer NA means that the abstract and introduction do not include the claims
548 made in the paper.
- 549 • The abstract and/or introduction should clearly state the claims made, including the
550 contributions made in the paper and important assumptions and limitations. A No or
551 NA answer to this question will not be perceived well by the reviewers.
- 552 • The claims made should match theoretical and experimental results, and reflect how
553 much the results can be expected to generalize to other settings.
- 554 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
555 are not attained by the paper.

556 **2. Limitations**

557 Question: Does the paper discuss the limitations of the work performed by the authors?

558 Answer: [\[Yes\]](#)

559 Justification: Limitations are discussed in the paper.

560 Guidelines:

- 561 • The answer NA means that the paper has no limitation while the answer No means that
562 the paper has limitations, but those are not discussed in the paper.
- 563 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 564 • The paper should point out any strong assumptions and how robust the results are to
565 violations of these assumptions (e.g., independence assumptions, noiseless settings,
566 model well-specification, asymptotic approximations only holding locally). The authors
567 should reflect on how these assumptions might be violated in practice and what the
568 implications would be.
- 569 • The authors should reflect on the scope of the claims made, e.g., if the approach was
570 only tested on a few datasets or with a few runs. In general, empirical results often
571 depend on implicit assumptions, which should be articulated.
- 572 • The authors should reflect on the factors that influence the performance of the approach.
573 For example, a facial recognition algorithm may perform poorly when image resolution
574 is low or images are taken in low lighting. Or a speech-to-text system might not be
575 used reliably to provide closed captions for online lectures because it fails to handle
576 technical jargon.
- 577 • The authors should discuss the computational efficiency of the proposed algorithms
578 and how they scale with dataset size.
- 579 • If applicable, the authors should discuss possible limitations of their approach to
580 address problems of privacy and fairness.
- 581 • While the authors might fear that complete honesty about limitations might be used by
582 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
583 limitations that aren't acknowledged in the paper. The authors should use their best
584 judgment and recognize that individual actions in favor of transparency play an impor-
585 tant role in developing norms that preserve the integrity of the community. Reviewers
586 will be specifically instructed to not penalize honesty concerning limitations.

587 **3. Theory Assumptions and Proofs**

588 Question: For each theoretical result, does the paper provide the full set of assumptions and
589 a complete (and correct) proof?

590 Answer: [\[NA\]](#)

591 Justification: There is no theoretical result in this paper.

592 Guidelines:

- 593 • The answer NA means that the paper does not include theoretical results.
- 594 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
595 referenced.
- 596 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 597 • The proofs can either appear in the main paper or the supplemental material, but if
598 they appear in the supplemental material, the authors are encouraged to provide a short
599 proof sketch to provide intuition.
- 600 • Inversely, any informal proof provided in the core of the paper should be complemented
601 by formal proofs provided in appendix or supplemental material.
- 602 • Theorems and Lemmas that the proof relies upon should be properly referenced.

603 4. Experimental Result Reproducibility

604 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
605 perimental results of the paper to the extent that it affects the main claims and/or conclusions
606 of the paper (regardless of whether the code and data are provided or not)?

607 Answer: [Yes]

608 Justification: The experiment can be easily reproduced based on the model description,
609 hyperparameter, and any well-known pre-training corpus.

610 Guidelines:

- 611 • The answer NA means that the paper does not include experiments.
- 612 • If the paper includes experiments, a No answer to this question will not be perceived
613 well by the reviewers: Making the paper reproducible is important, regardless of
614 whether the code and data are provided or not.
- 615 • If the contribution is a dataset and/or model, the authors should describe the steps taken
616 to make their results reproducible or verifiable.
- 617 • Depending on the contribution, reproducibility can be accomplished in various ways.
618 For example, if the contribution is a novel architecture, describing the architecture fully
619 might suffice, or if the contribution is a specific model and empirical evaluation, it may
620 be necessary to either make it possible for others to replicate the model with the same
621 dataset, or provide access to the model. In general, releasing code and data is often
622 one good way to accomplish this, but reproducibility can also be provided via detailed
623 instructions for how to replicate the results, access to a hosted model (e.g., in the case
624 of a large language model), releasing of a model checkpoint, or other means that are
625 appropriate to the research performed.
- 626 • While NeurIPS does not require releasing code, the conference does require all submis-
627 sions to provide some reasonable avenue for reproducibility, which may depend on the
628 nature of the contribution. For example
 - 629 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
630 to reproduce that algorithm.
 - 631 (b) If the contribution is primarily a new model architecture, the paper should describe
632 the architecture clearly and fully.
 - 633 (c) If the contribution is a new model (e.g., a large language model), then there should
634 either be a way to access this model for reproducing the results or a way to reproduce
635 the model (e.g., with an open-source dataset or instructions for how to construct
636 the dataset).
 - 637 (d) We recognize that reproducibility may be tricky in some cases, in which case
638 authors are welcome to describe the particular way they provide for reproducibility.
639 In the case of closed-source models, it may be that access to the model is limited in
640 some way (e.g., to registered users), but it should be possible for other researchers
641 to have some path to reproducing or verifying the results.

642 5. Open access to data and code

643 Question: Does the paper provide open access to the data and code, with sufficient instruc-
644 tions to faithfully reproduce the main experimental results, as described in supplemental
645 material?

646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697

Answer: [Yes]

Justification: Code will be released in camera-ready version. All of the data we use is public-available.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Hyperparameters are attached in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: For large language models, the variance between different runs is negligible. Moreover, the evaluation pipeline is deterministic.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 698
- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706
- 707
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
 - It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

708 8. Experiments Compute Resources

709 Question: For each experiment, does the paper provide sufficient information on the com-
710 puter resources (type of compute workers, memory, time of execution) needed to reproduce
711 the experiments?

712 Answer: [Yes]

713 Justification: The corresponding resources are stated in the paper.

714 Guidelines:

- 715
- 716
- 717
- 718
- 719
- 720
- 721
- 722
- The answer NA means that the paper does not include experiments.
 - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

723 9. Code Of Ethics

724 Question: Does the research conducted in the paper conform, in every respect, with the
725 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

726 Answer: [Yes]

727 Justification: We follow the NeurIPS Code of Ethics in the research.

728 Guidelines:

- 729
- 730
- 731
- 732
- 733
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

734 10. Broader Impacts

735 Question: Does the paper discuss both potential positive societal impacts and negative
736 societal impacts of the work performed?

737 Answer: [NA]

738 Justification: We work on fundamental research that has no direct societal impact.

739 Guidelines:

- 740
- 741
- 742
- 743
- 744
- 745
- 746
- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755
- 756
- 757
- 758
- 759
- 760
- 761
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

762 **11. Safeguards**

763 Question: Does the paper describe safeguards that have been put in place for responsible
764 release of data or models that have a high risk for misuse (e.g., pretrained language models,
765 image generators, or scraped datasets)?

766 Answer: [NA]

767 Justification: The paper does not pose safety risks.

768 Guidelines:

- 769
- 770
- 771
- 772
- 773
- 774
- 775
- 776
- 777
- 778
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

779 **12. Licenses for existing assets**

780 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
781 the paper, properly credited and are the license and terms of use explicitly mentioned and
782 properly respected?

783 Answer: [Yes]

784 Justification: We carefully follow the licenses of open-source code, data, and models.

785 Guidelines:

- 786
- 787
- 788
- 789
- 790
- 791
- 792
- 793
- 794
- 795
- 796
- 797
- 798
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

799 • If this information is not available online, the authors are encouraged to reach out to
800 the asset’s creators.

801 **13. New Assets**

802 Question: Are new assets introduced in the paper well documented and is the documentation
803 provided alongside the assets?

804 Answer: [NA]

805 Justification: The paper does not release new assets.

806 Guidelines:

- 807 • The answer NA means that the paper does not release new assets.
- 808 • Researchers should communicate the details of the dataset/code/model as part of their
809 submissions via structured templates. This includes details about training, license,
810 limitations, etc.
- 811 • The paper should discuss whether and how consent was obtained from people whose
812 asset is used.
- 813 • At submission time, remember to anonymize your assets (if applicable). You can either
814 create an anonymized URL or include an anonymized zip file.

815 **14. Crowdsourcing and Research with Human Subjects**

816 Question: For crowdsourcing experiments and research with human subjects, does the paper
817 include the full text of instructions given to participants and screenshots, if applicable, as
818 well as details about compensation (if any)?

819 Answer: [NA]

820 Justification: The paper does not involve crowdsourcing nor research with human subjects.

821 Guidelines:

- 822 • The answer NA means that the paper does not involve crowdsourcing nor research with
823 human subjects.
- 824 • Including this information in the supplemental material is fine, but if the main contribu-
825 tion of the paper involves human subjects, then as much detail as possible should be
826 included in the main paper.
- 827 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
828 or other labor should be paid at least the minimum wage in the country of the data
829 collector.

830 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
831 Subjects**

832 Question: Does the paper describe potential risks incurred by study participants, whether
833 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
834 approvals (or an equivalent approval/review based on the requirements of your country or
835 institution) were obtained?

836 Answer: [NA]

837 Justification: The paper does not involve crowdsourcing nor research with human subjects.

838 Guidelines:

- 839 • The answer NA means that the paper does not involve crowdsourcing nor research with
840 human subjects.
- 841 • Depending on the country in which research is conducted, IRB approval (or equivalent)
842 may be required for any human subjects research. If you obtained IRB approval, you
843 should clearly state this in the paper.
- 844 • We recognize that the procedures for this may vary significantly between institutions
845 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
846 guidelines for their institution.
- 847 • For initial submissions, do not include any information that would break anonymity (if
848 applicable), such as the institution conducting the review.