

AN EFFICIENT FRAMEWORK FOR LENGTH EXTENSION VIA DYNAMICALLY GROWING POSITIONAL EMBEDDING AND ROUTING ATTENTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Modeling long sequences is critical for numerous large-scale models. However, extending existing architectures to handle significantly longer sequences poses substantial technical and computational challenges. One inevitable issue is the overfitting of large models to positional encodings during pretraining, which limits their ability to generalize to unseen positional encoding scales. Additionally, extending sequence lengths requires extensive computational resources and time. Existing positional encoding methods often rely on carefully designed scaling factors but typically yield suboptimal results. To tackle these challenges, we propose **Cyclic, Randomly Truncated, and Dynamically Growing NTK Positional Embedding (CRG NTK)**, a data-augmentation-based technique that fully explores the RoPE encoding space, enabling models to adapt to various positional scales and achieve state-of-the-art extrapolation for the extension of lengths dominated by position encoding. Furthermore, we introduce **an efficient attention mechanism with a correlation-based routing strategy to enhance the fitting of the augmented positional encoding**, yielding superior performance and more efficient fine-tuning. With our approach, LLaMA-7B and Mistral-7B fine-tuned at 16K context length achieve extrapolation factors of at least $128\times$ on simple tasks and maintain stable perplexity over $32\times$ sequence length extensions and saves at least 16 times the GPU training resources compared to the existing optimal method. Experiments also show that correlation routing can achieve good performance by further filtering out large amounts of noise in long sequences.

1 INTRODUCTION

In various natural language processing (NLP) tasks, such as document-level sentiment analysis Behdenna et al. (2018), long document summarization Koh et al. (2022), and code generation Rozière et al. (2024), the ability to effectively model long-sequence dependencies is crucial. This capability allows for capturing complex relationships over sequences spanning hundreds or thousands of tokens, which is essential for tasks where contextual information is dispersed. Consequently, extending the context window enables large language models (LLMs) to perform tasks that shorter context windows cannot handle and potentially enhances performance across a variety of NLP tasks.

However, extending the context window of LLMs Touvron et al. (2023a); Jin et al. (2023) is a critical research direction, yet it faces numerous challenges. Foremost among these is the influence of positional encoding mechanisms on the context window extension. During pre-training, models typically operate within fixed-length windows using predefined positional encodings. This often leads to overfitting to the positional distribution within the training range, thereby limiting the model’s ability to handle unseen positional information beyond this range. Furthermore, the resource requirements for extending the context window are significant, involving not only substantial storage space but also excessive computation time due to the need for fine-tuning on extremely long sequences.

To address the limitations of positional encodings, recent work Press et al. (2022); Chi et al. (2022); Li et al. (2024); Zheng et al. (2024) has extensively analyzed relative mechanisms, particularly Rotary Positional Embedding (RoPE) Su et al. (2023). A common approach to extending RoPE involves scaling positional indices, effectively mapping out-of-distribution positions back into the pre-trained

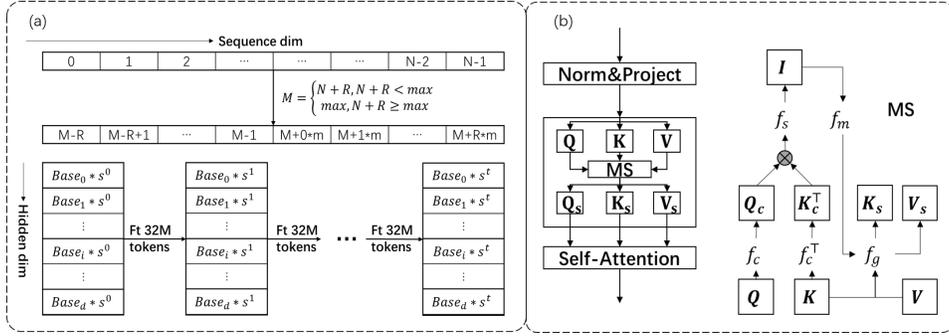


Figure 1: (a) During fine-tuning, the CRG NTK (Cyclic, Randomly Truncated, and Dynamically Growing NTK Positional Embedding) mechanism is employed for context length extension. N : Original sequence length. R and $m = 1$: Randomly sampled starting index for cyclic truncation. s^k : scale for Dynamically Growing. "Ft 32M tokens": fine-tuning with 32 million tokens. (b) Efficient fine-tuning architecture using Merge and Select Attention (MS Attention) integrated with LoRA. I : Indexes for selection. f_s, f_m : Select and merging functions. f_c, f_c^T, f_g : Compress and gather functions in MS Attention.

range. For example, Positional Interpolation (PI) Chen et al. (2023) rescales indices to fit within the original context, while NTK-aware Ntk (2023) methods scale the base frequency to modestly enhance extrapolation. YaRN Peng et al. (2023) improves on this by separating frequencies into high, mid, and low bands, applying different scaling strategies: leaving high frequencies unchanged, scaling low frequencies with PI, and adjusting mid frequencies using NTK-aware methods. LongRoPE Ding et al. (2024) further optimizes scaling ratios and anchor positions through a search algorithm. Despite these efforts, most methods rely on carefully crafted scaling heuristics and still fall short in capacity for extrapolation and performance. Moreover, context extension also faces hardware bottlenecks: large model sizes demand significant compute and memory. To mitigate this, LongLoRA Chen et al. (2024) integrates efficient Attention, DeepSpeed Rasley et al. (2020), and LoRA Hu et al. (2021) to reduce fine-tuning time and storage requirements.

To overcome these limitations, we propose a data-augmentation-based positional encoding method designed to improve the generalization ability of positional encodings. Our approach employs dynamically growing scaling ratios, allowing positional encodings to scan a broader solution space to achieve stronger extrapolation capability. Based on the above method, we provide a degree of theoretical interpretation and empirical support for the extrapolation capacity of models: **the extrapolation limit is approximately bounded by the model depth**. This extrapolation-based context extension strategy not only significantly reduces the resource requirements for expanding the model's context length, but also substantially advances the study of foundation models, such as understanding the transmission of information across layers.

We further optimize the LoRA configuration and design an efficient attention mechanism aligned with positional encoding strategies. Our relevance-based routing selectively weights value tokens based on relevance-based routing mechanism, reducing computation while enhancing the model's dependence on accurate positional cues. This allows effective regression even with partial positional information. We argue that simply increasing sequence length does not unconditionally improve model performance. As the sequence length grows, the beneficial information gain diminishes, while the negative impact of noise becomes increasingly pronounced. Although Softmax provides limited noise suppression, it is inadequate for long sequences. In contrast, our routing mechanism explicitly filters irrelevant tokens, enabling efficient attention under controlled noise and improving long-sequence performance.

In summary, our method tackles the core challenges of context window extension by introducing advances in both positional encoding and efficient Attention. Through theoretical analysis and empirical validation, we demonstrate superior extrapolation capabilities and performance. Leveraging CRG NTK positional encoding, we achieve extrapolation up to $32\times$ the fine-tuning length on a 32-layer model—approaching the hypothesized limit. Then, combining with MS Attention our method surpasses LongRoPE with at least $16\times$ less GPU fine-tuning time and comparable performance. Furthermore, our approach achieves over $128\times$ extrapolation on simple retrieval tasks. These results offer new insights and practical strategies for efficiently scaling context windows in LLMs.

2 PRELIMINARY

2.1 TRANSFORMER

The Llama2 Touvron et al. (2023b) and Mistral Jiang et al. (2023) models are based on the Transformer architecture, which consists of the core modules self-attention and feed-forward network (FFN). The computation process of self-attention Vaswani et al. (2017) is as follows: $O = \text{softmax}(QK^\top)V$, where Q , K , and V are obtained from $X \in \mathbb{R}^{n \times d}$ using embedding weights W_q , W_k , and W_v respectively. The final output O is then passed through W_o to obtain the final output of the attention. Subsequently, the entire Transformer layer is completed through the FFN.

2.1.1 FIXED-PATTERN AND DYNAMICALLY ROUTED SPARSE ATTENTION

This category of methods restricts queries (Q) to fixed or sequence-independent key (K) positions, which significantly reduces the computational cost. The general formulation is:

$$\tilde{K} = \text{Loc}(K, i_k(Q, K)), \quad \tilde{V} = \text{Loc}(V, i_k(Q, K)), \quad \text{output} = \text{Softmax}(Q\tilde{K})\tilde{V}, \quad (1)$$

where $\text{Loc}(\cdot)$ represents a selection function, and i_k denotes the selected indices. In fixed-pattern approaches, i_k is independent of Q and K . s is sliding window size and l is shift size:

- Sliding Window: $i_k \in (i - \frac{s}{2}, i + \frac{s}{2})$.
- Global Tokens: $i_k \in (0, s)$, commonly used in attention sink mechanisms.
- Swin Transformer Liu et al. (2021): alternates between windowed attention $i_k \in (i - \frac{s}{2}, i + \frac{s}{2})$ and shifted window attention $i_k \in ((i - \frac{s}{2} + l) \bmod n, (i + \frac{s}{2} + l) \bmod n)$.
- BigBird Zaheer et al. (2021): incorporates sliding window, global tokens, and randomly selected keys and values: $i_k \in (r_1, r_2, \dots, r_m)$.

Unlike fixed-pattern approaches, dynamically routed methods determine key (K) positions based on the importance of different regions, thereby achieving improved performance. In other words, the selected indices $i_k(Q, K)$ depend on Q and K .

For instance, Biformer Zhu et al. (2023) selects keys dynamically using: $i_k(Q, K) = \arg \max_{\text{top-}K}(\bar{Q}\bar{K}^\top)$, where $\bar{Q} \in \mathbb{R}^{n/s \times d}$ and $\bar{K} \in \mathbb{R}^{n/s \times d}$ are compressed segment representations of Q and K , respectively. More detailed analysis we put in the appendix E.1.

2.2 POSITIONAL ENCODING

Positional encoding can be divided into relative and absolute positional encodings. The most widely used method for relative positional encoding is the Rotary Positional Encoding (RoPE) Su et al. (2023). The encoding formula is given as follows, where $\theta_i = \text{base}^{-\frac{2i}{d}}$ and base is the fundamental frequency, usually taken as 10000 or some other integer value:

$$\mathbf{p}_c(m, \text{base}) = [\cos(m\theta_0) \quad \cos(m\theta_0) \quad \dots \quad \cos(m\theta_{\frac{d}{2}-1}) \quad \cos(m\theta_{\frac{d}{2}-1})] \quad (2)$$

$$\mathbf{p}_s(m, \text{base}) = [\sin(m\theta_0) \quad -\sin(m\theta_0) \quad \dots \quad \sin(m\theta_{\frac{d}{2}-1}) \quad -\sin(m\theta_{\frac{d}{2}-1})] \quad (3)$$

The output for the encoding of the m -th token is given by:

$$f_p(\mathbf{x}_m, \mathbf{p}(m)) = [x_0 \quad x_1 \quad \dots \quad x_{d-1} \quad x_d] \odot \mathbf{p}_c(m) + [x_1 \quad x_0 \quad \dots \quad x_d \quad x_{d-1}] \odot \mathbf{p}_s(m). \quad (4)$$

Positional Interpolation Chen et al. (2023) is performed by scaling the above position, where the index of the m -th position is transformed as follows. Where λ is the scaling factor, and f represents the vector of cosine or sine in RoPE.

$$\mathbf{p}_f^{(PI)}(m) = [f(\frac{m}{\lambda}\theta_0) \quad f(\frac{m}{\lambda}\theta_0) \quad \dots \quad f(\frac{m}{\lambda}\theta_{\frac{d}{2}-1}) \quad f(\frac{m}{\lambda}\theta_{\frac{d}{2}-1})]$$

NTK Positional Encoding Ntk (2023): For NTK positional encoding, the base is rescaled as:

$$\text{base}' = \text{base} \times \left(\lambda^{\frac{d}{d-2}} \right), \theta'_i = \text{base}'^{-\frac{2i}{d}}, \mathbf{p}_f^{(NTK)}(m) = [f(m\theta'_0), f(m\theta'_0), \dots, f(m\theta'_{\frac{d}{2}-1}), f(m\theta'_{\frac{d}{2}-1})]$$

162 **YaRN** Ntk (2023): Other positional encoding methods extend from the two encoding strategies
 163 described above. For example, YaRN Peng et al. (2023) posits that high-frequency information is
 164 more important, so it retains the high-frequency portion of the positional encoding unchanged, while
 165 applying interpolation to the low-frequency components.

$$166 \quad r(i) = \frac{L}{\theta_i} = \frac{L}{2\pi b' \frac{2i}{d}}, h_y(\theta_i) = \begin{cases} \frac{\theta_i}{\lambda}, & \text{if } r(i) < \alpha, \\ \theta_i, & \text{if } r(i) > \beta, \\ \frac{\beta - r(i)}{\beta - \alpha} \frac{\theta_i}{\lambda} + \frac{r(i) - \alpha}{\beta - \alpha} \theta_i, & \text{otherwise.} \end{cases} \quad (5)$$

$$170 \quad \mathbf{p}_f^{(\text{YaRN})}(m) = [f(mh_y(\theta_0)) \quad f(mh_y(\theta_0)) \quad \dots \quad f(mh_y(\theta_{\frac{d}{2}-1})) \quad f(mh_y(\theta_{\frac{d}{2}-1}))] \quad (6)$$

172 where α and β are designated by the author as 1 and 32, $s = \frac{L'}{L}$, L is original context size, L' is
 173 current length.

174 **LongRoPE**: LongRoPE Ding et al. (2024) addresses the issue of attention sink Xiao et al. (2024)
 175 by preserving the positional encoding up to a specific position \hat{n} , and performing interpolation for
 176 positions beyond \hat{n} . The interpolation ratio λ_i is determined via a search algorithm.

$$178 \quad h_l(m, \theta_i) = \begin{cases} \frac{\theta_i}{\lambda_i}, & \text{if } m \geq \hat{n}, \\ \theta_i, & \text{if } r(i) < \beta, \end{cases} \quad (7)$$

$$180 \quad \mathbf{p}_f^{(\text{LRoPE})}(m) = [f(mh_l(m, \theta_i)) \quad f(mh_l(m, \theta_i)) \quad \dots \quad f(mh_l(m, \theta_{\frac{d}{2}-1})) \quad f(mh_l(m, \theta_{\frac{d}{2}-1}))]. \quad (8)$$

183 3 METHODS

184 3.1 CYCLIC, RANDOMLY TRUNCATED, AND DYNAMICALLY GROWING NTK METHOD

185 Due to the pretraining of models on fixed-size windows, they tend to overfit to positional information
 186 within this range. To better handle unseen positional information, one solution is to fine-tune the
 187 model using larger context windows, thereby improving its generalization ability. Another approach
 188 is to scale the positional information during the fine-tuning, bringing it back into the range of the
 189 pretrained window size. This method can greatly enhance the positional encoding capability learned
 190 during pretraining. Furthermore, combining these two strategies is the most common way—using
 191 larger context windows and scaling the positional encoding during fine-tuning, enables the extension
 192 of even longer context windows.

193 This paper briefly analyzes the reasons for the above approach; please refer to the appendix D for
 194 more details. We break down the attention score calculation into the attention scores of sequences
 195 without positional information and the adjustments made by the relative receptive field change
 196 through positional information. To make the analysis more concise and easily scalable, we linearize
 197 the positional information of the Attention formula as follows:

$$200 \quad e^{\mathbf{q}_n \mathbf{k}_m^\top} (1 + f_p(\mathbf{q}_n, \mathbf{p}(n)) f_p(\mathbf{k}_m, \mathbf{p}(m)))^\top - \mathbf{q}_n \mathbf{k}_m^\top$$

202 The error of the above formulae with respect to the original Attention is $O\left(\frac{n}{\text{base}^{\frac{2i}{d}}}\right)$, and it has been
 203 shown experimentally that its extension for position encoding is identical to the original Attention.
 204 With the decomposition of the above formulae, the position information is concentrated only in the
 205 linear terms $\mathbf{p}(n) f_p(\mathbf{k}_m, \mathbf{p}(m))^\top$. A simple analysis of the forward process shows that the equality
 206 of the attention scores for the unfine-tuned extrapolation-only and the learned extrapolation length
 207 through learning indicates that successful extrapolation can be achieved, and this equation can be
 208 simply deduced to be similar to that of equation 9. This approach yields only one particular solution,
 209 and we perform a more general and rigorous derivation and analysis through the lens of optimisation
 210 in the appendix D.2. In a word, the goal is to make the following equation hold true, L is desired
 211 length:

$$212 \quad \frac{(L-1)/s_1}{(\text{base}/s_2)^{\frac{i}{2d}}} \in \left[\frac{(0)}{\text{base}^{\frac{i}{2d}}}, \frac{(M-1)}{\text{base}^{\frac{i}{2d}}} \right] \quad (9)$$

214 where M represents the sequence length during training or fine-tuning, s_1 and s_2 are scalar *position*
 215 and *frequency* scaling factors respectively (choices of s_1, s_2 are discussed below), $i \in \{0, \dots, d/2 - 1\}$ indexes the RoPE frequency.

Search over base frequencies and adjusted bases. We assume there exists a family of base frequencies $\{\text{base}_i\}$ such that, by fine-tuning the model on sequences constructed with these bases, the model can restore performance at the desired length L_t . Concretely, we consider *adjusted* bases $\text{base}_i^{(k)} = \frac{\text{base}_i}{s_i^{(k)}}$, where $s_i^{(k)}$ are per-frequency adjustment factors (either hand-designed, searched, or learned). Prior works (e.g. NTK, YaRN, LongRoPE) and appendix D.2 provide strong empirical and theoretical motivation for such frequency-wise adjustments; this motivates scanning the base-space to find $\{s_i^{(k)}\}$ that recover the pretrained behaviour.

Solution-space scanning via a power-law growth schedule. To efficiently explore the space of feasible base adjustments, we adopt a *power-law dynamic growth ratio* schedule for the scale parameter(s). The motivation is two-fold: Modern LMs exhibit some inherent extrapolation ability: when positional encodings are smoothly scaled, the model often preserves much of its short-range behaviour at longer contexts. Thus a smoothly varying family of scales can be learned incrementally. Under a power-law schedule each adjacent scale differs only slightly, so the positional information learned at one scale largely transfers to the next. Hence scanning a small number of scales with exponential spacing suffices to cover the continuous solution band (see Figure 2).

Practically, we enumerate a small set of scale factors $\{s^{(k)}\}$ with exponential spacing (e.g. $s^{(k)} = s^{(0)} \cdot a^k$ for some $a > 1$), fine-tune at each scale for a short time, and take the union of the covered phase ranges to ensure the target L_t is represented.

Mitigating catastrophic forgetting. The scanning process is designed to be *incremental* and *inheritable*: the positional representations learned at one scale serve as warm-starts for adjacent scales, which greatly reduces the amount of new information to be learned and mitigates catastrophic forgetting. Empirically, the difference in positional embeddings between neighboring scales is small (see Figure 3), so the optimization pressure per step is modest and convergence is fast.

Random shifts, cyclic phases and mitigation of Attention Sink. To further improve generalization, we introduce two practical regularizers: **Random shifts**: before applying the positional map, we add a random offset r (sampled per-batch) to the position indices. This is analogous to data augmentation: it forces the model to be robust to local phase offsets. **Cyclic modular mapping**: we optionally map positions using a modulo operator to a finite cycle length max . Concretely, we replace a raw position m with $(m + r) \bmod \text{max}$. When max is sufficiently large and shares factors with the RoPE frequencies, the positional encoding exhibits cyclic phases; this enables the model to *simulate* very long (even effectively infinite) relative positions by folding them into the learned phase space.

The final positional encoding formula is as follows:

$$h_c(m, i, s^{(k)}) = ((m + r) \% \text{max}) \left(\text{base}_i \times \left((w_i \times s^{(k)})^{\frac{d}{d-2}} \right) \right)^{-\frac{2i}{d}}$$

$$\mathbf{p}_c = [\cos(h_c(m, 0, s^{(k)})) \quad \cos(h_c(m, 0, s^{(k)})) \quad \dots \quad \cos(h_c(m, \frac{d}{2} - 1, s^{(k)})) \quad \cos(h_c(m, \frac{d}{2} - 1, s^{(k)}))] \quad (10)$$

$$\mathbf{p}_s = [\sin(h_c(m, 0, s^{(k)})) \quad -\sin(h_c(m, 0, s^{(k)})) \quad \dots \quad \sin(h_c(m, \frac{d}{2} - 1, s^{(k)})) \quad -\sin(h_c(m, \frac{d}{2} - 1, s^{(k)}))] \quad (11)$$

where scale factor $s^{(k)} = s^{(0)} \cdot a^k$ that after a certain number of fine-tuning steps multiply by a , r is randomly generated, max is set manually, typically to twice the target extension length and w_i is learnable parameter. By fine-tuning the LLM using the above methodology, we were able to achieve extrapolation capabilities close to a multiple of the model layer size, i.e., using $16K$ length and Full Attention fine-tuning, and were able to maintain a stable level of perplexity when tested using Full Attention in $16K \times l$ length contexts.

Despite comprehensive solution space coverage, our approach faces two primary challenges: (1) As scaling ratios increase dynamically, positional encodings may imperfectly align with the inherent locality bias of sequence modeling tasks, where current tokens preferentially attend to nearby contexts. Counterintuitively, we observe performance improvements with larger scaling ratios, suggesting that expanded contextual windows introduce beneficial information for regression tasks after extensive training. In order to further improve the scaled performance, we consider in the appendix the inclusion of YaRN-like segmentation to make it more consistent with the sequence modelling bias. (2) At

270 extreme sequence lengths, full attention computation becomes dominated by irrelevant tokens, with
 271 noise overwhelming useful signal.

272 To address these limitations, we integrate an efficient attention mechanism with relevance-based
 273 routing. Our key insight is that prediction-relevant tokens constitute a sparse subset of the complete
 274 sequence. By employing fixed-length windows selected via relevance-based filtering, we maintain
 275 near-optimal performance while achieving significant computational efficiency. This approach
 276 effectively mitigates noise accumulation while preserving essential contextual information.
 277

278 3.2 MERGE SELECTION

279 We propose a method for implementing more general and efficient sparse attention through correlation
 280 selection and merging mechanisms, as shown in Figure 4. This method consists of two main steps:
 281 selection of relevant regions and merging of these regions. The proposed method consists of two
 282 main steps: selection and merging. The pseudo-code is 2

283
 284 **Step 1: Selection.** First, the Q , K , and V tensors with shape (b, h, n, d) are segmented into regions,
 285 resulting in tensors Q_s with shape (b, h, n_{sq}, s_q, d) , K_s , and V_s with shape (b, h, n_{sk}, s_k, d) , where
 286 s_q and s_k denote the segment size. Each region is represented by a semantic token or an average
 287 compressed token, yielding Q'_s and K'_s . A dot product or another similarity metric is then applied
 288 between Q'_s and K'_s to analyze the relevance between Q_s and K_s regions. To prevent information
 289 leakage, we apply a mask and control the number of selected tokens. This process results in the
 290 indices of the top- k most relevant K_s regions for each Q_s region, denoted as $selectindex$ with shape
 291 $(b, h, n_{sq}, topk)$, where n_s is the number of Q_s regions and $topk$ is the number of selected K_s and
 292 V_s regions.

293 **Step 2: Merging.** In the merging step, the indices obtained from the first step and Q_s are per-
 294 muted and combined, effectively merging the selected regions. Q_s is merged according to a spec-
 295 ified number of segments $m = merges$, resulting in Q_{ms} with shape $(b, h, n_{ms}, m \cdot s, d)$, where
 296 $n_{ms} = n_s/m$. Similarly, $selectindex$ is split, permuted, and merged into $mselectindex$ with shape
 297 $(b, h, n_s/m, m, top - k) \rightarrow (b, h, n_{ms}, top - k, m) \rightarrow (b, h, n_{ms}, top - k \cdot m)$. Due to previous
 298 top- k operation, the permutation can ensure that each row is sorted by relevance, with the first m
 299 indices corresponding to the most relevant K_s region for each of the merged Q_s regions, the next
 300 m indices corresponding to the second most relevant K_s region, and so on. After obtaining the
 301 merged indices $mselectindex$, performs unique operation while maintaining relevance order. The
 302 top- n indices, denoted as $qmselectindex$, are selected as the final indices corresponding to the merged
 303 Q regions and their relevant K_s regions.

304 More detailed analysis is given in the Appendix E.1. Merging ensures sufficient K/V token coverage
 305 per Q , where larger Q regions sharing K/V tokens improve efficiency over smaller isolated regions.
 306 This strategy also enhances selection precision through localized relevance filtering.

307 While inspired by Biformer Zhu et al. (2023) and Routing Transformer Roy et al. (2021), our method
 308 uniquely enables: 1. **Dynamic region scaling:** Flexible $Q/K/V$ partitioning (Appendix E.1) adapts
 309 to varying attention ranges and token retention needs. 2. **Compressible merging:** Adjustable
 310 merged region sizes balance computational efficiency and contextual breadth. 3. **Layer-aware**
 311 **extrapolation:** Combined with CRG NTK-enhanced positional encoding, our attention mechanism
 312 selectively regresses positional subsets rather than full sequences, amplifying extrapolation limits.

313 We hypothesize these limits scale linearly with layer count—e.g., a 32-layer Llama2-7B model
 314 achieves 32 \times extrapolation beyond its training sequence length. This aligns with Swin Transformer’s
 315 shifted-window paradigm, where layer stacking expands the effective receptive field multiplicatively
 316 (current layer tokens attending to prior layer counterparts). See the appendix D.10 for a brief
 317 theoretical discussion, and empirical also validation supports this hypothesis.

318 4 EXPERIMENT

319 4.1 EXPERIMENTAL SETUP

320
 321 Our experiments were conducted on only one A100 80GB GPUs, using the Llama2-7B Touvron et al.
 322 (2023c) model and Mistral-7B-v0.1 with the attention mechanism replaced by our MS Attention as
 323

described in Section 3.2. The training approach used the efficient fine-tuning method described in Section E.3, with CRG NTK position interpolation applied.

Fine-tuning Steps and Parameters: The fine-tuning parameters included the linear transformations for key and query (W_k and W_q), as well as the embedding and normalization parameters. W_k and W_q take LoRA for fine-tuning with rank 16. The training approach used an autoregressive method, where the objective was to generate the next token. The loss function used was cross-entropy, and the optimizer was AdamW Loshchilov & Hutter (2019) with parameters $\beta_1 = 0.9$, $\beta_2 = 0.95$ and $lr = 1e - 4$.

Dataset: The dataset used for fine-tuning was the Redpajama Computer (2023) dataset, which is the same dataset used in LongLora. Evaluation was performed using the widely used long text dataset PG19 Rae et al. (2019) and proofpile. Additionally, the model’s performance was assessed using the passkey task Mohtashami & Jaggi (2023) and LongBench Bai et al. (2024).

4.2 MAIN RESULTS

Table 1: Passkey Retrieval Task Performance. FT=16K with FA indicates fine-tuning using Full Attention, others use MS Attention.

Model Configuration	2K	16K	32K	64K	128K	256K	512K	1024K	2048K	4096K
Baselines										
Llama2-7B-LongLora (ft=100K, PI)	1.0	1.0	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
Llama2-7B-LongRoPE (ft=256K, PI)	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.6	0.0
Our Method (Llama2-7B-MS)										
FT=16K, PI, scale=16	1.0	1.0	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
FT=16K, NTK, scale=4096	1.0	1.0	1.0	1.0	1.0	0.7	0.0	0.0	0.0	0.0
FT=16K, CR-NTK, scale=4096	1.0	1.0	1.0	1.0	1.0	1.0	0.0	0.0	0.0	0.0
FT=16K, CR-NTK, scale=16384	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.0	0.0	0.0
FT=16K, CRG-NTK, scale=4096 × 2^k	1.0									
FT=16K with FA, CRG-NTK, scale=4096 × 2^k	1.0									
Cross-Architecture Validation										
Mistral-7B-MS (FT=16K, CRG-NTK, scale=4096 × 2 ^k)	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0

Table 2: Perplexity Results on PG19 Dataset. Models were fine-tuned with Multi-Scale Attention and evaluated with either MS Attention or Full Attention. Lower values indicate better performance.

Model Configuration	16K	32K	64K	128K	256K	512K	1024K	2048K
Baselines								
Llama2-7B-LongLora (FT=32K, PI)	7.35	7.22	-	-	-	-	-	-
Llama2-7B-LongRoPE (FT=256K)	7.37	-	6.64	6.31	-	-	-	-
Our Method (Llama2-7B-MS)								
FT=16K, PI, scale=16	6.50	6.25	6.53	-	-	-	-	-
FT=16K, CR-NTK	7.86	7.92	7.97	8.23	8.68	-	-	-
FT=16K, CRG-NTK (test: Full Attention)	7.20	7.25	7.29	7.43	8.36	9.72	16.98	36.82
FT=16K, CRG-NTK (test: MS Attention)	6.92	6.85	7.09	7.31	7.56	7.78	7.96	8.32
Extended Fine-tuning (Llama2-7B-MS)								
FT=64K, CRG-NTK (test: Full Attention)	7.24	6.45	6.19	6.23	6.47	6.76	7.38	7.92
FT=64K, CRG-NTK (test: MS Attention)	6.62	6.60	6.09	6.22	6.56	6.88	6.89	6.92
Cross-Architecture Validation								
Mistral-7B-MS (FT=16K, CRG-NTK, test: FA)	8.20	8.15	8.29	8.56	9.76	12.12	18.96	44.82

Table 3: Extrapolation Perplexity for Models with Different Number of Layers.

Model	layers	16K	32K	64K	128K	256K	448K	512K	640K
Llama3.2-1B-MS (ft=16K, CRG NTK, test:FA)	16	13.10	12.86	13.19	13.92	$\boxed{15.77}_{\times 16}$	52.37	64.52	94.78
Llama3.2-3B-MS (ft=16K, CRG NTK, test:FA)	28	10.24	9.06	9.39	9.56	9.89	$\boxed{11.64}_{\times 28}$	20.14	62.96
Llama2-7B-MS (ft=16K, CRG NTK, test:FA)	32	7.10	7.07	7.16	7.40	8.26	-	$\boxed{9.96}_{\times 32}$	14.56
Llama2-13B-MS (ft=16K, CRG NTK, test:FA)	40	6.55	6.32	6.38	6.44	6.56	-	6.98	$\boxed{8.97}_{\times 40}$

4.2.1 LENGTH EXTENSION: FINE-TUNING

We use the CRG NTK and MS attention to fine-tune the LLM and enhance its generalization of positional information. First, for the architecture, during fine-tuning, we replace the attention mechanisms in Llama and Mistral with our MS Attention that segment size is 16, merge count is 128

378 and selection number is 512. For parameter fine-tuning, since only the Q and K tokens incorporate
 379 positional information, we use LoRA to fine-tune the mapping weights of W_Q and W_K with rank 16,
 380 while other parameters are only fine-tuned for embeddings and normalization layers. Regarding the
 381 choice of positional encoding, we apply CRG NTK method, starting with an initial scaling factor (set
 382 to scale = 4096 in this work), and increase this factor fourfold after fine-tuning a small fixed data
 383 volume (32M tokens in our case).

384 1. For the passkey task with a length of 16K for fine-tuning, as shown in Table 1, using standard NTK
 385 positional encoding combined with our algorithm enables length extensions of up to 16x. Moreover,
 386 employing random cropping and circular modular (CR NTK) sampling stabilizes this extension.
 387 Finally, by incorporating dynamic growth (CRG NTK), we can scale up to 4M and possibly arbitrary
 388 length. For every 4x increase in base scale, our approach can solve passkey tasks with a 2x length
 389 extension, potentially due to the square root relationship between positional encoding and the base
 390 value.

391 A potential issue with this approach is that the positional encoding ratio during inference remains
 392 fixed. For instance, when training with a length of 16K, the NTK positional encoding starts with
 393 scale = 4096 and the maximum position_ids = 32K. During inference across all passkey task lengths
 394 (0-256K), we use a fixed base scaling rate of 1024 for prediction, achieving 100% performance.
 395 After extending by 16x, the scaling factor becomes scale = 4096×16 , and inference using a fixed
 396 base scaling rate of 4096 can solve passkey tasks up to 1M in length. Similarly, extending by the
 397 corresponding multiples and using a fixed base scaling rate of 65536 extends passkey tasks up to 4M.

398 2. For the long-text perplexity experiments, we evaluated our model on the PG19 and proofpile
 399 datasets, as presented in Table 2 and Table 12. Based on the experimental results, positional
 400 interpolation generally yields better performance, as it effectively preserves the learned distributions.
 401 In contrast, the NTK-based method may disrupt this distribution to some extent, leading to a
 402 performance loss, but it enables much longer extrapolation.

403 From the results, with 16K fine-tuning, we tested using Full Attention, and the performance on a
 404 512K length test was close to the original perplexity, achieving an extrapolation of nearly 32 times,
 405 which corresponds to the number of layers in Llama7B. Moreover, in Table 3 we show the results of
 406 the perplexity test for models with different number of layers, which shows that the extrapolation
 407 ability is very much related to the number of layers of the model. Moreover, we believe that after
 408 reaching a certain length, the gain from introducing useful information becomes much smaller than
 409 the effect of noise in the sequence. Therefore, we use efficient attention with correlation routing to
 410 further enhance the filtering capability, thus maintaining or even improving performance, as show in
 411 Table 2: test with MS.

412 4.2.2 IMPACT OF SELECTION AND MERGE

413 Further validation shows that our efficient Attention enhances the filtering capability and improves
 414 performance. As shown in Table 10, with 16K fine-tuning, we set different patch selection precision,
 415 selection quantity (select), and merging quantity (merge). We found that using better precision,
 416 selecting fewer items, and merging fewer items leads to better performance. However, this comes at
 417 the cost of reduced operational efficiency and decreased generalization. For example, with patch=8
 418 and select=32, we achieved a perplexity of 6.09 on PG19, which is far better than other optimal
 419 methods, such as LongLora, which achieved a perplexity of 6.86 with 16K fine-tuning.

420 4.3 MEMORY USAGE AND FINE-TUNING TIME

421 Table 4 compares the fine-tuning time of our method with LongLoRA and LongRoPE. Our approach
 422 significantly reduces both the memory overhead and fine-tuning time. Due to the strong extrapolation
 423 capability of our method, it requires only short-length fine-tuning to achieve the same context window
 424 length as existing methods. As shown in the table, we fine-tune with a sequence length of 50K for
 425 only 29 hours, yet successfully generalize to a 1M-length context window. For LLaMA2-7B with 32
 426 layers, the theoretical extrapolation limit implies that a training length of 32K is sufficient. However,
 427 since performance degradation may occur near the extrapolation boundary, we include a small margin
 428 by slightly increasing the training length. For LongRoPE its expansion to 1M length requires 8
 429 A100s for 5 days of search (data from its paper Ding et al. (2024)) and needs to be fine-tuned on 128K
 430 length, so its total time is $5 \times 24h \times 8 + 98h$, where 98h comes from the fine-tuning time required to
 431

Table 4: Training hours comparison on LLaMA2-7B for various positional encoding methods to realize a certain context window length.

Model	Up to 64K	Up to 256K	Up to 512K	Up to 1024K	Up to 2048K
LongLora	52.4h×8	OOM	OOM	OOM	OOM
LongRoPE	–	3×24h×1	3×24h×2	5×24h×8 + 98h	5×24h×8 + 202h
CRG NTK	2.6h×1	5.9h×1	13.8h×1	29.2h×1	64.7h×1

perform 128K for our configurations, and the LongRoPE method, since it does not use LoRA and MS Attention would only be longer.

Table 5: Performance comparison of different 7B parameter models on short standard benchmarks.

Model	ARC-c	HellaSwag	MMLU	TruthfulQA
Original LLaMA2-7B	53.1	78.6	45.6	39.0
Together	47.6	76.1	43.3	39.2
Code LLaMA	42.4	64.8	40.1	37.1
YaRN ($s = 16$)	52.4	78.7	42.4	38.2
YaRN ($s = 32$)	52.2	78.5	41.8	37.4
LongRoPE-2048k (ft=128k)	52.9	76.5	43.4	38.8
LongRoPE-2048k (ft=256k)	51.0	75.3	39.6	37.3
LLaMA2-7B-MS (ft=16K, CRGNTK)	51.9	77.8	41.9	40.1

4.4 ABLATION

In addition to the ablation comparisons of CRG NTK with other positional encodings and Full Attention with MS Attention shown in Tables 1 and 2, we also studied the impact of Lora on our performance, as shown in Table 6 and 9. The results indicate that full fine-tuning provides a slight performance improvement, but our efficient fine-tuning is very close to the full fine-tuning results. We also demonstrate that our MS Attention has better dynamic scaling extrapolation ability and forgetting at very long lengths, whereas Full Attention fits better at a fixed large scale. As a result, performance is slightly better when using FA for testing compared to the MS Attention fine-tuning method. However, our method outperforms FA when tested with different ratios, allowing for greater inheritance, which is beneficial for dynamic growth methods, as show in Table 8.

4.5 OTHERS

Additionally, in Appendix A.3, we highlight the advantages of our efficient Attention in pre-training. Our MS Attention method is capable of breaking the translation invariance of Full Attention and intrinsically achieving high-multiple length extrapolation. For example, we were able to achieve excellent performance using MS Attention and a context window length of 8K for training on PG19, and then tested on a length of 64K without a perplexity explosion, as shown in Tables 14 and Tables 16. In addition, we were able to achieve performance comparable to the current optimal efficient Attention by pre-training on Imagenet, as shown in Table 15. Finally we include the results of the tests on LongBench, as shown in Table 11.

5 CONCLUSION

In this paper, we build upon the LongLoRA framework, utilizing a mechanism of selection and merging within the Attention mechanism (referred to as MS Attention). By employing our approach on $1 \times$ A100 80GB GPUs, we achieve the same level of length extension for Llama2-7B as LongLoRA and LongRoPE does on $8 \times$ A100 80GB GPUs. Specifically, we extend the context length of Llama2-7B to 2M tokens, significantly reducing the resource requirements for handling long sequences. Moreover, We also indicate the possible extrapolation limit of the model: the number of layers of the model. Finally, the flexibility of our MS Attention mechanism allows for adjustable selection size and selection quantity, which, combined with restricting the attention range of each token, enables adaptable fitting and convergence rates.

REFERENCES

- 486
487
488 Ntk-aware scaled rope, 2023. URL [https://www.reddit.com/r/LocalLLaMA/](https://www.reddit.com/r/LocalLLaMA/comments/141z7j5/ntkaware_scaled_rope_allows_llama_models_to_have/)
489 [comments/141z7j5/ntkaware_scaled_rope_allows_llama_models_to_](https://www.reddit.com/r/LocalLLaMA/comments/141z7j5/ntkaware_scaled_rope_allows_llama_models_to_have/)
490 [have/](https://www.reddit.com/r/LocalLLaMA/comments/141z7j5/ntkaware_scaled_rope_allows_llama_models_to_have/).
- 491 Shengnan An, Yifei Li, Zeqi Lin, Qian Liu, Bei Chen, Qiang Fu, Weizhu Chen, Nanning Zheng, and
492 Jian-Guang Lou. Input-tuning: Adapting unfamiliar inputs to frozen pretrained models, 2022.
493
- 494 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du,
495 Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. Longbench: A bilingual,
496 multitask benchmark for long context understanding, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2308.14508)
497 [2308.14508](https://arxiv.org/abs/2308.14508).
- 498 Salima Behdenna, Fatiha Barigou, and Ghalem Belalem. Document level sentiment analysis: a
499 survey. *EAI endorsed transactions on context-aware systems and applications*, 4(13):e2–e2, 2018.
500
- 501 Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer,
502 2020.
503
- 504 Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending context window of
505 large language models via positional interpolation, 2023.
- 506 Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. Longlora:
507 Efficient fine-tuning of long-context large language models, 2024.
508
- 509 Ta-Chung Chi, Ting-Han Fan, Peter J Ramadge, and Alexander Rudnicky. Kerple: Ker-
510 nelized relative positional embedding for length extrapolation. In S. Koyejo, S. Mo-
511 hamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural*
512 *Information Processing Systems*, volume 35, pp. 8386–8399. Curran Associates, Inc.,
513 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/](https://proceedings.neurips.cc/paper_files/paper/2022/file/37a413841a614b5414b333585e7613b8-Paper-Conference.pdf)
514 [file/37a413841a614b5414b333585e7613b8-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/37a413841a614b5414b333585e7613b8-Paper-Conference.pdf).
- 515 Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas
516 Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy
517 Colwell, and Adrian Weller. Rethinking attention with performers, 2022.
518
- 519 Together Computer. Redpajama: an open dataset for training large language models, 2023. URL
520 <https://github.com/togethercomputer/RedPajama-Data>.
521
- 522 Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and
523 memory-efficient exact attention with io-awareness, 2022.
- 524 Yiran Ding, Li Lina Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan Yang,
525 and Mao Yang. Longrope: Extending llm context window beyond 2 million tokens, 2024.
526
- 527 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
528 and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.
529
- 530 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
531 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
532 L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas
533 Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b, 2023. URL [https://arxiv.](https://arxiv.org/abs/2310.06825)
534 [org/abs/2310.06825](https://arxiv.org/abs/2310.06825).
- 535 Hongye Jin, Xiaotian Han, Jingfeng Yang, Zhimeng Jiang, Chia-Yuan Chang, and Xia Hu.
536 Growlength: Accelerating llms pretraining by progressively growing training length, 2023.
537
- 538 Huan Yee Koh, Jiaxin Ju, Ming Liu, and Shirui Pan. An empirical survey on long document
539 summarization: Datasets, models, and metrics. *ACM Comput. Surv.*, 55(8), dec 2022. ISSN
0360-0300. doi: 10.1145/3545176. URL <https://doi.org/10.1145/3545176>.

- 540 Shanda Li, Chong You, Guru Guruganesh, Joshua Ainslie, Santiago Ontanon, Manzil Zaheer, Sumit
541 Sanghai, Yiming Yang, Sanjiv Kumar, and Srinadh Bhojanapalli. Functional interpolation for
542 relative positions improves long context transformers, 2024. URL [https://arxiv.org/
543 abs/2310.04418](https://arxiv.org/abs/2310.04418).
- 544 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
545 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the
546 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10012–10022, October 2021.
547
- 548 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- 549 Amirkeivan Mohtashami and Martin Jaggi. Landmark attention: Random-access infinite context
550 length for transformers, 2023. URL <https://arxiv.org/abs/2305.16300>.
- 551 Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. Yarn: Efficient context window
552 extension of large language models, 2023.
553
- 554 Ofir Press, Noah A. Smith, and Mike Lewis. Train short, test long: Attention with linear biases
555 enables input length extrapolation, 2022. URL <https://arxiv.org/abs/2108.12409>.
- 556 Jack W Rae, Anna Potapenko, Siddhant M Jayakumar, Chloe Hillier, and Timothy P Lillcrap.
557 Compressive transformers for long-range sequence modelling. *arXiv preprint*, 2019. URL
558 <https://arxiv.org/abs/1911.05507>.
- 559 Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System op-
560 timizations enable training deep learning models with over 100 billion parameters. In *Pro-
561 ceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery &
562 Data Mining*, KDD '20, pp. 3505–3506, New York, NY, USA, 2020. Association for Com-
563 puting Machinery. ISBN 9781450379984. doi: 10.1145/3394486.3406703. URL [https:
564 //doi.org/10.1145/3394486.3406703](https://doi.org/10.1145/3394486.3406703).
- 565 Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse
566 attention with routing transformers. *Transactions of the Association for Computational Linguistics*,
567 9:53–68, 2021. doi: 10.1162/tacl.a.00353. URL [https://aclanthology.org/2021.
568 tacl-1.4](https://aclanthology.org/2021.tacl-1.4).
- 569 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
570 Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov,
571 Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre
572 Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas
573 Scialom, and Gabriel Synnaeve. Code llama: Open foundation models for code, 2024.
574
- 575 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced
576 transformer with rotary position embedding, 2023. URL [https://arxiv.org/abs/2104.
577 09864](https://arxiv.org/abs/2104.09864).
- 578 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
579 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand
580 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language
581 models, 2023a.
582
- 583 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
584 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cris-
585 tian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu,
586 Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
587 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel
588 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,
589 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,
590 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,
591 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh
592 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen
593 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,
Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models,
2023b. URL <https://arxiv.org/abs/2307.09288>.

- 594 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
595 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, Dan Bikel, Lukas Blecher, Cris-
596 tian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu,
597 Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
598 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel
599 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,
600 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,
601 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,
602 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh
603 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen
604 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,
605 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models,
606 2023c.
- 607 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
608 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information*
609 *Processing Systems*, pp. 5998–6008, 2017.
- 610 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming lan-
611 guage models with attention sinks, 2024. URL <https://arxiv.org/abs/2309.17453>.
- 612
613 Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago
614 Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: Trans-
615 formers for longer sequences. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin
616 (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 17283–17297. Cur-
617 ran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2020/file/c8512d142a2d849725f31a9a7a361ab9-Paper.pdf)
618 [paper/2020/file/c8512d142a2d849725f31a9a7a361ab9-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/c8512d142a2d849725f31a9a7a361ab9-Paper.pdf).
- 619 Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon,
620 Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: Transformers for
621 longer sequences, 2021. URL <https://arxiv.org/abs/2007.14062>.
- 622
623 Chuanyang Zheng, Yihang Gao, Han Shi, Minbin Huang, Jingyao Li, Jing Xiong, Xi-
624 aozhe Ren, Michael Ng, Xin Jiang, Zhenguo Li, and Yu Li. Dape: Data-adaptive
625 positional encoding for length extrapolation. In A. Globerson, L. Mackey, D. Bel-
626 grave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural In-*
627 *formation Processing Systems*, volume 37, pp. 26659–26700. Curran Associates, Inc.,
628 2024. URL [https://proceedings.neurips.cc/paper_files/paper/2024/](https://proceedings.neurips.cc/paper_files/paper/2024/file/2f050fa9f0d898e3f265d515f50ae8f9-Paper-Conference.pdf)
629 [file/2f050fa9f0d898e3f265d515f50ae8f9-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/2f050fa9f0d898e3f265d515f50ae8f9-Paper-Conference.pdf).
- 630 Lei Zhu, Xinjiang Wang, Zhanhan Ke, Wayne Zhang, and Rynson W.H. Lau. Biformer: Vision
631 transformer with bi-level routing attention. In *Proceedings of the IEEE/CVF Conference on*
632 *Computer Vision and Pattern Recognition (CVPR)*, pp. 10323–10333, June 2023.
- 633 Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr:
634 Deformable transformers for end-to-end object detection, 2021.
- 635
636
637
638
639
640
641
642
643
644
645
646
647

A ADDITIONAL EXPERIMENTS

STATEMENT ON THE USE OF LARGE LANGUAGE MODELS

In the preparation of this paper, Large Language Models (LLMs) were used to assist with specific text and formatting tasks. The applications included:

- **Text Translation and Polishing:** Polishing English text to improve fluency and academic rigor.
- **Format Conversion:** Converting part of the formulas and tables content from other formats into \LaTeX code.

It is crucial to emphasize that all core ideas, theoretical derivations, experimental designs, result analyses, and final scientific conclusions were independently generated by the authors. The LLM served solely as a tool to enhance writing efficiency. All its outputs were rigorously reviewed, modified, and integrated by the authors to ensure accuracy and consistency with the paper’s core ideas.

A.1 ABLATION EXPERIMENTS

This subsection shows the results of the ablation experiments.

Model Configuration	2K	4K	8K	16K	32K
Llama-7B-MS-no.LoRA	7.06	6.96	6.78	6.56	6.93
Llama-7B-MS-LoRA	7.61	7.42	7.17	7.03	7.06

Configuration	Memory (MB)	Time (s/it)
With DeepSpeed & LoRA	36,786	— 6.37
Without DeepSpeed	31,524	— 6.19
Without LoRA	35,074	— 7.28
Without DeepSpeed & LoRA	40,132	— 6.38

Table 8: PPL on the PG19 (50 samples) with MS Attention or not. FT steps: 1000, FT length: 16384, ntk scale: 16384.

Model Configuration	16K	64K	128K	256K (scale=16384/32768/8192)
Llama-7B-MS-ntk16384	9.58	9.55	22.09	55.76/90.92/181.20
Llama-7B-FA-ntk16384	9.57	9.36	21.54	53.92/103.03/298.86
Llama-7B-MS-PI16	9.23	9.18	-	-
Llama-7B-FA-PI16	9.18	9.04	-	-

Table 9: Ablation experiment for PEFT. PEFT steps: 1000, Evaluation DataSet: PG19 validation.

PEFT	Training Context Length w/wo setting	Evaluation Context Length		
		32768	16384	8192
qkvo	16384-m16-s64	8.96	8.33	8.32
ko	16384-m16-s64	8.55	8.36	8.43

Table 10: PPL for Llama2-7B with MS Attention on PG19.

Parameters	PPL
MS-patch16-select512-merge64	7.17
MS-patch16-select64-merge4	6.95
MS-patch16-select16-merge4	6.83
MS-patch8-select32-merge8	6.18
MS-patch8-select32-merge4	6.09

A.2 LONGBENCH AND PROOFPILE

Table 11 shows Few-shot Learning and Code Completion Evaluation on LongBench: These tasks do not require chat instruct fine-tuning. We have observed significant improvements in "trec", "triviaqa", "lcc", and "repobench-p" using our method on Llama2-7B-4k comparing with Llama2-7B-chat-4k, especially outperforming other models in the LCC task.

Table 11: Few-shot Learning and Code Completion Evaluation on LongBench

Model	TREC	TriviaQA	SAMSum	LCC	RepoBench-P
GPT-3.5-Turbo-16k	68	91.4	41.7	54.7	53.6
LongChat-v1.5-7B-32k	63.5	82.3	34.2	53	55.3
XGen-7B-8k	65.5	77.8	25.3	38.6	38.6
InternLM-7B-8k	52	77.8	21.2	44.1	28.8
ChatGLM2-6B-32k	62.5	78.7	36.3	55.6	49.9
Vicuna-v1.5-7B-16k	71.5	86.2	40.8	51	43.5
ChatGLM3-6B-32k	79	87.1	38.2	57.66	54.76
Llama2-7B-chat-4k	61.5	77.8	40.7	52.4	43.8
Llama-7B-32k-MS-PI	72.6	85.7	40.6	61.95	49.09
Llama-7B-16k-MS-CRGNTK	72.6	87.5	39.3	65.24	54.76
Mistral-7B-16k-MS-CRGNTK	65.6	88.3	38.3	64.4	54.1

Table 12: Perplexity results on proofpile dataset. Fine-tuning with MS Attention; Testing with MS Attention or Full Attention.

Model	4K	8K	32K	64K	128K	256K
Llama2-7B-LongLora (ft=32K, PI)	3.01	2.78	2.58	2.57	-	-
Llama2-7B-LongRoPE (ft=256K)	3.71	3.50	2.60	2.36	2.26	1.88
Llama2-7B-MS (ft=16K, CRG NTK, test:FA)	3.12	2.76	2.56	2.58	2.60	2.68
Llama2-7B-MS (ft=16K, CRG NTK, test:MS)	2.96	2.64	2.32	2.31	2.26	2.28

A.3 MS ATTENTION PRETRAINING ON PG AND IMAGENET AND THE EXTRAPOLATION EFFECT

Tables 14 and 15 show the results of pretraining with MS Attention on PG19 and ImageNet, where our method outperforms various baselines. For PG19, test data is truncated when it is larger than the length of the test.

We further conducted experiments on automatic length extrapolation. After training on sequences of 8K and 16K tokens, our method shows stable perplexity even when extending inference to 8x the training length. Table 16 presents the results.

B RELATED WORK

Efficient Attention Mechanisms To fully exploit the inherent sparsity and positional relationships between tokens, a significant body of research has focused on developing efficient attention mechanisms. These mechanisms reduce the computational complexity of attention operations by focusing on a subset of tokens at each step, thus processing long sequences more efficiently or reducing resource consumption. Existing methods first preserve local features and then use various strategies to attend to more distant tokens. For instance, BigBirdZaheer et al. (2020) and PerformerChoromanski et al. (2022) use random patterns, LongformerBeltagy et al. (2020) and DETRZhu et al. (2021) use fixed patterns, while BiformerZhu et al. (2023) and Routing TransformerRoy et al. (2021) utilize relevance routing mechanisms. Our proposed relevance selection and merging mechanism adapts flexibly to various scenarios and is compatible with FlashAttention2Dao et al. (2022), achieving more efficient and general sparse attention.

Table 13: Comprehensive evaluation of our method compared to baselines on LLaMA2-7B and LLaMA2-13B models. PPL: Perplexity (lower is better), FA: Full Attention, MS: Multi-Scale Attention.

Model	Method	FT Length	PPL 8K	PPL 32K	PPL 64K	PPL 128K	Method (13B)	ARC-c	HellaSwag	MMLU	TruthfulQA
LLaMA2-7B	Codellama NTK	16K	3.71	2.74	2.55	2.71	LLaMA 2 Code LLaMA YaRN (s = 16) YaRN (s = 32) CRG NTK (Ours)	59.4	82.1	55.8	37.4
	YaRN	64K	3.51	2.65	2.42	> 10					
	Ours (FA/MS)	16K	3.12	2.56	2.58 / 2.32	2.60 / 2.26					
LLaMA2-13B	Codellama NTK	-	3.54	2.63	2.41	2.54	58.1	82.3	52.8	37.8	
	YaRN	-	3.25	2.50	2.29	> 10					
	YaRN	-	3.29	2.53	2.31	2.24					
	Ours (FA/MS)	-	3.05	2.42	2.30 / 2.22	2.31 / 2.16					

Table 14: Perplexity (PPL) results on the PG19.

Model	Params	PPL (PG19)
TransformerXL	-	36.3
Routing Transformer	-	33.3
Landmark Attention-200M	200M	14.55
Selection-Merging Attention-200M	200M	10.89

Table 15: Top-1 accuracy results on ImageNet.

Model	FLOPs (G)	Top-1 Accuracy (%)
VVT-T-12.9M	2.0	79.4
Swin Transformer-29M	4.2	81.3
Biformer-13.1M	2.2	81.4
Selection-Merging Attention-13.1M	2.1	82.1

Table 16: Perplexity across different sequence lengths for Select-Merge Attention. Test data is truncated when it is larger than the length of the test.

Method	8K	16K	32K	64K	128K
Select-Merge Attention (8K)	10.9	10.57	12.57	14.16	-
Select-Merge Attention (16K)	9.82	9.56	10.93	12.16	16.28

Positional Encoding Another research direction aimed at extending sequence length in LLMs focuses on positional encoding techniques, such as positional interpolation and extrapolation. Most pretrained models are trained on fixed-length sequences with fixed positional encodings, leading to performance degradation when extended to unknown positions. Therefore, numerous studies have analyzed the impact of positional encodings and modified them through interpolation or extrapolation to extend to longer sequences. For example, Position InterpolationChen et al. (2023), NTK-awareNtk (2023), YarnPeng et al. (2023), and LongRopEDing et al. (2024) mitigate the effects of pretrained positional encodings by using interpolation with different scales based on frequency importance, effectively extending sequence modeling lengths.

Efficient Fine-Tuning Efficient fine-tuning of LLMs has become a critical research direction, especially for handling long sequences. Techniques like Input-tuningAn et al. (2022) and LoRAHu et al. (2021) have shown significant promise in this area. Building on the LongLoraChen et al. (2024) method, we further optimize by using fewer parameters for fine-tuning, aiming to reduce the computational and memory overheads associated with fine-tuning large models on extended sequences while maintaining their performance on downstream tasks.

C POSITIONAL AWARENESS AND BREAKING TRANSLATION INVARIANCE

C.1 POSITIONAL AWARENESS:

Our algorithm(3.2, 4, E.2) leverages semantic token-based routing and selection mechanisms to capture boundary effects, thus disrupting the translation invariance inherent in global Attention mechanisms.

Semantic Token Routing and Selection: The routing process in our algorithm employs semantic tokens, which represent the semantics of a region, carrying contextual information similar to a sliding window in convolutional neural networks (CNNs). Due to the contextual variance surrounding each token, the derived semantic tokens differ, leading to the routing of different KV tokens. This results in the selection of distinct interpolation variables and, consequently, varying outputs. This mechanism allows our approach to effectively capture boundary effects, unlike traditional global Attention which maintains translation invariance.

C.2 EXTRAPOLATION THROUGH FINETUNING AND POSITION INTERPOLATION:

Our method can be fine-tuned during training and later applied to full Attention during inference, enabling significantly higher extrapolation factors. Specifically, our MS Attention mechanism can achieve up to twice the extrapolation compared to traditional methods, with potential reasons outlined as follows:

- **Improved Relative Positional Awareness:** The primary reason behind this extrapolation capability is our selection mechanism, which better extends the model’s awareness of relative positions. In full Attention, tokens near the boundaries rely on all preceding tokens for regression, potentially leading

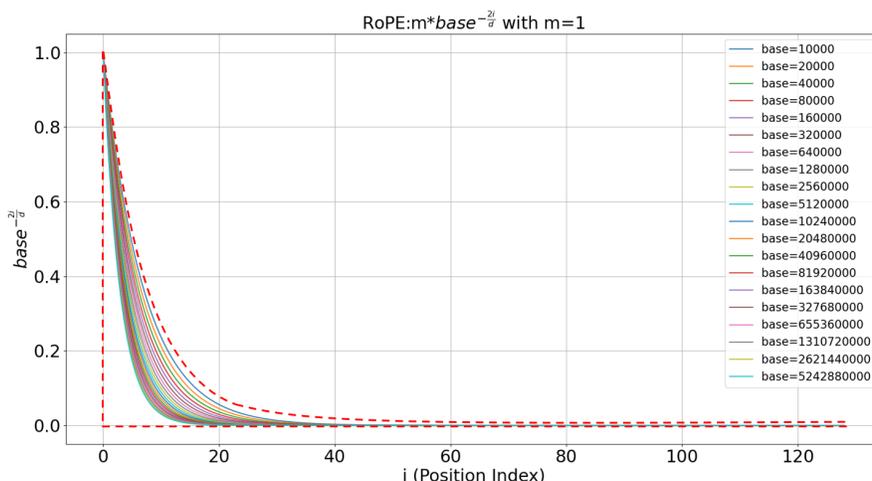


Figure 2: Dynamic scanning PE solution space.

the model to integrate all previous positional information. In contrast, our selection mechanism integrates a subset of positions into the final positional information, leading to a more accurate representation of relative positions. This method allows the model to assert the correct positional information even with different combinations of preceding positions.

- **Extrapolation Limits:** The reason for successful extrapolation up to twice the training length, but not beyond, lies in the finetuning process itself. During finetuning with a length l , the model learns the correct positional information and its integration within this range. Therefore, the model can correctly extrapolate up to $2l$, but further extrapolation may fail as it exceeds the recognized range.

For instance, when finetuning LLaMA2-7B with a length of $l = 16K$, the positional range for ‘Position_ids’ is set from 1 to 16K. By increasing the initial position range, e.g., choosing starting positions within the 1-64K range ($l_1 = 64K$), and setting the interpolation ratio for position interpolation as $80K/4K = 20$ or higher, we successfully achieve 100% accuracy on the 80K-length passkey task. Similarly, by exposing the model to even longer positions and using our adaptive positional interpolation encoding, we set $l_1 = 128K$ and successfully extrapolate to 144K. This analysis suggests that our method, even with simple position interpolation, has the potential to achieve near-infinite length extrapolation.

C.3 EXTRAPOLATION THROUGH FINETUNING AND NTK-BASED POSITIONAL ENCODING:

Combining our MS Attention with NTK-based positional encoding enables extrapolation by factors exceeding tenfold. This approach is explained using the radix theory proposed by the authors of RoPE:

- **Learning Relative Magnitudes:** Our MS Attention accurately learns the relative magnitudes of positional encodings. By incrementally increasing the radix base size, the model learns to represent the relative magnitudes across different radices. Additionally, by varying the starting positions, the model further refines its understanding of relative magnitudes within the same radix.

- **Future Work:** Future extensions will incorporate relative shifts between Q tokens and K tokens, allowing the model to sense deviations between different positions of Q and K.

This combination of MS Attention and NTK-based encoding showcases a significant potential for enhancing extrapolation capabilities, ultimately pushing the boundaries of positional understanding in large language models.

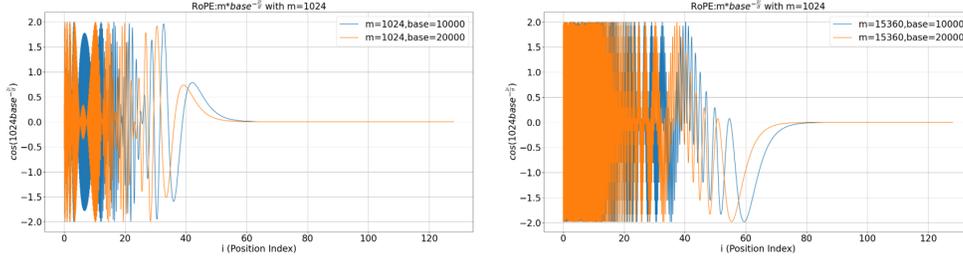


Figure 3: Differences in information encoded in different scale positions.

D POSITION INFORMATION LEARNING PROCESS ANALYSIS:

D.1 FORWARD PASS

In the computation process, the steps of Attention typically follow this sequence:

1. First, we introduce two components, P_c and P_s , to model cosine and sine positional encoding:

$$K_P = K \odot P_c + K_s \odot P_s$$

$$Q_P = Q \odot P_c + Q_s \odot P_s$$

The matrices P_c and P_s are structured as:

$$P_c = \begin{bmatrix} \cos \theta_0 & \cos \theta_0 & \dots & \cos \theta_{\frac{d}{2}-1} \\ \cos 2\theta_0 & \cos 2\theta_0 & \dots & \cos 2\theta_{\frac{d}{2}-1} \\ \vdots & \vdots & \ddots & \vdots \\ \cos n\theta_0 & \cos n\theta_0 & \dots & \cos n\theta_{\frac{d}{2}-1} \end{bmatrix}$$

$$P_s = \begin{bmatrix} \sin \theta_0 & \sin \theta_0 & \dots & \sin \theta_{\frac{d}{2}-1} \\ \sin 2\theta_0 & \sin 2\theta_0 & \dots & \sin 2\theta_{\frac{d}{2}-1} \\ \vdots & \vdots & \ddots & \vdots \\ \sin n\theta_0 & \sin n\theta_0 & \dots & \sin n\theta_{\frac{d}{2}-1} \end{bmatrix}$$

Moreover, we apply a learned matrix W_s to K_s as follows:

$$K_s = K \times W_s = K \times \begin{bmatrix} 0 & -1 & & 0 \\ 1 & 0 & & \\ & & \ddots & \\ 0 & & & 0 & -1 \\ & & & 1 & 0 \end{bmatrix}_{n \times n}$$

2. Compute similarity scores:

$$\text{Scores} = \frac{Q_P K_P^\top}{\sqrt{d_k}}$$

3. Apply Softmax:

$$\text{Weights} = \text{softmax}(\text{Scores})$$

4. Compute the output:

$$\text{Output} = \text{Weights} \times V$$

D.2 THEORETICAL ANALYSIS OF POSITIONAL ENCODING EXTRAPOLATION BASED ON LINEAR APPROXIMATION

D.2.1 NOTATION AND SETUP

Let d denote the model hidden dimension (assumed even) and $m := d/2$ the number of 2D rotational blocks used by RoPE-style position encoding. Let $\text{base} > 1$ denote the RoPE base. For block index $i \in \{0, \dots, m-1\}$ we define the block frequency parameter

$$\theta_i = \text{base}^{-2i/d},$$

and the phase at position p as

$$\phi_i(p) = p\theta_i.$$

We consider a Transformer trained on sequences of length M (positions $0, \dots, M-1$) with parameters Θ minimizing

$$\min_{\Theta} \mathcal{L}(\Theta; \mathcal{D}_M).$$

We are interested in extrapolation to length $L > M$, and in finding (possibly small) parameter adjustments $\Delta\Theta$ or deterministic transforms on positions/frequencies that preserve the loss behavior:

$$\mathcal{L}(\Theta + \Delta\Theta; \mathcal{D}_L) \approx \mathcal{L}(\Theta; \mathcal{D}_M).$$

D.2.2 LINEAR APPROXIMATION OF ATTENTION AND POSITIONAL TERMS

We adopt a first-order / linearized viewpoint of the attention exponential term. For queries \mathbf{q}_n and keys \mathbf{k}_m we approximate:

$$e^{\mathbf{q}_n \mathbf{k}_m^\top} \approx e^{\mathbf{q}_n \mathbf{k}_m^\top} \left(1 + f_p(\mathbf{q}_n, \mathbf{p}(n)) f_p(\mathbf{k}_m, \mathbf{p}(m))^\top - \mathbf{q}_n \mathbf{k}_m^\top \right),$$

where $f_p(\cdot, \cdot)$ denotes the position-dependent linear contribution. For RoPE,

$$f_p(\mathbf{q}_n, \mathbf{p}(n)) f_p(\mathbf{k}_m, \mathbf{p}(m))^\top = (Q_i P_c(n) + Q_{i,s} P_s(n)) (K_j P_c(m) + K_{j,s} P_s(m))^\top,$$

with $P_c(n), P_s(n)$ the cos/sin matrices evaluated at position n .

Proof. The original attention computation can be expressed as:

$$e^{\frac{1}{\sqrt{d}} \sum_{j=0}^d q_{nj} \cdot k_{sj}} \sum_s e^{\frac{1}{\sqrt{d}} \sum_{j=0}^d q_{nj} \cdot k_{sj} \cdot (\cos[(n-s)\theta_j] - 1) + \sum_{j=0}^d \frac{(-1)^{j+1}}{\sqrt{d}} q_{nj} \cdot k_{s, (j+\frac{d}{2}) \% d} \cdot (\sin[(n-s)\theta_j])} \quad (12)$$

Let us define the residual term:

$$r(\theta_1, \dots, \theta_{\frac{d}{2}-1}) = \frac{1}{\sqrt{d}} \sum_{j=0}^d q_{nj} \cdot k_{sj} \cdot (\cos[(n-s)\theta_j] - 1) + \sum_{j=0}^d \frac{(-1)^{j+1}}{\sqrt{d}} q_{nj} \cdot k_{s, (j+\frac{d}{2}) \% d} \cdot (\sin[(n-s)\theta_j]) \quad (13)$$

Since $\theta_j = \text{base}^{-\frac{2j}{d}}$, when base is large, we have the approximations:

$$\cos \theta_j = 1 - \frac{\theta_j^2}{2} + O(\theta_j^4) \quad (14)$$

$$\sin \theta_j = \theta_j + O(\theta_j^3) \quad (15)$$

Following the CRG-NTK method which has been proven to extend context window length and achieve excellent performance, we extend base to very large values (e.g., $\text{base} = 2^{40} \times 10000$). In this regime, $\cos \theta_j$ and $\sin \theta_j$ become very small, particularly for dimensions with larger j values.

Since $\frac{1}{\sqrt{d}} \sum_{j=0}^d q_{nj} \cdot k_{sj}$ is typically on the order of 1, we can apply Taylor expansion to $e^{r(\theta_1, \dots, \theta_{\frac{d}{2}-1})}$:

$$e^{r(\theta_1, \dots, \theta_{\frac{d}{2}-1})} = 1 + r(\theta_1, \dots, \theta_{\frac{d}{2}-1}) + O(r(\theta_1, \dots, \theta_{\frac{d}{2}-1})^2) \quad (16)$$

□

972 D.2.3 GRADIENT LINEAR-DOMINANCE

973
974 **Theorem D.1** (Linear dominance of positional gradients). *Under the linearized attention model*
975 *above and standard smoothness assumptions on activations, the gradient of the loss with respect to*
976 *positional encoding parameters P (or frequency parameters θ_i) is dominated by the linear positional*
977 *term:*

$$978 \frac{\partial \mathcal{L}}{\partial P} \approx \sum_{m,n} \frac{\partial \mathcal{L}}{\partial A_{mn}} \frac{\partial}{\partial P} [f_p(\mathbf{q}_n, \mathbf{P}(n)) f_p(\mathbf{k}_m, \mathbf{P}(m))^\top],$$

979
980 *where A_{mn} denotes pre-softmax attention scores. Consequently, positional modifications that*
981 *preserve these linear terms approximately preserve gradient directions/magnitudes.*
982

983 *Sketch.* The proof follows by substituting the first-order expansion into the attention expression,
984 differentiating the scalar loss $\mathcal{L}(A)$ via chain rule, and observing that higher-order derivatives
985 involving the exponential factor or the quadratic term $\mathbf{q}_n \mathbf{k}_m^\top$ are of smaller magnitude under mild
986 norm bounds on \mathbf{q}, \mathbf{k} . Hence the linear position-dependent contribution dominates the derivative w.r.t.
987 positional parameters. \square
988

989 D.3 GRADIENT CONTINUITY AND FREQUENCY SENSITIVITY

990
991 **Lemma D.1** (Frequency-wise gradient expression). *For RoPE-like encoding, the gradient w.r.t. the*
992 *i -th frequency parameter θ_i can be written as*

$$993 \frac{\partial \mathcal{L}}{\partial \theta_i} = \sum_{m,n} \frac{\partial \mathcal{L}}{\partial A_{mn}} \cdot \frac{\partial A_{mn}}{\partial \theta_i},$$

994
995 *with (up to the block-specific linear coefficients)*

$$996 \frac{\partial A_{mn}}{\partial \theta_i} = -(m-n) \left[Q_i^{(m)} K_i^{(n)} \sin((m-n)\theta_i) + Q_{i,s}^{(m)} K_{i,s}^{(n)} \cos((m-n)\theta_i) \right].$$

997
998 *Sketch.* Differentiate the block-wise rotated representations $R(\phi_i(p))$ with respect to θ_i where
999 $\phi_i(p) = p\theta_i$. By the product and chain rules, derivatives bring down the multiplicative factor p
1000 (hence the difference $m-n$ appears after combining query/key contributions), and the trigonometric
1001 derivative yields sine/cosine terms as shown. \square
1002
1003
1004

1005 **Remark D.1.** *Lemma D.1 highlights that gradients for different frequency components are weighted*
1006 *by the difference index $(m-n)$ and trigonometric factors; thus, phases beyond the training interval*
1007 *may flip signs or amplify gradients unless controlled.*
1008

1009 D.4 PHASE-CONSISTENCY AND GRADIENT-MATCHING CONDITIONS

1010
1011 **Definition D.1** (Effective training phase set). *The set of phases exposed during training length M is*

$$1012 \Phi_M = \{ m\theta_i \bmod 2\pi : m = 0, \dots, M-1, i = 0, \dots, m-1 \}.$$

1013
1014 *The convex hull $\text{conv}(\Phi_M)$ denotes the (real-valued) interval(s) spanned by these phases (we treat*
1015 *phases lifted to the real line under the principal branch approximation).*

1016 **Theorem D.2** (Phase extrapolation condition). *A necessary condition for successful extrapolation (so*
1017 *that gradients and attention statistics behave similarly on length L as on M) is that for extrapolated*
1018 *positions $p \in \{M, \dots, L-1\}$ the adjusted phases $\phi'_i(p)$ lie in $\text{conv}(\Phi_M)$. Equivalently, under a*
1019 *position-frequency scaling*

$$1020 p\theta_i \mapsto \phi'_i(p) = \frac{p}{s_1} \cdot \frac{\theta_i}{s_2^{2i/d}} = \frac{p}{s_1 s_2^{2i/d}} \theta_i,$$

1021
1022
1023 *we must have for the maximal position $p = L-1$ that*

$$1024 0 \leq \frac{L-1}{s_1 s_2^{2i/d}} \leq M-1. \quad (17)$$

Proof. If the scaled phases lie outside the training-phase interval, trigonometric factors in Lemma D.1 will take values not seen during training and can cause sign changes or amplitude differences in $\partial A_{mn}/\partial\theta_i$. Hence keeping $\phi'_i(p)$ inside $\text{conv}(\Phi_M)$ (and in particular for $p = L - 1$) ensures the phase-dependent terms remain in the range observed during training. Solving $\phi'_i(L - 1) \leq (M - 1)\theta_i$ yields equation 17. \square

D.5 VARIATIONAL DERIVATION OF SCALING FACTORS

We now pose a quantitative objective to choose (s_1, s_2) . Define the gradient-discrepancy objective

$$\mathcal{J}(s_1, s_2) = \sum_{i=0}^{m-1} \left\| \frac{\partial \mathcal{L}}{\partial \theta_i}(L; s_1, s_2) - \frac{\partial \mathcal{L}}{\partial \theta_i}(M) \right\|_2^2.$$

Under the linear approximation (Theorem D.1) the mapping $(s_1, s_2) \mapsto \partial \mathcal{L}/\partial \theta_i(L; s_1, s_2)$ is approximately multiplicative:

$$\frac{\partial \mathcal{L}}{\partial \theta_i}(L; s_1, s_2) \approx \frac{1}{s_1 s_2^{2i/d}} \frac{\partial \mathcal{L}}{\partial \theta_i}(M).$$

Thus minimizing \mathcal{J} w.r.t. (s_1, s_2) reduces approximately to minimizing

$$\tilde{\mathcal{J}}(s_1, s_2) = \sum_{i=0}^{m-1} \left\| \left(\frac{1}{s_1 s_2^{2i/d}} - 1 \right) \frac{\partial \mathcal{L}}{\partial \theta_i}(M) \right\|_2^2.$$

Proposition D.1 (Stationary condition for multiplicative approximation). *A stationary (and in some cases minimizer) solution of $\tilde{\mathcal{J}}$ satisfies approximately*

$$\frac{1}{s_1 s_2^{2i/d}} \approx 1 \quad \text{for those } i \text{ with large } \left\| \frac{\partial \mathcal{L}}{\partial \theta_i}(M) \right\|_2.$$

Hence optimal (s_1, s_2) tends to match the most “energetic” frequency components.

Sketch. Take partial derivatives of $\tilde{\mathcal{J}}$ w.r.t. s_1, s_2 and set to zero. The resulting equations emphasize frequencies with large pre-factor norms; when few frequencies dominate, one can approximately match those frequencies, leading to the multiplicative conditions for those indices. \square

D.6 FINAL ALGEBRAIC FORM OF THE EXTRAPOLATION CONSTRAINT

Using $\theta_i = \text{base}^{-2i/d}$, condition (17) can be rewritten as

$$\frac{L - 1}{s_1} \leq (M - 1) s_2^{2i/d} \iff \frac{(L - 1)/s_1}{(\text{base}/s_2)^{2i/d}} \leq \frac{M - 1}{\text{base}^{2i/d}}.$$

Therefore the (conservative) target constraint for extrapolation becomes

$$\boxed{\frac{(L - 1)/s_1}{(\text{base}/s_2)^{\frac{2i}{d}}} \in \left[0, \frac{M - 1}{\text{base}^{\frac{2i}{d}}} \right]} \quad (18)$$

for every frequency index i that we require to be “covered” by training. Note that equality corresponds to exact matching of the maximal phase; the interval allows any scaled maximal phase within the training-covered range.

D.7 ANALYSIS OF TYPICAL EXTRAPOLATION STRATEGIES

Position interpolation (PI). Set $s_2 = 1$. Then condition equation 17 reduces to $(L - 1)/s_1 \leq M - 1$. Choosing

$$s_1 = \frac{L - 1}{M - 1} \approx \frac{L}{M}$$

maps the target maximum position to the training maximum position and (under multiplicative gradient approximation) yields exact gradient-amplitude matching across all frequencies (Theorem D.1). PI is therefore the uniform, parameter-free approach to preserve gradient amplitudes.

Frequency/base scaling (NTK-aware-like). Set $s_1 = 1$. For each frequency i , exact matching requires

$$s_2^{2i/d} = \frac{L-1}{M-1} \implies s_2 = \left(\frac{L-1}{M-1}\right)^{\frac{d}{2i}}.$$

Since the RHS depends on i , a single scalar s_2 cannot match all frequencies. In practice, choosing $s_2 > 1$ biases the transform to preserve higher-frequency gradient amplitudes (large i), which is sometimes desired when high-frequency structure carries important sequence-local information.

Hybrid and optimization-based calibration. Treat (s_1, s_2) as tunable and directly optimize $\mathcal{J}(s_1, s_2)$ on a small calibration set (or via projected gradient descent). This hybrid method can trade off which frequency bands to prioritize according to their gradient energy, and is theoretically justified by the variational derivation above.

D.8 ALGORITHM: PROJECTED GRADIENT CALIBRATION ON (s_1, s_2)

Algorithm 1 Calibrate (s_1, s_2) via projected gradient descent

- 1: initialize $s_1 \leftarrow 1, s_2 \leftarrow 1$
 - 2: **for** $t = 1, 2, \dots$ **do**
 - 3: compute $\Delta_i = \partial\mathcal{L}/\partial\theta_i(L; s_1, s_2) - \partial\mathcal{L}/\partial\theta_i(M)$ for all i
 - 4: $s_1 \leftarrow s_1 - \eta_1 \frac{\partial}{\partial s_1} \sum_i \|\Delta_i\|^2$
 - 5: $s_2 \leftarrow s_2 - \eta_2 \frac{\partial}{\partial s_2} \sum_i \|\Delta_i\|^2$
 - 6: project $s_1 \leftarrow \max(s_1, L/M), s_2 \leftarrow \max(s_2, 1)$
 - 7: **if** $\sum_i \|\Delta_i\|^2 \leq \varepsilon$ **then break**
 - 8: **end if**
 - 9: **end for**
-

Under standard smoothness and boundedness assumptions on the mapping $(s_1, s_2) \mapsto \partial\mathcal{L}/\partial\theta$, sufficiently small learning rates η_1, η_2 ensure monotone decrease of the objective and convergence to a local minimum.

D.9 CONCLUDING REMARKS

- The key analytic object is the scaled phase $\phi'_i(p) = \frac{p}{s_1 s_2^{2i/d}} \theta_i$. Ensuring $\phi'_i(L-1)$ lies in the training-phase interval $[0, (M-1)\theta_i]$ for the set of important frequencies i yields inequality equation 18.
- Position interpolation offers a simple, parameter-free method that preserves gradient magnitudes across frequencies in the multiplicative linear approximation. Frequency (base) scaling adds a degree of freedom to prioritize certain frequencies but cannot match all frequencies exactly with a single scalar.
- Practical schemes combine interpolation, frequency-aware scaling, and small calibration fine-tuning on (s_1, s_2) to obtain robust extrapolation; the above derivations give clear, testable predictions for such hybrid methods.
- All theoretical claims are made under a first-order / linearized approximation of attention; nevertheless the derived conditions are both interpretable and useful as initializations or constraints for empirical calibration.

D.9.1 STRATEGIES FOR GENERALIZATION:

1. Increase M: Increasing the length of the original sequence M allows us to cover a broader range of positions.
2. Adjust Token Positions: Using a method like Position Interpolation (PI), we can adjust the token positions to fall within the range that the model has learned.

3. Increase Base: Using methods such as NTK-aware RoPE, we can increase the base value. Since the positional base is inherently discrete across dimensions, it is only necessary to adjust a subset of the positional encodings correctly to achieve proper adjustment. An advantage of this method is the reuse of previously learned positional bases. For example, bases like $16^{\frac{2}{3}}$ and $256^{\frac{1}{3}}$ can share learned weights, allowing us to generalize to various sequence lengths dynamically.

4. Dimensional Generalization: Similarly, we can extend the generalization mechanism across sequence dimensions. By selectively using parts of the position encodings that fall within the learned range, we can efficiently cover new positions while maintaining computational efficiency.

5. Scaling Numerator and Denominator: We can also consider scaling both the numerator and denominator of position encodings to achieve a more generalized weighting over positional information.

D.10 EXTRAPOLATION CAPACITY LIMIT: MODEL DEPTH

We first analyze potential causes from the perspective of the backward learning process. For the entire training corpus, which we assume to be a sequence of length L_S , we partition it into chunks of the pre-training length. Although these chunks are not explicitly modeled as a continuous sequence (i.e., they are not concatenated for end-to-end sequence modeling), the linear superposition of backpropagated gradients across chunks allows the model to capture partial linear relationships between them.

The gradient information propagates backward through each layer, enabling the model to learn most inter-chunk dependencies. Theoretically, between any two adjacent layers, the model can learn information equivalent to a concatenated sequence of two chunks. Consequently, the ultimate extrapolation limit scales linearly with the number of model layers. In the following, we analyze the specific inverse process to illustrate how the above process is implemented:

The gradients of the MLP and the linear layer

1. Linear Layer

Given the linear projection:

$$\mathbf{Y} = \mathbf{X}\mathbf{W} + \mathbf{b}$$

where: $\mathbf{X} \in \mathbb{R}^{n \times d_{in}}$ is the input, $\mathbf{W} \in \mathbb{R}^{d_{in} \times d_{out}}$ is the weight matrix, $\mathbf{b} \in \mathbb{R}^{d_{out}}$ is the bias.

The gradients during backpropagation (given upstream gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{Y}}$) are:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \mathbf{X}^\top \frac{\partial \mathcal{L}}{\partial \mathbf{Y}}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{b}} = \sum_{i=1}^n \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Y}} \right)_i, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{X}} = \frac{\partial \mathcal{L}}{\partial \mathbf{Y}} \mathbf{W}^\top$$

Thus, it is clear:

$$\frac{\partial \mathcal{L}_{\text{cat}(S_m, S_n)}}{\partial \mathbf{W}} = \mathbf{X}_{S_m}^\top \frac{\partial \mathcal{L}_{S_m}}{\partial \mathbf{Y}} + \mathbf{X}_{S_n}^\top \frac{\partial \mathcal{L}_{S_n}}{\partial \mathbf{Y}}$$

2. MLP (Multilayer Perceptron)

With two linear transformations and activation:

$$\mathbf{Z} = \sigma(\mathbf{X}\mathbf{W}_1 + \mathbf{b}_1), \quad \mathbf{Y} = \mathbf{Z}\mathbf{W}_2 + \mathbf{b}_2$$

where σ denotes the activation function (e.g., GELU).

(a) Second Layer Gradients (identical to linear layer):

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_2} = \mathbf{Z}^\top \frac{\partial \mathcal{L}}{\partial \mathbf{Y}}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{b}_2} = \sum_{i=1}^n \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Y}} \right)_i$$

(b) First Layer Gradients:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{Z}} = \frac{\partial \mathcal{L}}{\partial \mathbf{Y}} \mathbf{W}_2^\top, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{X}} = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Z}} \odot \sigma'(\mathbf{X}\mathbf{W}_1 + \mathbf{b}_1) \right) \mathbf{W}_1^\top$$

Here \odot denotes element-wise multiplication, and σ' is the derivative of the activation function.

(c) First Layer Parameter Gradients:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} = \mathbf{X}^\top \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Z}} \odot \sigma'(\mathbf{X}\mathbf{W}_1 + \mathbf{b}_1) \right), \quad \frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} = \sum_{i=1}^n \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Z}} \odot \sigma'(\mathbf{X}\mathbf{W}_1 + \mathbf{b}_1) \right)_i$$

From this formula, it can be seen that for sequences whose dimensions are in the Sequence direction, the linear superposition of their gradients is the same as using the inverse gradient of the entire sequence, i.e., these two middle layers completely model the splicing of chunks

The gradients of the Attention layer

Given the attention computation:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V}$$

where: $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{n \times d_k}$ are queries and keys, $\mathbf{V} \in \mathbb{R}^{n \times d_v}$ are values, n is sequence length, d_k is key dimension.

3. Attention

1. Output Gradient: Given upstream gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{O}}$ where $\mathbf{O} = \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$:

2. Attention Weights: Let $\mathbf{A} = \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}$ and $\mathbf{P} = \text{softmax}(\mathbf{A})$

Then:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{P}} = \frac{\partial \mathcal{L}}{\partial \mathbf{O}} \mathbf{V}^\top$$

3. Softmax Gradient:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{A}} = \mathbf{P} \odot \left(\frac{\partial \mathcal{L}}{\partial \mathbf{P}} - \sum_{i=1}^n P_{ij} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{P}} \right)_{ij} \right)$$

4. Query/Key Gradients:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{Q}} = \frac{1}{\sqrt{d_k}} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{A}} \right) \mathbf{K}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{K}} = \frac{1}{\sqrt{d_k}} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{A}} \right)^\top \mathbf{Q}$$

5. Value Gradient:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{V}} = \mathbf{P}^\top \frac{\partial \mathcal{L}}{\partial \mathbf{O}}$$

For simplicity, we analyze the last two layers to learn the gradient information of one chunk first and then perform Attention on the latter chunk, with the error of using two chunks spliced together for Attention to illustrate that the model is able to learn a part of the information of the spliced Attention, thus realizing a certain degree of extrapolation ability.

$$\Delta_o = \text{Attention}(\mathbf{Q}_{S_n}, \mathbf{K}_{\text{cat}(S_m, S_n)}, \mathbf{V}_{\text{cat}(S_m, S_n)}) - \text{Attention}(\mathbf{Q}_{S_n}^{\text{new}}, \mathbf{K}_{S_n}^{\text{new}}, \mathbf{V}_{S_n}^{\text{new}})$$

$$\Delta_o = \text{softmax} \left(\frac{\mathbf{Q}_{S_n} \mathbf{K}_{S_m}^\top}{\sqrt{d_k}} \right) \mathbf{V}_{S_m} + \text{softmax} \left(\frac{\mathbf{Q}_{S_n} \mathbf{K}_{S_n}^\top}{\sqrt{d_k}} \right) \mathbf{V}_{S_n} - \left(\text{softmax} \left(\frac{\mathbf{Q}_{S_n}^{\text{new}} \mathbf{K}_{S_n}^{\text{new} \top}}{\sqrt{d_k}} \right) \mathbf{V}_{S_n}^{\text{new}} \right)$$

where:

$$\mathbf{Q}_{S_n}^{\text{new}} = \mathbf{Q}_{S_n} - \frac{\partial \mathcal{L}}{\partial \mathbf{Q}_{S_m}}, \quad \mathbf{K}_{S_n}^{\text{new}} = \mathbf{K}_{S_n} - \frac{\partial \mathcal{L}}{\partial \mathbf{K}_{S_m}}, \quad \mathbf{V}_{S_n}^{\text{new}} = \mathbf{V}_{S_n} - \frac{\partial \mathcal{L}}{\partial \mathbf{V}_{S_m}},$$

For ease of analysis we freeze the gradients of Q and V (gradients with Q and V will have smaller errors) and observe only the output after learning K.

$$\mathbf{O}_{S_n}^{\text{new}} = \left(\text{softmax} \left(\frac{\mathbf{Q}_{S_n} \left(\mathbf{K}_{S_n} - \frac{1}{\sqrt{d_k}} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{A}} \right)^\top \mathbf{Q}_{S_m} \right)^\top}{\sqrt{d_k}} \right) \right) \mathbf{V}_{S_n}$$

Observing only the propagation of information, it can be seen that the learned Attention is able to approximate the Attention with Q with the information of the previous sequence S_m , thus completing a portion of the complete Attention information, and then each chunk, with each layer of attention, is able to realize an approximate Attention of the L-fold sequence, thus realizing a certain amount of extrapolation. Finally more detailed and rigorous proofs we leave for future work to investigate.

D.10.1 OTHER EXTRAPOLATION ERROR ANALYSIS:

To analyze the error during extrapolation, we assume a linear relationship between different segments, where each part is responsible for the corresponding position transformation. Specifically, the part of the gradient involving P_c is responsible for handling the ground-truth P_c . Let the positional difference between seen and unseen sequences be:

$$Q_Z - Q_L = XW_{QZ} - X \left(W_Q + \frac{\partial \mathcal{L}}{\partial W_Q} \right)$$

Now, consider a fine-tuning setup where we train on sequences of length M . For the first M tokens, the model performs without loss. However, for tokens beyond M , the extrapolation requires additional adjustments. For the first M tokens, the gradient is:

$$X_{:M}W_{QZ} = X_{:M} \left(W_Q + \frac{\partial \mathcal{L}}{\partial W_Q} \right) + L(V, K)$$

For the tokens beyond M , the extrapolation error is:

$$(X_M:W_{QZ}) \odot P_c - (X_M:(W_Q + \frac{\partial \mathcal{L}}{\partial W_Q})) \odot P_c$$

Letting it generalize continuously over a fixed length enables the continuous reduction of the above error L, thus reducing the extrapolation of the final position information.

E ANALYSIS OF ATTENTION WITH CORRELATION-AWARE SELECTION AND MERGING:

E.1 APPROXIMATION AND CONVERGENCE OF OUR MS ATTENTION TO VARIOUS EFFICIENT ATTENTION METHODS

Our proposed Select-Merge Attention (MS Attention) can approximate and converge to the optimal solutions of many efficient Attention mechanisms that utilize a subset of key-value (KV) pairs. Below, we describe the computation process using the notations introduced in the algorithm description:

Selection:

$$A_s = Q'_s K_s'^T, \quad \text{Idx} = \text{topkIndex}(A_s)$$

Merge:

$$Q_s = \text{merge}_q(Q'_s), \quad \text{Idx} = \text{filter}(\text{merge}_q(\text{Idx})), \quad KV_s = \text{Select}(KV, \text{Idx})$$

Final Attention:

$$O = \text{Attention}(Q_s, K_s, V_s)$$

E.1.1 LANDMARK ATTENTION

Landmark Attention adjusts the attention scores by incorporating landmark tokens, with the output expressed as:

$$O = (\text{softmax}(QK^T) \cdot \text{repeat}(\text{softmax}(QG^T), \text{blocksize}, \text{dim} = -1))V$$

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Algorithm 2 Correlation-based Selection & Merging for Sparse Attention

Require: $Q, K, V \in \mathbb{R}^{B \times H \times N \times D}$ \triangleright batch B , heads H , tokens N , dim D

Require: segment sizes s_q, s_k ; number of Q -segments $n_s \leftarrow \lceil N/s_q \rceil$; $topk$; merge factor m (divisor of n_s)

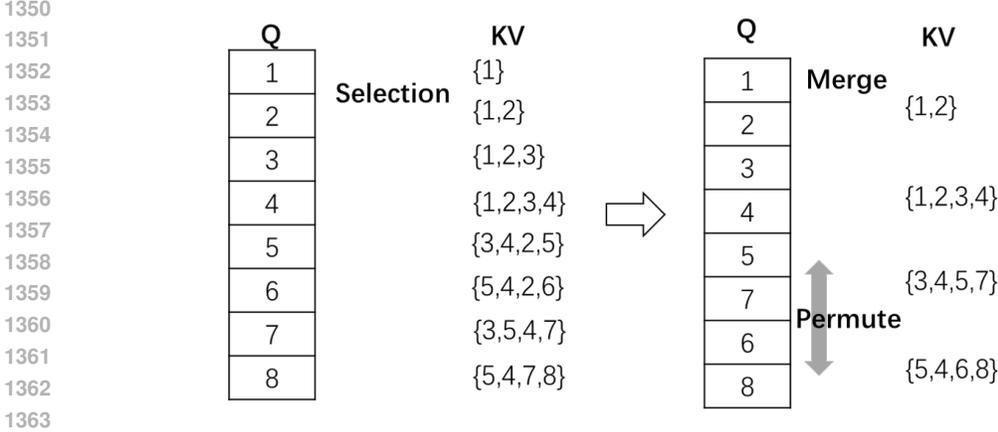
Require: compress function $\text{compress}(\cdot)$ (e.g. mean / learned semantic token)

Require: similarity function $\text{sim}(\cdot, \cdot)$ (e.g. dot-product)

Require: mask function $\text{mask}(\cdot)$ (to prevent leakage); optional final select length top_final

Ensure: Output $\in \mathbb{R}^{B \times H \times N \times D}$ (approximate attention output)

- 1: **Segmentation:**
 - 2: Split Q into $Q_s \in \mathbb{R}^{B \times H \times n_s \times s_q \times D}$ by segments \triangleright last segment padded if needed
 - 3: Split K, V into $K_s \in \mathbb{R}^{B \times H \times n_k \times s_k \times D}, V_s$ similarly
 - 4: **Compress region representatives:**
 - 5: $Q'_s \leftarrow \text{compress}(Q_s), Q'_s \in \mathbb{R}^{B \times H \times n_s \times D}$ \triangleright one token per Q -segment
 - 6: $K'_s \leftarrow \text{compress}(K_s), K'_s \in \mathbb{R}^{B \times H \times n_k \times D}$
 - 7: **Compute relevance / selection:**
 - 8: $S \leftarrow \text{sim}(Q'_s, K'_s)$ $\triangleright S \in \mathbb{R}^{B \times H \times n_s \times n_k}$
 - 9: $S \leftarrow \text{mask}(S)$ \triangleright mask out invalid / future / leakage entries
 - 10: $select_idx \leftarrow \text{topk_indices}(S, k = topk)$ $\triangleright (B, H, n_s, topk)$: top- k K_s indices per Q_s
 - 11: **(Optional) Permutation for sharing/clustered KV (pretrain/infer):**
 - 12: Compute permutation π from S (group similar Q_s) \triangleright optional, improves KV sharing
 - 13: Apply π to Q_s and $select_idx$ (and record inverse π^{-1})
 - 14: **Merging Q segments:**
 - 15: $n_{ms} \leftarrow n_s/m$ \triangleright number of merged Q regions
 - 16: Reshape / merge Q_s into $Q_{ms} \in \mathbb{R}^{B \times H \times n_{ms} \times (m \cdot s_q) \times D}$ by concatenating every m adjacent Q_s
 - 17: Optionally compress merged Q_{ms} to one representative per merged region (if desired)
 - 18: **Transform and merge selected indices:**
 - 19: Reshape $select_idx$ from $(B, H, n_s, topk)$ to $(B, H, n_{ms}, m, topk)$ by grouping every m adjacent Q_s rows
 - 20: Permute to $(B, H, n_{ms}, topk, m)$ and then reshape to $(B, H, n_{ms}, topk \cdot m)$ \triangleright now each merged Q region has $topk \cdot m$ candidate K_s
 - 21: For each merged row do: $mselect_idx \leftarrow \text{unique_preserve_order}(\text{candidates})$ \triangleright remove duplicates while keeping relevance order
 - 22: $qmselect_idx \leftarrow \text{top}_n(mselect_idx, n = top_final)$ \triangleright final top- n K_s groups per merged Q region; shape $(B, H, n_{ms}, top_final)$
 - 23: **Gather relevant K/V segments:**
 - 24: For each merged Q region, gather the corresponding K_s and V_s regions indexed by $qmselect_idx$ and similarly V_{ms}^g :
$$K_{ms}^g \in \mathbb{R}^{B \times H \times n_{ms} \times top_final \times s_k \times D}$$
 - 25: **Compute attention on merged blocks:**
 - 26: **for** each merged region index $r = 1 \dots n_{ms}$ **do**
 - 27: Flatten $Q_{ms}[\dots, r]$ to queries of shape $(B, H, qlen = m \cdot s_q, D)$
 - 28: Flatten $K_{ms}^g[\dots, r]$ and $V_{ms}^g[\dots, r]$ across selected segments to $(B, H, klen, D)$ where $klen = top_final \cdot s_k$
 - 29: Compute local attention outputs:
$$O_{ms}[\dots, r] \leftarrow \text{Attention}(Q_{ms}[\dots, r], K_{ms}^g[\dots, r], V_{ms}^g[\dots, r])$$
 - 30: **end for**
 - 31: **Unmerge and scatter back to token positions:**
 - 32: Split each $O_{ms}[\dots, r]$ back to the m original Q_s regions and to individual tokens, then reorder according to π^{-1} if permutation was applied and get $\hat{O} \in \mathbb{R}^{B \times H \times N \times D}$.
 - 33: Attention computation process can be fully parallelised.
 - 34: **return** \hat{O}
-



1364 **Figure 4: Overview of Merging and Selection Attention Mechanism (MS Attention).** The MS
 1365 Attention mechanism involves two main steps. In the first step, the QKV tensors are split into regions,
 1366 and a single token is used to represent each region. Subsequently, the regional representatives are
 1367 used to compute dot products or other similarity measures to select the most relevant KV regions
 1368 for each Q region. For example, Q regions 5, 6, 7 and 8 select KV regions {3, 4, 2, 5}, {5, 4, 2, 6},
 1369 {3, 5, 4, 7} and {5, 4, 7, 8} respectively. In the second step, for each Q region, tokens are merged
 1370 with their adjacent or related tokens after permuting. The union of the selected KV regions is taken,
 1371 and the top-n regions are chosen. For example, combine Q regions 6 and 8, along with their selected
 1372 KV regions, resulting in {5, 4, 2, 7, 6, 8}. To ensure tensor consistency, we select the top k regions
 1373 from the merged set. If $k = 4$, the final selection is {5, 4, 6, 8}. The reason for retaining 6 and 8
 1374 is because by default we believe that the local region is important and must be preserved, and that
 1375 the local region does not perform scoring calculations. Finally, each merged Q region performs
 1376 self-attention with its selected relevant KV regions.

1377

1378 When setting the split size of Q to 1 and the split size of K to blocksize, our MS Attention also
 1379 calculates $\text{softmax}(QK_s^T)$ and $\text{softmax}(QK_s'^T)$, which fully enables the adjustment of attention
 1380 scores using semantic tokens. This operation can be approximated under other settings as well.

1381 Alternatively, if the score adjustment is performed without directly multiplying $\text{softmax}(QK_s'^T)$,
 1382 treating semantic tokens as regular KV tokens during Attention computation, the output can be
 1383 expressed as:

$$1384 \frac{\text{sumexp}(QK_s^T) \times \text{softmax}(QK_s^T)V_s}{\text{sumexp}(QK_s^T) + \text{sumexp}(QK_s'^T)} + \frac{\text{sumexp}(QK_s'^T) \times \text{softmax}(QK_s'^T)V_s'}{\text{sumexp}(QK_s^T) + \text{sumexp}(QK_s'^T)}$$

1385 where V_s' is a linear combination of V_s or V . This approach can also directly adjust the attention
 1386 coefficient of V_s in the output, allowing convergence to the optimal solution based on the task.

1390 E.1.2 BiFORMER

1391 The BiFormer computation process is as follows:

$$1392 A_r = Q_r(K_r)^T, \quad I_r = \text{topkIndex}(A_r), \quad KV_g = \text{Select}(KV, I_r)$$

$$1393 O = \text{Attention}(Q, K_g, V_g)$$

1394 When the Merge size is set to 1, our method fully degenerates to BiFormer. To achieve lower
 1395 complexity, BiFormer requires setting the initial region Q_r, K_r relatively large and then compressing
 1396 it as a region representative. This approach reduces the accuracy of selection.

1397 When merge size is greater than 1, our method performs more fine-grained partitioning and selec-
 1398 tion. If BiFormer's selection is optimal, our algorithm can also converge to this optimal selection.
 1399 Overall, compared to BiFormer, our method has a larger convergence space, and under the same fine
 1400 partitioning and selection, our method has relatively lower computational complexity.

1404 E.1.3 ROUTING ATTENTION METHODS

1405 For the Routing Transformer, the following update rule is applied:

$$1406 \mu \leftarrow \lambda\mu + \frac{(1-\lambda)}{2} \sum_{i:\mu(Q_i)=\mu} Q_i + \frac{(1-\lambda)}{2} \sum_{j:\mu(K_j)=\mu} K_j$$

1407 This can be rewritten as:

$$1408 \mu \leftarrow \lambda\mu + \frac{(1-\lambda)}{2} \operatorname{argmax}_{qu}(\mu Q^T)Q + \frac{(1-\lambda)}{2} \operatorname{argmax}_{qu}(\mu K^T)K$$

1409 where

$$1410 \operatorname{argmax}_{qu}(\cdot) = \begin{cases} 1 & \operatorname{argmax}(\cdot, \dim = -2) \\ 0 & \text{otherwise} \end{cases}$$

$$1411 \operatorname{Idx}_Q = \operatorname{topkIndex}(\mu Q^T), \quad \operatorname{Idx}_K = \operatorname{topkIndex}(\mu K^T)$$

1412 Using the triangle inequality of the metric, the above approximation can select Q -related KV tokens with:

$$1413 \operatorname{Idx} = \operatorname{topkIndex}(Q\mu^T\mu K^T)$$

1414 Our semantic tokens in each Attention step also cluster similarly:

$$1415 m_s = \operatorname{Softmax}(m_s K_s^T)K_s$$

1416 or

$$1417 m_s = m_s + \operatorname{softmax}(Q_s K_s^T)V_s$$

1418 Our selection process can be written as:

$$1419 \operatorname{Idx} = \operatorname{topkIndex}(m_s W_Q W_K^T m_s^T)$$

1420 where m_s is the regional semantic token of X_s (e.g., average or $m_s = \operatorname{Softmax}(m_s K_s^T)K_s$), equivalent to further clustering.

1421 Therefore, during selection, the cluster center can be used to represent the tokens within the relevant cluster, similar to the Routing Transformer. Our method, by only using the cluster center for selection, introduces minimal quantization error. However, because the selection quantity is sufficient, the loss is negligible. In contrast, the Routing Transformer loses some related Q and K tokens to ensure regular shape, introducing non-negligible error.

1422 E.1.4 SWIN TRANSFORMER

1423 The Swin Transformer utilizes local Attention and shifted local Attention. Our selection mechanism can completely converge to Swin Transformer when it is optimal.

1424 - Local Attention:

$$1425 \operatorname{softmax}(Q_{i:j} K_{i:j}^T) V_{i:j}$$

1426 When this is the optimal solution, our selection mechanism will automatically select tokens within the local region.

1427 - Shifted Local Attention:

$$1428 \operatorname{softmax}(Q_{i:j} K_{i+r:j+r}^T) V_{i+r:j+r}$$

1429 When this is optimal, our selection mechanism will automatically select tokens within the corresponding region ($KV_{i+r:j+r}$), where r is the cyclic shift offset.

1430 In a similar manner, our method can approximate or cover many variants of the above Transformers.

1431 To compare our method with these efficient Transformers, we conducted experiments on the PG19 datasets, as shown in the tables 17.

1432 On the PG19 dataset, our method outperforms Landmark Attention in terms of both performance and efficiency, improving PPL by 4. The efficiency advantage becomes more pronounced as the input length increases.

Table 17: Perplexity (PPL) and Training Memory Consumption (MB) on PG19. The 20186/22496 means using flashatten or not. The 1×8192 means $bsz \times seqLen$.

Model	PPL (PG19)	1×8192	1×4096	1×2048
MS Attention-60M	13.9	-	-	-
MS Attention-200M	9.32	20186/22496	13152/14300	9976/10346
Landmark Attention-200M	14.55	≈ 40960 (OOM)	≈ 40960 (OOM)	17938

E.2 DETAILED PARAMETER SETTINGS

Our algorithm can encompass the majority of Q-attention to KV-subset methods by adjusting parameters such as the QKV segment size, the number of selected top-K high-relevance KV regions, and the number of merged Q regions. Below is a detailed discussion on the selection of these hyperparameters:

E.2.1 QKV SEGMENT SIZE SELECTION:

The QKV segment size is crucial and typically ranges from 8 to 128. This parameter must be chosen by balancing computational complexity and performance. In future work, we plan to incorporate Triton operators in the QK routing step to achieve linear spatial complexity, thereby mitigating the current limitations:

- (1) During the selection and routing process, a semantic token represents each region, and the relevance between Q and K semantic tokens is measured. KV tokens related to the Q region are then selected based on this relevance. A larger QK region size results in a coarser semantic representation, leading to less precise KV token selection. To minimize this loss, more KV tokens must be selected in the second step, thereby reducing information loss due to coarse semantics. Thus, smaller segment sizes are preferred, though incorporating Triton operators in the QK routing step is anticipated to provide linear spatial complexity.
- (2) Alternatively, a larger segment size can be set, and more KV tokens can be selected in the second step. This approach may increase algorithmic complexity as more KV tokens are likely to be used for interpolation within the same Q region. However, this issue can be alleviated through the third step of merging, where the selected KV tokens can be further merged and selected according to relevance, filtering out irrelevant information.
- (3) Additionally, the segment size of Q regions is independent of the segment size of KV regions. A relationship between them should only be established when specific tasks and complexity constraints require it.

E.2.2 TOP-K HIGH-RELEVANCE KV REGION SELECTION:

The number of selected top-K high-relevance KV regions is generally determined by factors such as the model’s pre-training length and the amount of task-specific data. It may also need to be adjusted based on the granularity of the segmentation from the first step. Current large model training methods suggest that selecting 1K-8K KV tokens for interpolation is robust:

- (1) The above represents a general approach, while scenario-specific selection yields higher performance and efficiency. For example, if the task’s data volume is small, a relatively small number of high-relevance regions can be selected to maximize overfitting and memorize the critical data. Conversely, as the data volume increases, more KV tokens may be required for autoregressive prediction of the next token, necessitating the selection of more high-relevance regions.
- (2) In addition to task-based selection, the granularity of segmentation from the first step must also be considered. If segmentation precision is insufficient, more KV tokens should be selected to prevent the loss of critical information. Subsequent merging can then be used to jointly select high-relevance information, improving space and time complexity by sharing KV tokens across Q regions.

1512 - (3) For length fine-tuning, the model’s parameters have adapted to the interpolation degree of the
 1513 pre-training length, which uses a specific number of KV tokens. Hence, the same magnitude of KV
 1514 tokens must be selected during fine-tuning to avoid overfitting.

1515 This parameter is highly flexible and can be adapted to various scenarios, often requiring some
 1516 experience for optimal settings. Future work may involve incorporating a loss function to control the
 1517 selection space, allowing the algorithm to automatically select the appropriate size based on the loss.

1519 E.2.3 MERGING Q REGIONS:

1521 The number of merged Q regions generally depends on the segmentation precision and space
 1522 complexity. In scenarios where space complexity is not a major concern, sizes ranging from $\frac{topk}{4C_Q}$ to
 1523 $\frac{topk}{C_Q}$ can be chosen. The number of merged KV tokens is typically selected from a range of $topk$ to
 1524 $2topk$ KV tokens. Generally, sharing high-relevance KV regions across multiple Q regions results
 1525 in lower space complexity. Moreover, during the merging process, high-relevance KV regions are
 1526 selected with flexibility, enhancing performance through methods such as:

- 1527 - a. Simple unique and sort operations for selection.
- 1528
- 1529 - b. Unique operations followed by a secondary segmentation based on high-relevance scores, with
 1530 allocation and selection according to region clustering.
- 1531

1532 E.2.4 ALGORITHM COMPLEXITY ANALYSIS

1533 (1) Let the segmentation sizes for the Q and KV regions be C_Q and C_K , respectively. The time and
 1534 space complexity of routing Q and K using dot products is $O\left(\frac{N^2}{C_Q C_K}\right)$.

1535 (2) In the second step, the selection of high-relevance regions, denoted by C_S , results in space
 1536 complexity of $O\left(\frac{N}{C_Q} C_S\right)$ due to the storage of necessary indices.

1537 (3) The merging step, with a merge size of C_M , involves combining the selected indices. Filtering
 1538 algorithms can be introduced during merging; in this work, we employ unique and high-score
 1539 selection. The space complexity for storing KV tokens after selection and merging is $O\left(\frac{N^2}{C_Q C_M} C_S\right)$.

1540 (4) The final attention operation has a time complexity of $O\left(2\frac{N}{C_Q C_M} C_Q C_M C_S\right) = O(NC_S)$. Hence,
 1541 only the segmentation size and selection quantity affect the algorithm’s complexity.

1542 Without merging, the selected KV tokens have a complexity of $O\left(\frac{N^2}{C_Q} C_S\right)$, whereas with merging,
 1543 the complexity is $O\left(\frac{N^2}{C_Q C_M} C_S\right)$. The reduction in complexity due to merging is because the
 1544 selection process is only carried out after merging, storing only the selected indices before that.

1551 E.3 REDUCED LONGLORA

1552 We propose a method that selectively fine-tunes only the W_k and W_o weights, achieving results nearly
 1553 identical to fine-tuning the entire attention mechanism. Specifically, since $QK^T = XW_qW_k^T X^T$,
 1554 updating only W_k yields $W_qW_k^T$, effectively replicating the effect of simultaneously updating W_k
 1555 and W_q . This approach is particularly effective when W_q is full rank. Similarly, fine-tuning W_o
 1556 follows the same rationale, as the final output is a linear mapping of W_vW_o . Since only the Q and K
 1557 tokens undergo positional encoding, the weights to be fine-tuned are limited to W_Q and W_K . Based
 1558 on the above analysis, it suffices to only fine-tune W_Q or W_K . In this work, for the fine-tuning of
 1559 LLaMA, we adopt a low-rank fine-tuning of W_Q and W_K , while for Mistral, we apply low-rank
 1560 learning of W_Q for sequence length extension.

1561 Additionally, since query tokens can be considered well-fitted through extensive training, learning
 1562 the linear mapping for their corresponding key tokens is reasonable. Furthermore, as the attention
 1563 mechanism becomes heterogeneous, updating the final classification head is a direct approach.
 1564 However, considering the large number of parameters, altering the initial embedding An et al. (2022)
 1565 to achieve a certain degree of equivalence is also a feasible consideration.

1566 E.3.1 THE STRATEGY FOR INCORPORATING POSITIONAL INFORMATION DURING
1567 PRE-TRAINING
1568

1569 Our algorithm leverages semantic token routing, where each semantic token encapsulates contextual
1570 information from its surrounding tokens, similar to a sliding window mechanism in convolution.
1571 Since different tokens are surrounded by distinct contexts, the resulting semantic tokens are unique,
1572 leading to different routed KV tokens and interpolated variables. As a result, we do not assign
1573 positional information during the selection process. Instead, after the selection process, RoPE
1574 positional encoding is applied to the K and Q tokens to enhance performance.
1575

1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619