# LM-Infinite: Zero-Shot Extreme Length Generalization for Large Language Models

**Chi Han[1] [\*], Qifan Wang[2], Hao Peng[1], Wenhan Xiong[3], Yu Chen[4] [†], Heng Ji[1], Sinong Wang[3]**

[1] University of Illinois Urbana-Champaign, [2] Meta, [3] GenAI Meta, [4] Anytime AI
[1]{chihan3,haopeng,hengji}@illinois.edu,
[23]{wqfcr,xwhan,sinongwang}@meta.com,
[4]ychen@anytime-ai.com

## Abstract

Today's large language models (LLMs) typically train on short text segments (e.g., <4K tokens) due to the quadratic complexity of their Transformer architectures. As a result, their performance suffers drastically on inputs longer than those encountered during training, substantially limiting their applications in real-world tasks involving long contexts such as encoding scientific articles, code repositories, or long dialogues. Through both theoretical analysis and empirical investigation, this work identifies three major factors contributing to this length generalization failure. Our theoretical analysis reveals that commonly used techniques like using a sliding-window attention pattern or relative positional encodings are inadequate to address them. Answering these challenges, we propose LM-Infinite, a simple and effective method for enhancing LLMs' capabilities of handling long contexts. LM-Infinite is highly flexible and can be used with most modern LLMs off-the-shelf. *Without any parameter updates*, it allows LLMs pre-trained with 2K or 4K-long segments to generalize to up to 200M length inputs while retaining perplexity. It also improves performance on downstream tasks such as Passkey Retrieval and Qasper in the zero-shot setting. LM-Infinite brings substantial efficiency improvements: it achieves $2.7\times$ decoding speed up and $7.5\times$ memory saving over the original model. Our codes are released at https://github.com/Glaciohound/LM-Infinite.

## 1 Introduction

Large language models (LLMs) have recently advanced the state-of-the-art across various natural language processing tasks. They typically train on text segments of fewer than 4K tokens (Touvron et al., 2023b; Team, 2023), primarily due to the computational overhead quadratic in the input

lengths of their Transformer architectures. As a result, they face challenges in generalization to inputs that are excessively longer than what they are trained on and suffer substantial deterioration in their performance (Tworkowski et al., 2023; Chen et al., 2023a). This limits their applicability in tasks that require long-range contexts, such as encoding scientific articles, source code repository generation, or long-context dialogues.

Extensive efforts have been devoted to addressing this length generalization challenge. Relative positional encodings such as RoPE (Su et al., 2021) and Alibi (Press et al., 2021) have been widely adopted by state-of-the-art LLMs, which calculate attention based on inter-token distance instead of absolute positions, hoping to avoid model failures due to unseen absolute position embeddings. Moreover, although applying a sliding-window attention pattern on the Transformer architecture can reduce the memory overhead (Beltagy et al., 2020; Ding et al., 2023; Zaheer et al., 2020), they are not directly applicable to pre-trained models for length generalization without further training. Through both theoretical analysis and empirical investigation, §3 pinpoints three primary factors underlying the length generalization failures: (1) the challenge of handling unseen distances among tokens, (2) the difficulty of attending to unseen numbers of tokens, and (3) implicitly encoded absolute positional information in initial tokens. These challenges can make LLMs' computational features, such as attention logits and hidden vectors, deviate from the training distribution, leading to failures of length generalization. Existing techniques fall short of addressing these underlying issues.

Answering these challenges, we propose LM-Infinite, a simple and effective method to enhance Transformer LLMs' capabilities for modeling long contexts *without parameter updates*. LM-Infinite consists of two major components designed to alleviate the three factors above. (1) a $\Lambda$-shaped

---

attention mask and (2) a ceiling on attention distances. The former forces the model to attend to only the beginning of the sequence and the most recent tokens within a pre-defined window, ignoring the rest. The latter component caps the relative distance values to the maximum the model has seen during training. It can also optionally re-introduce top-$k$ tokens in the middle to achieve better performance in some downstream tasks. LM-Infinite is highly flexible and applies to any off-the-shelf LLMs that use relative positional encoding and does *not* require any finetuning.

Our experiments thoroughly evaluate LM-Infinite on a variety of tasks and LLMs. On ArXiv (academic papers) and OpenWebText2 (Reddit posts) LM-Infinite facilitates *zero-shot* generalization for a wide range of LLMs to texts up to 200M tokens, retaining the language modeling perplexity and generation quality. *Without any parameter updates*, LM-Infinite improves scores compared with the original model and truncation baselines on downstream tasks including Passkey Retrieval (Mohtashami and Jaggi, 2023) and Qasper (Dasigi et al., 2021), which are two established benchmarks for long-context evaluation. We observe a 37.2% gain on Passkey Retrieval and a 1.2% gain on Qasper in the zero-shot setting. LM-Infinite also brings substantial efficiency improvements: it achieves $2.7\times$ decoding speed up and $7.5\times$ GPU memory saving over the original LLMs.

## 2 Background and Related Work

### 2.1 Relative Positional Encodings

The traditional absolute positional encodings provide the absolute position information, usually with the help of a sequence of vectors called *position embeddings* (Vaswani et al., 2017; Kenton and Toutanova, 2019; Ke et al., 2020). They, however, have trouble when the model encounters unseen positions in inputs longer than the training length. Relative positional encodings aim to address the limitations of previous-generation positional encoding methods and consider the relative distances between tokens instead of the absolute positions. Examples include a learned attention logit bias in T5 (Raffel et al., 2020), Transformer-XL (Dai et al., 2019), Skyformer (Chen et al., 2021), Sketching (Chen et al., 2022) and Sandwich (Chi et al., 2023), a fixed linear attention decay (Press et al., 2021), and rotating query and key sequences based on distances such as RoPE (Su et al., 2021; Li et al.,

2023), CAPE (Likhomanenko et al., 2021) and XPos (Sun et al., 2022; Ding et al., 2023). Despite some promising empirical evidence, length generalization failures are still widely observed when directly applied to large language models (Kaiokendev, 2023). In what follows, we briefly discuss two widely used relative positional encoding methods. They lay out the necessary context for our onward discussion and experiments.

**Rotary Position Embedding (RoPE; Su et al., 2021)** It rotates the key and query vectors based on positions before computing the inner product. Specifically, each vector $\mathbf{x}$ (either key or query) is split into pairs of elements $\{(x_0, x_1), (x_2, x_3), \cdots\}$, with each pair interpreted as a 2-dimensional vector. RoPE then rotates the vector $(x_a, x_{a+1})$ of token $i$ with angle $\theta_{a,i} = i\omega_a$, where $\omega_a$ is the rotating speed associated with dimension pair $(a, a+1)$. After rotation, the 2-D vector becomes $\begin{pmatrix} \cos i\omega_a & -\sin i\omega_a \\ \sin i\omega_a & \cos i\omega_a \end{pmatrix} \begin{pmatrix} x_i \\ x_{i+1} \end{pmatrix}$. They show that the inner product between rotated query $\mathbf{q}_i$ and rotated key $\mathbf{k}_j$ is solely determined by $\mathbf{q}_i, \mathbf{k}_j$, and their relative distance $i - j$. We always have $i \geq j$ due to the causal attention mask.

**AliBi (Press et al., 2021)** It offsets all attention logits between tokens $i, j$ by a linear term $-m(i - j)$ and become $\mathbf{q}_i^\top \mathbf{k}_j - m(i - j)$. To this end, the MPT-7B codes implement an offset matrix as an additive term in attention logits.

### 2.2 Efforts Towards Length Generalization

In light of generalization failures observed in LLMs, one straightforward solution is to finetune LLMs on longer text sequences (Chen et al., 2023a; Tworkowski et al., 2023; Tao et al., 2023; Kiyono et al., 2021; Anil et al., 2022). These approaches do not address the underlying causes of length generalization failures and require massive training resources. Other solutions propose to grant LLMs access to longer contexts without really reading them in full (Zhou et al., 2023; Bueno et al., 2022; Mohtashami and Jaggi, 2023; Yang et al., 2023). Augmenting LLMs with retrieval-based memories (Wu et al., 2021; Guu et al., 2020; Borgeaud et al., 2022; Khandelwal et al., 2019; Kaiser et al., 2016; Yogatama et al., 2021) also make LLMs applicable to a large database. These designs, however, usually need finetuning and are not directly compatible with the existing LLMs. Our work, in contrast, facilitates zero-shot length generalization. Another

(a) Attention Logit Explosion at Long Distances   (b) Attention Entropy Explosion at Long Lengths   (c) Starting Tokens Occupy Distinct Areas
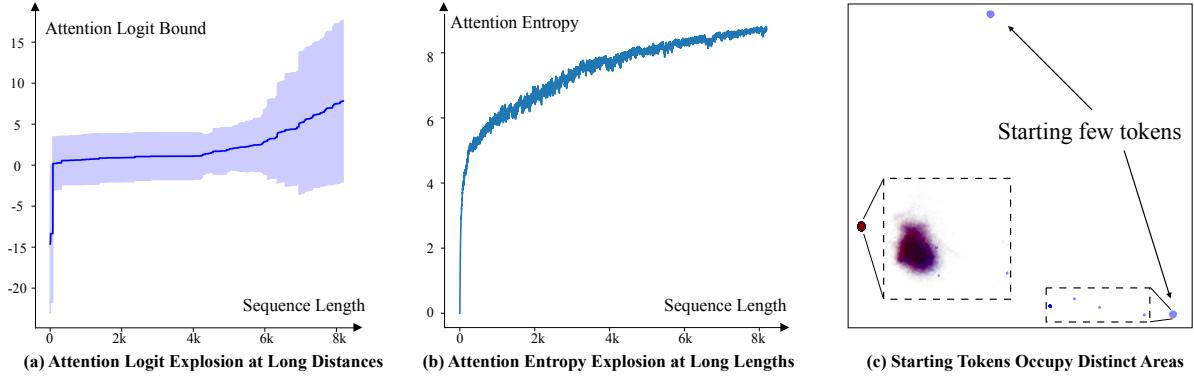
Figure 1: We identify three factors underlying the length generalization failure in LLMs in §3. (a) Factor 1: Unseen distances between tokens cause attention logits to explode. (b) Factor 2: An unseen number of tokens can cause attention entropy to increase beyond the training range as the length increases. (c) Factor 3: Starting few tokens occupy a distinct feature region and should not be discarded. The two blue regions at the upper center and lower right correspond to the initial tokens that are highly concentrated but also very far from later tokens. The lower-left region contains the thousands of overlapping dots corresponding to the later tokens.

similar work (Ratner et al., 2023) increases context length with attention patterns without further training. However, it is limited to the in-context learning setting.

## 3 Why do Transformer LLMs Fail to Generalize to Long Contexts?

Through a series of theoretical and experimental investigations, this section aims to identify the potential causes underlying current LLMs' failure in length generalization. Our discussion assumes Transformer-based LLMs that use relative positional encodings, as this design is widely adopted in today's LLMs (Touvron et al., 2023b; Team, 2023). We use Llama-2 (Touvron et al., 2023b), which is pre-trained with 4K-length segments, for investigation. On sequences longer than the training length, we will show that the *unseen inter-token distances*, the *increasing number of attended tokens*, and the *implicitly encoded position information of the starting tokens* can all make certain computational features out of the training distribution. As deep models can be sensitive to input distribution shifts, these factors need to be handled for LLMs to generalize to unseen lengths.

**Factor 1: challenges in handling unseen distances among tokens** With relative positional encoding, the impact of positions on the attention weight between two tokens depends solely on their relative distance. As the sequence length grows exceedingly long, some distance values will surpass those seen during training. We make the following informal theoretical claim:

**Theorem 1.** *(Informal) For an attention mechanism using relative positional encoding, the attention logits must explode to infinities to differentiate previously unseen distances apart as the sequence length increases.*

The formal theorem and its proof can be found in Appendix C. We also empirically verify this on Llama-2 on the ArXiv dataset truncated down to 8K length. We extract the attention logits of all attention heads and their maximum attention logits on different sequence lengths in Figure 1(a). It shows the average and variance among attention heads. We see that the attention logits increase to substantially larger values when the sequence length exceeds the training length of 4K. To mitigate this issue, we conjecture that **it may help to cap the relative distance values to the maximum that the model has seen during training (i.e., a distance ceiling)**. However, as we will see from the proposition below, addressing logit explosion leads to another challenge.

**Factor 2: attending to unseen numbers of tokens** On longer sequences, tokens at later positions must distribute attention weights across a larger context. We then make the following claim that, if attention logits are bounded but the number of tokens to attend to is not limited, it can cause the attention entropy to increase beyond the training range:

**Proposition 1.** *If the attention logits are bounded, as the sequence becomes longer, the attention entropy grows to infinity.*

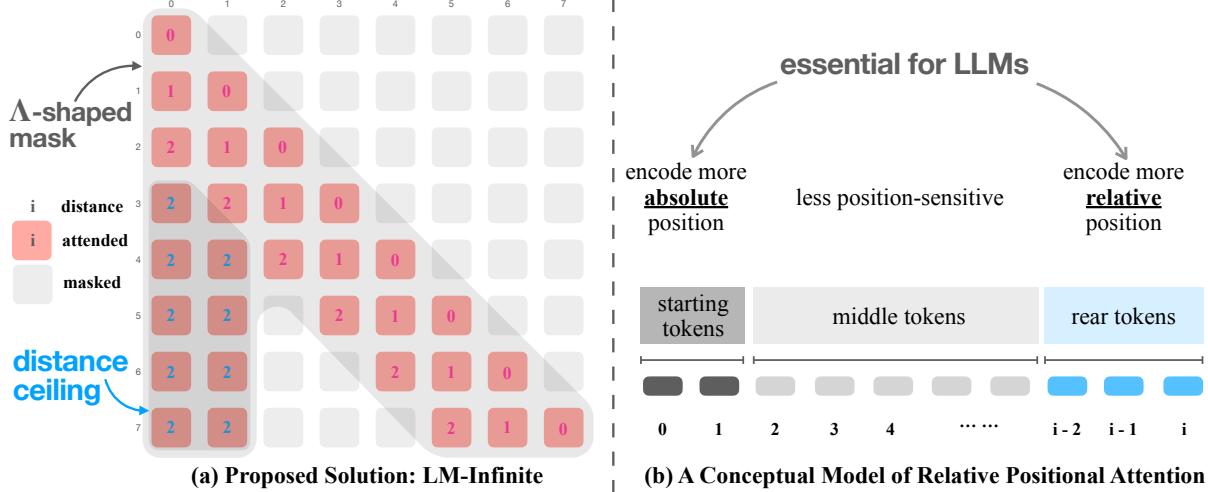A formal statement as well as the proof can be

(a) Proposed Solution: LM-Infinite

(b) A Conceptual Model of Relative Positional Attention

Figure 2: (a) LM-Infinite is a plug-and-play solution for various LLMs, consisting of a $\Lambda$-shaped mask and a distance ceiling during attention. For clarity, this figure shows a toy scenario where $L_{\text{pre-train}}$ and $n_{\text{starting}}$ are both 2. (b) We also provide a conceptual model for understanding how relative position encoding works.

found in Appendix D. This conclusion is further empirically verified by plotting the attention entropy against context lengths in Figure 1(b). The curve shows an ever-increasing attention entropy. This trend, although increasing logarithmically, is still harming the LLMs' performance, as we will illustrate in the ablation study in §5.2 and Figure 5. This suggests that **we should bound the attention context size** to ensure that the attention entropy stays in seen ranges during pre-training and avoid degenerated outputs. A simple windowed attention, where each token only attends to the nearest tokens within a distance, might handle factors 1 and 2. This is similar to the block-diagonal attention mask used in XPos (Sun et al., 2022) and Long-former (Beltagy et al., 2020). However, as we will show in the next paragraph, this introduces another factor that can also fail LLMs.

**Factor 3: starting tokens occupy a distinct feature space** Perhaps counter-intuitively:

**Observation 1.** *Even without explicit absolute positional embeddings, attention outputs of the first few tokens can occupy a distinct representational space compared to other positions. Therefore, when passed to later layers, these starting tokens have distinct value vectors from their lower-layer outputs.*

This follows from Theorem 1 in Kazemnejad et al. (2023), which proves that the absolute positions can be implicitly encoded in the outputs of tokens of a single attention layer, even *without* positional encodings. In their construction, the starting tokens'

signals are the strongest and easiest to distinguish from other tokens. As relative positional encoding is strictly more expressive than no positional encoding setting in Kazemnejad et al. (2023) (e.g., by letting all distances have the same attention function), the same conclusion applies to relative positional encoding as well.

As an empirical verification, we take the hidden states output by the second layer of Llama-2 and plot a Principal Component Analysis (PCA) projection into a 2-d plane in Figure 1(c). More figures for other layers can be found in §E. The dots correspond to the first 4096 tokens in 32 sequences, with blue ones corresponding to the initial tokens and red tokens being the tail ones. The two blue regions at the upper center and lower right correspond to the initial highly concentrated tokens (whose positions are around 0~25) and are very far from later tokens. The lower-left region contains the remaining overlapping tokens in a sequence (zoomed in to another box). The plot shows that the vector representations of the initial tokens concentrate on regions in the feature space that are distinct from the remaining tokens. This fresh finding reveals a fundamental flaw of the sliding-window attention pattern, which restricts the attention to the most recent tokens within a predefined window size, a widely adopted baseline recently (Beltagy et al., 2020; Ding et al., 2023; Zaheer et al., 2020). As attention is essentially a weighted average over the value vectors, sliding-window attention discards the initial tokens, keeping the attention output from reaching the regions that value vectors of the initial

tokens occupy. This enforces the model to handle a different region during the computation, introducing additional generalization challenges. As a straightforward solution to this issue, **the initial tokens need to be kept in the attention computation.**

## 4 Our proposal: LM-Infinite

Inspired by the analyses and take-away messages in the previous section, we propose LM-Infinite to achieve zero-shot length generalization for LLMs. An overview of LM-Infinite is illustrated in Figure 2(a). This simple solution consists of two components: a $\Lambda$-shaped attention mask and a distance ceiling. Besides, re-introducing the middle top-$k$ tokens is optional for enhanced downstream performance.

**$\Lambda$-shaped attention mask**   It contains two attention spans: the *starting* one allows each token to attend to the first $n_{\text{starting}}$ tokens if they come before the current one; the ending one allows each token to attend to most recent $L_{\text{pretrain}}$ tokens. $L_{\text{pretrain}}$ is the maximum length during training. Other tokens are ignored. In ablation studies in §A, we find that choosing $n_{\text{starting}} \in [5, 100]$ generally achieves equally good performance. Note that $n_{\text{starting}} = 0$ (i.e., attending only to the most recent tokens) substantially hurts the performance. This resolves Factors 2 and 3 in §3 by both limiting the number of tokens under attention and ensuring the starting few tokens are attended.

**Distance ceiling**   LM-Infinite further bounds the "effective distance" to $L_{\text{pretrain}}$. This only affects the starting few tokens when attended by tokens at later positions. Specifically, in relative positional encoding, the original attention logit is $w(\mathbf{q}, \mathbf{k}, d)$, where $d$ is the distance between two tokens. Now we modify it as $w(\mathbf{q}, \mathbf{k}, d')$ where $d' = \min(d, L_{\text{pretrain}})$. Figure 2(a) shows an illustrative example where the distance ceiling is $L_{\text{pretrain}} = 2$. This addresses Factor 1 in §3 by bounding the distance value in attention calculation.

**Optionally attending to top-$k$ tokens in the middle**   LM-Infinite can optionally attend to $k$ tokens in the middle with the largest attention logits. This is particularly useful in downstream tasks where information in the middle tokens matters (§5.3). Here the $k$ tokens are selected independently for each attention head in layers higher than $h$-th layer, and have an attention distance of $d = \frac{1}{2}L_{\text{pre-train}}$.

These hyperparameters are selected based on a held-out Passkey Retrieval validation set, where we set $k = 5$ and $h = 5$, with more details in Appendix A. Our selection of $k$ and $h$ does not depend on task-specific tuning, and in our experiments, we apply this same set of hyperparameters in all other downstream tasks and achieve consistent improvements. These intermediate tokens do not hurt performance. Rather, in the evaluation of downstream tasks in §5.3, intermediate tokens are more useful and selectively attending to top-$k$ tokens brings substantial performance improvements with little impact on the efficiency. For LLM generation and inference, however, we find the intermediate tokens *unnecessary* to attend to for LM-Infinite to achieve good perplexity or generation quality.

LM-Infinite's $\Lambda$-shaped mask is conceptually similar to the attention patterns derived from heuristics (Beltagy et al., 2020; Ding et al., 2023; Zaheer et al., 2020). However, we formally show in §3 Factor 3 that these previous approaches theoretically cannot generalize to unseen lengths but require parameter updates. This inherent limitation motivates the other two components in LM-Infinite to achieve zero-shot length generalization.

**Implementation details**   LM-Infinite is applicable in all Transformer models with relative positional encoding. One only needs to replace the attention function in each Transformer layer with LM-Infinite without any parameter updates. The $\Lambda$-shaped attention mask is relatively straightforward to implement. In RoPE, attention logits in the ending attention span follow the original calculation. In the starting attention span (excluding its overlap with the ending span), we keep all $\mathbf{k}$ vectors unrotated and rotate all $\mathbf{q}$ vectors to a fixed distance $L_{\text{pretrain}}$. Then the logits in two spans can be composed. Augmenting AliBi with LM-Infinite is also straightforward. We simply clip the offset matrix with a minimum value of $-|mL_{\text{pretrain}}|$ and apply the $\Lambda$-shaped attention mask.

**Discussion.**   In Figure 2(b), we show a conceptual model of how relative positional encoding functions. This conceptual model reflects the design choices of LM-Infinite. In this conceptual model, a long context can be roughly partitioned into 3 parts: The *starting tokens* encode strong absolute position information (Factor 3). As explained in §3, they are essential to attention to because their features occupy a distinct region in the feature space. As attention is essentially a weighted average over

Negative Log-Likelihood

LLaMA

Llama-2

MPT-7B

GPT-J-6B

MPT-7B + LM-Infinite
MPT-7B-Storywriter
Llama-2 + LM-Infinite
LLaMA + LM-Infinite
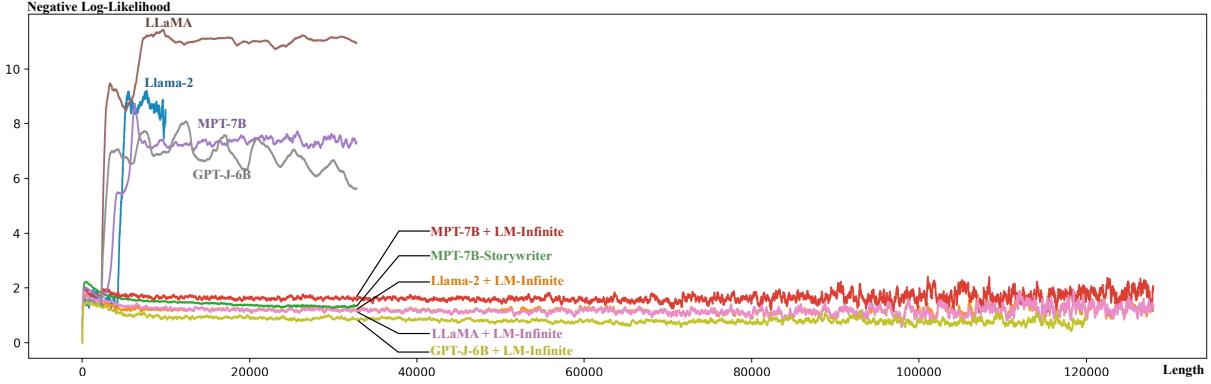GPT-J-6B + LM-Infinite

Length

Figure 3: LM-Infinite flattens the negative log-likelihood (NLL) curves of various LLMs on ArXiv dataset without any parameter updates. The trends are similar to MPT-7B-Storywriter, an explicitly fine-tuned LLM. Llama-2 outputs NaN values on long sequences so the curve is relatively shorter.

| Model | $L_{\text{pretrain}}$ | ArXiv | | | | | OpenWebText2 | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 2K | 4K | 8K | 16K | 32K | 2K | 4K | 8K | 16K |
| *Long-context Training/Finetuning* | | | | | | | | | | |
| Sandwich | 512 | 5.0 | 5.2 | 5.3 | - | - | 23.3 | 23.8 | 24.7 | - |
| XPos | 1K | 21.6 | 20.7 | - | - | - | - | - | - | - |
| LongLLaMA | 8K | 8.2 | 7.4 | - | 6.9 | - | - | - | - | - |
| MPT-7B-SW | 65K | 6.5 | 5.4 | 4.3 | 4.4 | 3.6 | 9.8 | 10.9 | 6.6 | **5.1** |
| *Vanilla* | | | | | | | | | | |
| MPT-7B | 4K | 5.5 | $2.5\times10^2$ | $1.1\times10^3$ | $1.7\times10^3$ | $1.6\times10^3$ | 8.3 | $1.3\times10^2$ | $1.9\times10^2$ | $1.3\times10^2$ |
| LLaMA | 2K | 3.8 | $1.0\times10^4$ | $6.0\times10^4$ | $6.8\times10^4$ | $4.9\times10^4$ | 6.2 | $6.6\times10^3$ | $4.6\times10^5$ | $4.4\times10^4$ |
| GPT-J-6B | 2K | 3.9 | $1.3\times10^3$ | $1.0\times10^3$ | $1.6\times10^3$ | $2.8\times10^2$ | 8.8 | $7.5\times10^2$ | $1.3\times10^3$ | $1.8\times10^3$ |
| Llama-2 | 4K | **3.4** | 3.8 | $8.5\times10^3$ | NaN | NaN | 6.2 | 5.8 | $6.5\times10^3$ | NaN |
| *LM-Infinite* | | | | | | | | | | |
| MPT-7B | 4K | 5.7 | 6.8 | 5.8 | 6.0 | 4.6 | 8.5 | 12.2 | 8.5 | 8.9 |
| LLaMA | 2K | 4.4 | 4.5 | 3.7 | 4.2 | **1.0** | 6.3 | 6.1 | 9.5 | 7.0 |
| GPT-J-6B | 2K | 3.8 | **3.1** | **3.0** | **3.1** | 2.1 | 8.8 | 8.5 | **6.5** | 7.4 |
| Llama-2 | 4K | 4.3 | 3.6 | 3.3 | 4.2 | 6.5 | **6.1** | **5.3** | 8.3 | 8.2 |

Table 1: Perplexity on ArXiv and OpenWebText2 test split. LLMs with LM-Infinite achieve the highest perplexity on 7 out of 9 columns while requiring no parameter updates. $L_{\text{pretrain}}$ indicates the lengths of the text segments that the models are trained on.

$\mathbf{v}_i$ vectors, without the starting few tokens, the attention output can not reach regions that $\mathbf{v}_i$ vectors of the initial tokens occupy. The *rear tokens* provides primarily their relative positions to the final tokens. Their importance probably arises from the "recency bias" (Peysakhovich and Lerer, 2023) learned by LLMs during pre-training. The *middle tokens* encode less position-sensitive information. As analyzed in Factor 2, including too many intermediate tokens does more harm than good to length generalization.

## 5 Evaluation

We evaluate LM-Infinite with LLaMA-7B (Touvron et al., 2023a), Llama-2-7b (Touvron et al., 2023b), MPT-7B (Team, 2023), and GPT-J-

6B (Wang and Komatsuzaki, 2021). LLaMA-7B and GPT-J-6B are pre-trained with 2K lengths and the other models are pre-trained with 4K lengths. LLaMA, Llama-2, and GPT-J use RoPE encoding, and MPT-7B uses Alibi encoding. MPT-7B-Storywriter (fine-tuned on long sequences) is used as one of the baselines.

### 5.1 Language Modeling with Extremely Long Context

We use ArXiv and OpenWebText2 corpora from the Pile dataset (Gao et al., 2020), which contain preprint papers from ArXiv and Reddit submissions, respectively. We evaluate with negative log-likelihood (NLL) and perplexity (exponential of NLL). Figure 3 plots the NLL curves on the ArXiv
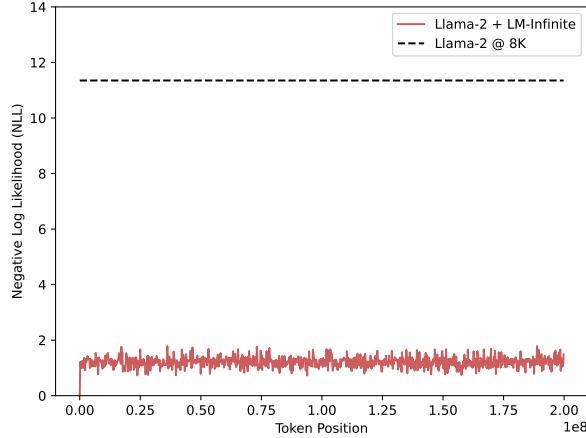
Figure 4: LM-Infinite generalizes Llama-2 to an extreme length of 200M. The dashed line is the NLL at 8K length of the vanilla Llama-2 model.
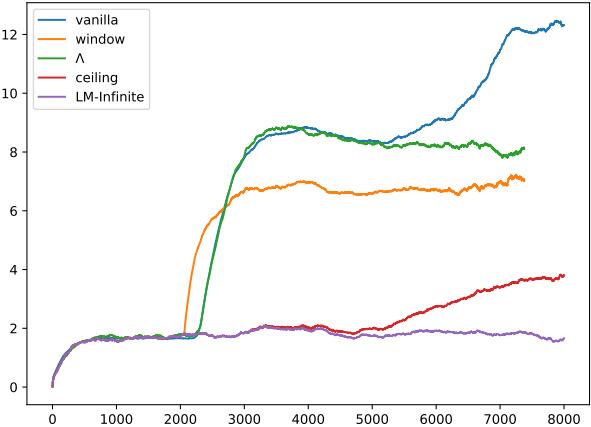


Figure 5: Ablation study on LLaMA in §5.2. x-axis is token position and y-axis is negative log-likelihood (NLL). The vanilla model (vanilla), using a windowed attention (window), using only a Λ-shaped attention mask (Λ), and only the ceiling the distance value (ceiling) all more or less suffer from perplexity explosion. Only LM-Infinite can retrain the performance while generalizing to unseen lengths.

dataset. Here, we break down the models' perplexity performance by positions so that the curve shows the NLL that the model achieves around that specific position, averaged across all evaluated sequences. Llama-2 outputs NaN probabilities on sequences that are slightly longer than 10K, thereby its shorter curve. All vanilla models run out of memory at ∼32K lengths.[1] The baselines' NLL quickly blows up when the tested sequences are longer than what they train on. With LM-Infinite, all models can generalize to sequences that are substantially longer than the lengths they are trained on, retaining the NLL performance. This validates our length failure factor analyses in §1. The longer ends of curves have larger variances because of fewer documents of those lengths. In Figure 4, we further evaluate LM-Infinite + Llama-2 on a sequence of **200M** tokens, which is constructed by sampling with replacement from the ArXiv dataset and concatenating all data. LM-Infinite shows the ability to remain stably low log-perplexity level over extreme lengths.

Table 1 summarizes the perplexity performance at a few milestone lengths (2K, 4K, 8K, 16K, and 32K) on ArXiv and OpenWebText2, which shows a similar trend. OpenWebText2 has very few data instances over a length of 32K, so we omit the column. With LM-Infinite, all models can generalize to unseen lengths, and LM-Infinite achieves best perplexity in 7 out of 9 cases. On LLaMA + LM-Infinite, the perplexity decreases as length increases and position becomes larger. Surprisingly, *without*

*any parameter update*, LM-Infinite outperforms many strong baselines that are trained on substantially longer text segments. As a direct comparison, MPT-7B+LM-Infinite achieves only slightly worse performance than its fine-tuned counterpart, MPT-7B-Storywriter. This confirms that LM-Infinite is a promising alternative to resource-consuming fine-tuning.

## 5.2 Ablation Study

Figure 5 provides an ablation study with the LLaMA model on the ArXiv dataset about why both components in LM-Infinite are essential for maintaining LLM functionality over the length of 8K. We compare LM-Infinite with its variants to show the efficacy of the design and also to validate the factors in §3. Among all curves, only LM-Infinite has relatively stable log-perplexity, meaning that components in LM-Infinite are all essential for successful length generalization. The vanilla LLM model (the "vanilla" curve) fails immediately with exploding NLL. If we only apply Λ-shaped mask (the "Λ" curve) and do not bound inter-token distance (Factor 1), the NLL still explodes immediately after pre-training lengths. The "ceiling" curve only applies the distance ceiling technique but not the Λ-shaped mask to limit the number of attended tokens. The performance still degenerates (evidenced by an ever-increasing NLL). This confirms that the existence of Factor 2, too many to-

---

[1] We run on a single A100 GPU with 80GB GPU memory.

| Model | Passkey Retrieval | | | | | | Qasper |
|---|---|---|---|---|---|---|---|
| | 6K | 8K | 10K | 12K | 16K | average | |
| **Original** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.2 |
| **Truncated** | 66.0 | 55.3 | 38.8 | 32.8 | 27.3 | 44.0 | 30.1 |
| **LM-Infinite** | **70.3** | **90.8** | **86.5** | **79.3** | **79.1** | **81.2** | **31.3** |

Table 2: Downstream evaluation on Passkey Retrieval and Qasper. LM-Infinite enables Llama-2 to consistently outperform both the original model and the baseline that truncates inputs to 4K. The truncation baseline drops excessive tokens altogether when the context is longer than the model's pretraining length, keeping only the most recent ones, which happens before the forward pass starts without changing the attention mechanism.

| Model | $L_{\mathrm{pretrain}}$ | BLEU | | | | | ROUGE | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 2K | 4K | 8K | 16K | 32K | 2K | 4K | 8K | 16K | 32K |
| *ArXiv* | | | | | | | | | | | |
| MPT-7B | 4K | 0.0 | 0.2 | 0.0 | 0.0 | 0.4 | 5.6 | 3.6 | 5.9 | 1.7 | 1.4 |
| MPT-7B-SW | 65K | 16.6 | 21.5 | 15.2 | 18.9 | 14.8 | 26.5 | 30.1 | 26.6 | 27.4 | **27.0** |
| MPT-7B + LM-Infinite | 4K | 16.1 | 20.2 | 12.6 | 13.9 | **19.7** | 23.8 | 24.9 | 24.1 | 29.0 | 26.6 |
| Llama-2 | 4K | 26.6 | 0.0 | 0.0 | 0.0 | 0.0 | 31.4 | 0.2 | 0.0 | 0.0 | 0.0 |
| Llama-2 + LM-Infinite | 4K | **26.9** | **23.6** | **23.9** | **24.8** | 18.4 | **31.8** | **30.9** | **28.8** | **29.2** | 20.4 |
| *OpenWebText2* | | | | | | | | | | | |
| MPT-7B | 4K | 0.9 | 0.9 | 1.0 | 1.0 | - | 7.5 | 6.6 | 6.4 | 6.8 | - |
| MPT-7B-SW | 65K | 8.4 | 6.1 | 7.5 | 8.4 | - | 21.0 | 19.3 | 18.5 | **22.0** | - |
| MPT-7B + LM-Infinite | 4K | 5.0 | 4.1 | 5.1 | 2.8 | - | 16.6 | 15.4 | 16.2 | 16.0 | - |
| Llama-2 | 4K | 8.8 | 0.0 | 0.0 | 0.0 | - | **22.4** | 0.2 | 0.0 | 0.0 | - |
| Llama-2 + LM-Infinite | 4K | **9.0** | **7.2** | **9.7** | **9.6** | - | 21.9 | **21.2** | **19.6** | 19.6 | - |

Table 3: Text generation quality on ArXiv and OpenWebText2. LM-Infinite consistently generalizes the generation quality to extreme lengths, achieving performance that is comparable to or better than the fine-tuned LLM, MPT-7B-Storywriter. Some 0 BLEU scores are caused by the poor generation quality from the vanilla LLMs as they generate mostly nonsensical texts.

kens, is still harming the LLMs' performance. The "window" curve shows a baseline with the sliding-window attention pattern, which only attends to the most recent tokens in a sliding window without altering the input text. It produces the second worst NLL values, which indicates a significant performance and fluency degradation. This confirms our theoretical analysis of factor 3. Due to its visibly much worse performance, we exclude it from other evaluations.

Another similar baseline to "window" is the truncation baseline, which drops excessive tokens altogether when the context is longer than the model's pre-training length, keeping only the most recent ones. This truncation process happens before the forward pass starts and essentially removes the truncated text from the input to the model without changing the attention mechanism. We compared this baseline in two places in the paper. In §5.3 and Table 2, LM-infinite outperforms this baseline on downstream tasks. In Section 5.4 and Figure 6, LM-

infinite achieves a better trade-off between computation complexity and generation quality than this baseline.

## 5.3 Downstream Evaluation

As LLMs are often deployed for downstream tasks, we evaluate how LM-Infinite performs on two long-input tasks under the zero-shot setting: Passkey Retrieval (Mohtashami and Jaggi, 2023) and Qapser (Dasigi et al., 2021). Passkey Retrieval buries a passkey at a random position in a long distraction text and, in the end, asks what the passkey is. Qasper is a question-answering dataset on scientific papers with a total of 1.5K testing question-answer pairs. We evaluate Llama-2-7b-chat, as its instruction tuning enables good task-solving ability (Bai et al., 2023), with middle top-5 tokens enabled on higher than 5-th layer (see §4 for definition and Appendix A for hyperparameter selection). Results are listed in Table 2. LM-Infinite consistently outperforms the baselines on both tasks, with
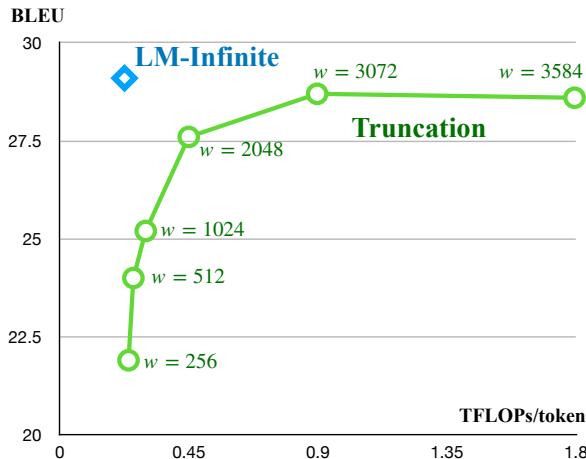
Figure 6: LM-Infinite achieves a better trade-off between computation complexity with generation quality than simple truncation.

a 37.2 percentage gain on Passkey Retrieval and a 1.2 percentage gain on the Qasper task. Passkey retrieval locates useful information uniformly in a sequence, so the performance of the truncated baseline largely depends on whether the remaining part covers the passkey. On Qasper, the top-$k$ attention is necessary for achieving good performance, which indicates that similarly important information in the middle needs to be attended to. This suggests that it can improve downstream task performance on long inputs without fine-tuning while the vanilla model immediately fails.

### 5.4 Generation Quality

We further evaluate LM-Infinite's generation quality on ArXiv and OpenWebText2 test sets, with BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004) (ROUGE-L). We let the LLMs generate 100 tokens after each milestone length and use the following 100 tokens in original texts as references. As the generation is time-consuming, we sample 100 long sequences for evaluation for each dataset. The results are summarized in Table 3. The trend is similar to that in the last section: *without* parameter updates, LM-Infinite successfully allows LLMs to retain their generation quality while generating sequences longer than training, comparable to the fine-tuned baselines such as MPT-7B-SW. The generation results from the vanilla LLMs are poor and contain mostly nonsensical texts, resulting in many close-to-zero scores. For some BLEU scores, it yields zero {2,3,4}-gram overlaps with the reference texts. As BLEU is weighted geometric mean over {1,2,3,4}-gram precisions, the

final BLEU scores for those columns are 0. Appendix Table 8 presents some generation output examples that can provide a good picture of the generation quality. We also evaluate the efficiency in Appendix G: with 32K-long sequences, LM-Infinite achieves 2.7× decoding speedup and 7.5× GPU memory saving.

A few example generations are shown in Appendix H. We also compare LM-Infinite with a simple truncation-based baseline by repeatedly truncating excessive contexts. However, as the generation length increases, frequent truncations and re-encoding of new contexts are required. The larger the truncation window is, the more context is kept, but the larger the computational overhead. We let the models generate 10k tokens on ArXiv. In Figure 6, it is clear that LM-Infinite achieves a substantially better quality-efficiency tradeoff. With similar computation, LM-Infinite outperforms the baseline by about 5 BLEU. To achieve a similar BLEU, LM-Infinite incurs only <25% computational overhead than the truncation baseline.

### 6 Conclusions and Future Work

This work proposes a zero-shot length generalization method for various off-the-shelf LLMs without parameter updates. Through theoretical analysis and empirical investigation, this work identifies three major factors contributing to this length generalization failure. Our theoretical analysis further reveals why truncating the attention window and relative positional encodings are inadequate to address them. Our solution, LM-Infinite, is a simple and effective method for enhancing LLMs' capabilities of handling long contexts. It allows LLMs pre-trained with 2K or 4K-long segments to generalize to up to 200M length inputs while retaining perplexity. It also improves performance on downstream tasks such as Passkey Retrieval and Qasper in the zero-shot setting. It brings substantial efficiency improvements: a 2.7× decoding speed up and a 7.5× memory saving over the original model. LM-Infinite's computational efficiency and ease of use allow researchers without enormous computational resources to use LLMs on long sequences. Future work can investigate if these techniques allow for more efficient and effective LLM pre-training and fine-tuning. Another direction is to apply LM-Infinite to applications such as long reasoning, long-dialogue, retrieval-augmented generation, or long literature generation.

## Limitations

This work evaluates a wide range of open-domain LLMs. However, without access to the source code of proprietary LLMs such as ChatGPT, the proposed method could not be evaluated on them. Furthermore, due to limited computational resources and time, the proposed method has not been evaluated on texts with even larger lengths, such as 1G. The model is designed on relative positional encoding Transformer models, which is the mainstream backbone for most modern LLMs. The question of how LM-Infinite can enable more efficient fine-tuning or pre-training can also be explored in future work.

## Acknowledgement

## References

Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. 2022. Exploring length generalization in large language models. *Advances in Neural Information Processing Systems*, 35:38546–38556.

Yushi Bai, Xin Lv, Jiajie Zhang, Yuze He, Ji Qi, Lei Hou, Jie Tang, Yuxiao Dong, and Juanzi Li. 2024. Longalign: A recipe for long context alignment of large language models. *arXiv preprint arXiv:2401.18058*.

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2023. Longbench: A bilingual, multi-task benchmark for long context understanding.

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR.

Mirelle Candida Bueno, Carlos Gemmell, Jeff Dalton, Roberto Lotufo, and Rodrigo Nogueira. 2022. Induced natural language rationales and interleaved markup tokens enable extrapolation in large language models. In *Proceedings of the 1st Workshop on Mathematical Natural Language Processing (MathNLP)*, pages 17–24.

Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. 2023a. Extending context window of large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*.

Yifan Chen, Qi Zeng, Heng Ji, and Yun Yang. 2021. Skyformer: Remodel self-attention with gaussian kernel and nyström method. In *Proc. Thirty-fifth Annual Conference on Neural Information Processing Systems (NeurIPS2021)*.

Yifan Chen, Qi Zeng, Heng Ji, and Yun Yang. 2022. Sketching as a tool for understanding and accelerating self-attention for long sequences. In *Proc. The 2022 Conference of the North American Chapter of the Association for Computational Linguistics - Human Language Technologies (NAACL-HLT2022)*.

Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2023b. Longlora: Efficient fine-tuning of long-context large language models. In *The Twelfth International Conference on Learning Representations*.

Ta-Chung Chi, Ting-Han Fan, Alexander Rudnicky, and Peter Ramadge. 2023. Dissecting transformer length extrapolation via the lens of receptive field analysis. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13522–13537.

Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988.

Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A Smith, and Matt Gardner. 2021. A dataset of information-seeking questions and answers anchored in research papers. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4599–4610.

Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang, Shaohan Huang, Wenhui Wang, and Furu Wei. 2023. Longnet: Scaling transformers to 1,000,000,000 tokens. *arXiv preprint arXiv:2307.02486*.

Yiran Ding, Li Lyna Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan Yang,

and Mao Yang. 2024. Longrope: Extending llm context window beyond 2 million tokens. *arXiv preprint arXiv:2402.13753*.

Harry Dong, Xinyu Yang, Zhenyu Zhang, Zhangyang Wang, Yuejie Chi, and Beidi Chen. 2024. Get more with less: Synthesizing recurrence with kv cache compression for efficient llm inference. *arXiv preprint arXiv:2402.09398*.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The Pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.

David Haussler. 2018. Decision theoretic generalizations of the pac model for neural net and other learning applications. In *The Mathematics of Generalization*, pages 37–116. CRC Press.

Xijie Huang, Li Lyna Zhang, Kwang-Ting Cheng, and Mao Yang. 2023. Boosting llm reasoning: Push the limits of few-shot learning with reinforced in-context pruning. *arXiv preprint arXiv:2312.08901*.

Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023. Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt compression. *arXiv preprint arXiv:2310.06839*.

Kaiokendev. 2023. Things iḿ learning while training superhot. https://kaiokendev.github.io/til#extending-context-to-8k.

Lukasz Kaiser, Ofir Nachum, Aurko Roy, and Samy Bengio. 2016. Learning to remember rare events. In *International Conference on Learning Representations*.

Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva Reddy. 2023. The impact of positional encoding on length generalization in transformers. *arXiv preprint arXiv:2305.19466*.

Guolin Ke, Di He, and Tie-Yan Liu. 2020. Rethinking positional encoding in language pre-training. In *International Conference on Learning Representations*.

Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186.

Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2019. Generalization through memorization: Nearest neighbor language models. In *International Conference on Learning Representations*.

Shun Kiyono, Sosuke Kobayashi, Jun Suzuki, and Kentaro Inui. 2021. Shape: Shifted absolute position embedding for transformers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3309–3321.

Kuang-Huei Lee, Xinyun Chen, Hiroki Furuta, John Canny, and Ian Fischer. 2024. A human-inspired reading agent with gist memory of very long contexts. *arXiv preprint arXiv:2402.09727*.

Dacheng Li, Rulin Shao, Anze Xie, Ying Sheng, Lianmin Zheng, Joseph E. Gonzalez, Ion Stoica, Xuezhe Ma, and Hao Zhang. 2023. How long can open-source llms truly promise on context length?

Tatiana Likhomanenko, Qiantong Xu, Gabriel Synnaeve, Ronan Collobert, and Alex Rogozhnikov. 2021. Cape: Encoding relative positions with continuous augmented positional embeddings. *Advances in Neural Information Processing Systems*, 34:16079–16092.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Xiaoran Liu, Hang Yan, Chenxin An, Xipeng Qiu, and Dahua Lin. 2023. Scaling laws of rope-based extrapolation. In *The Twelfth International Conference on Learning Representations*.

Kun Luo, Zheng Liu, Shitao Xiao, and Kang Liu. 2024. Bge landmark embedding: A chunking-free embedding method for retrieval augmented long-context large language models. *arXiv preprint arXiv:2402.11573*.

Kai Lv, Xiaoran Liu, Qipeng Guo, Hang Yan, Conghui He, Xipeng Qiu, and Dahua Lin. 2024. Longwanjuan: Towards systematic measurement for long text quality. *arXiv preprint arXiv:2402.13583*.

Amirkeivan Mohtashami and Martin Jaggi. 2023. Landmark attention: Random-access infinite context length for transformers. *arXiv preprint arXiv:2305.16300*.

Matanel Oren, Michael Hassid, Yossi Adi, and Roy Schwartz. 2024. Transformers are multi-state rnns. *arXiv preprint arXiv:2401.06104*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023. Yarn: Efficient context window extension of large language models. In *The Twelfth International Conference on Learning Representations*.

Alexander Peysakhovich and Adam Lerer. 2023. Attention sorting combats recency bias in long context language models. *arXiv preprint arXiv:2310.01427*.

David Pollard. 1990. Empirical processes: Theory and applications. In *NSF-CBMS Regional Conference Series in Probability and Statistics*, pages i–86. JSTOR.

Ofir Press, Noah Smith, and Mike Lewis. 2021. Train short, test long: Attention with linear biases enables input length extrapolation. In *International Conference on Learning Representations*.

Zexuan Qiu, Jingjing Li, Shijue Huang, Wanjun Zhong, and Irwin King. 2024. Clongeval: A chinese benchmark for evaluating long-context large language models. *arXiv preprint arXiv:2403.03514*.

Zhijie Qu. 2023. Gpt rotational position embedding for length extrapolation. In *Proceedings of the 2023 6th International Conference on Machine Learning and Natural Language Processing*, pages 166–170.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.

Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506.

Nir Ratner, Yoav Levine, Yonatan Belinkov, Ori Ram, Inbal Magar, Omri Abend, Ehud Karpas, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. Parallel context windows for large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6383–6402.

Siyu Ren and Kenny Q Zhu. 2024. On the efficacy of eviction policy for key-value constrained generative language model inference. *arXiv preprint arXiv:2402.06262*.

Ninglu Shao, Shitao Xiao, Zheng Liu, and Peitian Zhang. 2024. Flexibly scaling large language models contexts through extensible tokenization. *arXiv preprint arXiv:2401.07793*.

Kaiqiang Song, Xiaoyang Wang, Sangwoo Cho, Xiaoman Pan, and Dong Yu. 2023. Zebra: Extending context window with layerwise grouped local-global attention. *arXiv preprint arXiv:2312.08618*.

Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. 2021. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*.

Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav Chaudhary, Xia Song, and Furu Wei. 2022. A length-extrapolatable transformer. *arXiv preprint arXiv:2212.10554*.

Mingxu Tao, Yansong Feng, and Dongyan Zhao. 2023. A frustratingly easy improvement for position embeddings via random padding. *arXiv preprint arXiv:2305.04859*.

MosaicML NLP Team. 2023. Introducing mpt-7b: A new standard for open-source, commercially usable llms. Accessed: 2023-05-05.

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

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Szymon Tworkowski, Konrad Staniszewski, Mikołaj Pacek, Yuhuai Wu, Henryk Michalewski, and Piotr Miłoś. 2023. Focused transformer: Contrastive training for context scaling. *arXiv preprint arXiv:2307.03170*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. https://github.com/kingoflolz/mesh-transformer-jax.

Cunxiang Wang, Ruoxi Ning, Boqi Pan, Tonghui Wu, Qipeng Guo, Cheng Deng, Guangsheng Bao, Qian Wang, and Yue Zhang. 2024. Novelqa: A benchmark for long-range novel question answering. *arXiv preprint arXiv:2403.12766*.

Yuhuai Wu, Markus Norman Rabe, DeLesley Hutchins, and Christian Szegedy. 2021. Memorizing transformers. In *International Conference on Learning Representations*.

Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. Efficient streaming language models with attention sinks. In *The Twelfth International Conference on Learning Representations*.

Kevin Yang, Dan Klein, Nanyun Peng, and Yuandong Tian. 2023. DOC: Improving long story coherence with detailed outline control. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3378–3465, Toronto, Canada. Association for Computational Linguistics.

Zi Yang and Nan Hua. 2024. Attendre: Wait to attend by retrieval with evicted queries in memory-based transformers for long context processing. *arXiv preprint arXiv:2401.04881*.

Dani Yogatama, Cyprien de Masson d'Autume, and Lingpeng Kong. 2021. Adaptive semiparametric language models. *Transactions of the Association for Computational Linguistics*, 9:362–373.

Tao Yuan, Xuefei Ning, Dong Zhou, Zhijie Yang, Shiyao Li, Minghui Zhuang, Zheyue Tan, Zhuyu Yao, Dahua Lin, Boxun Li, et al. 2024. Lv-eval: A balanced long-context benchmark with 5 length levels up to 256k. *arXiv preprint arXiv:2402.05136*.

Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. 2020. Big bird: Transformers for longer sequences. *Advances in neural information processing systems*, 33:17283–17297.

Peitian Zhang, Zheng Liu, Shitao Xiao, Ninglu Shao, Qiwei Ye, and Zhicheng Dou. 2024a. Soaring from 4k to 400k: Extending llm's context with activation beacon. *arXiv preprint arXiv:2401.03462*.

Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, et al. 2024b. ∞ bench: Extending long context evaluation beyond 100k tokens. *arXiv preprint arXiv:2402.13718*.

Zhenyu Zhang, Runjin Chen, Shiwei Liu, Zhewei Yao, Olatunji Ruwase, Beidi Chen, Xiaoxia Wu, and Zhangyang Wang. 2024c. Found in the middle: How language models use long contexts better via plug-and-play positional encoding. *arXiv preprint arXiv:2403.04797*.

Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark Barrett, et al. 2024d. H2o: Heavy-hitter oracle for efficient generative inference of large language models. *Advances in Neural Information Processing Systems*, 36.

Wangchunshu Zhou, Yuchen Eleanor Jiang, Peng Cui, Tiannan Wang, Zhenxin Xiao, Yifan Hou, Ryan Cotterell, and Mrinmaya Sachan. 2023. Recurrentgpt: Interactive generation of (arbitrarily) long text. *arXiv preprint arXiv:2305.13304*.

Dawei Zhu, Nan Yang, Liang Wang, Yifan Song, Wenhao Wu, Furu Wei, and Sujian Li. 2023. Pose: Efficient context window extension of llms via positional skip-wise training. In *The Twelfth International Conference on Learning Representations*.

## A Implementation Details

In this section, we introduce some implementation details of LM-Infinite as well as the hyperparameter selection.

### A.1 Computational resources

All experiments are run on single A100 GPUs with 80GB GPU memory each. The 200M length generalization runs for 20 hours. The downstream tasks take 3∼7 hours to evaluate each. Our work does not involve any training or fine-tuning. All models are loaded with Huggingface[2] code repository. Rouge and BLEU scores are loaded from evaluate[3] package. Datasets and models are used with permission from their licenses.

### A.2 The size of starting attention span

We vary the value of $n_{\text{starting}}$ and find LM-Infinite to be tolerant with it taking a wide range of values without sacrificing the NLL values. Specifically, we evaluate it on sequences of 16k length in the ArXiv validation set and calculate the average NLL.

| | | | $n_{\text{starting}}$ | | | |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 10 | 100 | 1000 | 2000 |
| 6.43 | 1.03 | 1.03 | 1.02 | 1.02 | 1.81 | 4.96 |

Table 4: Effect on LM-Infinite's NLL by varying $n_{\text{starting}}$.

### A.3 Reintroducing Top-$k$ Middle Tokens

This optional technique involves optionally attending to $k$ tokens in the middle with the largest attention logits. Here, the $k$ tokens are selected independently for each attention head and only apply to layers higher than $h$-th layer. These tokens have an attention distance of $d = \frac{1}{2}L_{\text{pre-train}}$. We select these hyper-parameters based on a held-out validation set of Passkey Retrieval. On Llama-2, we use $k = 5$ and $h = 5$. As an ablation study, we vary each hyper-parameter and observe its effects on Passkey Retrieval accuracy.

On the Qasper dataset, for both vanilla models and LM-Infinite, we use 6K sub-sequence of inputs as prompts and use a systematic prompt format described in Llama-2 paper (Touvron et al., 2023b).

---

[2]https://huggingface.cohttps://huggingface.co
[3]https://huggingface.co/docs/evaluate/index

| | | | $k$ | | | |
|---|---|---|---|---|---|---|
| 1 | 3 | 5 | 10 | 20 | 50 | 200 |
| 0.69 | 0.81 | 0.81 | 0.79 | 0.8 | 0.79 | 0.73 |

Table 5: Effect of varying $k$.

| **attention distance** | | | | |
|---|---|---|---|---|
| 512 | 1024 | 2048 | 3072 | 4096 |
| 0.78 | 0.79 | 0.81 | 0.68 | 0.63 |

Table 6: Effect of attention distance of the middle tokens.

## B Additional Related Work

After our preprint, there have been papers that cite our work and investigate zero-shot or few-shot length generalization of LLMs. As many absolute or relative position encoding methods are based on periodic functions, (Qu, 2023; Ding et al., 2024; Liu et al., 2023; Jiang et al., 2023) propose to apply LLMs (fine-tuned or not) on decreased period frequencies (which is equivalent to interpolating position indices) to adapt LLMs to longer sequences. Some other papers finetune LLMs with designed attention patterns (Oren et al., 2024; Zhang et al., 2024a) on long contexts, using neural tangent kernel (Peng et al., 2023), or with low-rank adaptation(LoRA) (Chen et al., 2023b). (Yang and Hua, 2024) instead proposes a wait-to-attend mechanism to extend length limits for memory-based Transformers. Other ways of key-value cache selection/eviction methods are investigated in (Ren and Zhu, 2024; Dong et al., 2024; Zhang et al., 2024d). Similarly, (Huang et al., 2023; Lee et al., 2024) tackles long context by learning to dynamically prune, select, or summarize contexts. Alternatively, context compression methods (Shao et al., 2024) propose to learn to compress long contexts into shorter embedding sequences. Some work proposes alternative position encodings (Song et al., 2023; Zhang et al., 2024c; Zhu et al., 2023) or landmark token embeddings (Luo et al., 2024) that enable extendable context limits. (Xiao et al., 2024) is a later concurrent work to ours with a similar approach to LLM length generalization. Unlike our work, they feed a sequence to an LLM token-by-token which limits their extreme length generalization (4M v.s. 200M of ours), and more importantly, they do not show improvements on downstream tasks without pre-training an LLM from scratch.

| | | | $h$ | | | |
|---|---|---|---|---|---|---|
| 0 | 4 | 5 | 6 | 8 | 16 | 24 |
| 0.81 | 0.94 | 0.94 | 0.91 | 0.46 | 0.45 | 0.46 |

Table 7: Effect of varying $h$.

Finally, there is a lot of new benchmarks (Qiu et al., 2024; Zhang et al., 2024b; Yuan et al., 2024; Wang et al., 2024; Lv et al., 2024; Bai et al., 2024) proposed to evaluate the long-context performance of LLMs.

## C  Formal Statement of Theorem 1

Let us denote the logit function with relative position encoding as $w(\mathbf{q}, \mathbf{k}, d) \in \mathbb{R}$. It maps the query $\mathbf{q}$, key $\mathbf{k}$, and their distance $d$, to an attention logit. The final attention weights are usually calculated by a softmax operation. For example, given $n$ tokens with indices $(1, \cdots, n)$, the attention by the last token on a preceding token at position $i$ is:

$$\text{Attn}(\text{token}_n, \text{token}_i) = \frac{e^{w(\mathbf{q}_n, \mathbf{k}_i, n-i)}}{\sum_{j=1}^n e^{w(\mathbf{q}_n, \mathbf{k}_j, n-j)}} \quad (1)$$

Then the formal theorem of Theorem 1 is as follows:

**Theorem 2.** *(Formal) Let $\mathbf{q}$ and $\mathbf{k}$ be random vectors sampled from training distributions $\mathcal{D}_\mathbf{q}$ and $\mathcal{D}_\mathbf{k}$, respectively, where $\mathcal{D}_\mathbf{q}$ and $\mathcal{D}_\mathbf{k}$ are the trained distributions for $\mathbf{q}$ and $\mathbf{k}$, respectively. We use the pseudo-dimension $\dim_P(\cdot)$ defined in (Pollard, 1990), which measures the representation capacity of a function family. Assume that the set of distance-based logit functions $\mathcal{H} = \{w(\cdot, \cdot, d) | d \in \mathbb{N}\}$ has bounded pseudo-dimension $\dim_P(\mathcal{H}) = r$[4]. Let us also define the distinguish-ability of two distances $d$ and $d'$ under $w$ as their expected squared difference: $\mu_w(d, d') = \mathbb{E}_{\mathbf{q} \sim \mathcal{D}_\mathbf{q}, \mathbf{k} \sim \mathcal{D}_\mathbf{k}} (w(\mathbf{q}, \mathbf{k}, d) - w(\mathbf{q}, \mathbf{k}, d'))^2$. We assume that $w$ is not limited to recognizing only a finite group of distances, otherwise, all distances longer than a threshold will become almost the same as shorter distances. Formally, for any $n$, there is a partition of $[0..n]$ into $\alpha(n)$ groups so that, $\mu_w(d, d') \le \epsilon$ for any $d, d'$ from the same group. $\alpha(n) \in \mathbb{N}$ is non-decreasing and unbounded function. Then we have:*

$$\sup_{\mathbf{q}, \mathbf{k}, d \le n} |w(\mathbf{q}, \mathbf{k}, d)| \ge \left( \frac{\alpha(n)}{2} \right)^{\frac{1}{2r}} \frac{\epsilon}{4e}.$$

---

[4]This is true for most current techniques. See discussions in Appendix F

We first borrow a lemma from (Haussler, 2018), which we paste below. Note that a cover size $\mathcal{N}(\epsilon, \mathcal{H}, \mu)$ is defined as the minimum cardinal of a cover-set $\mathcal{H}'$ so that any element of $h \in \mathcal{H}$ will be within $\epsilon$ distance to at least one element $h' \in \mathcal{H}'$.

**Lemma 3.** *Let $\mathcal{H}$ be a function family mapping from space $X$ to range $[0, B]$, and its pseudo-dimension $\dim_P(\mathcal{H}) = r$. Then for any probabilistic measure $P$ on $X$, and $\epsilon \in [0, B]$, we have that the $\epsilon$ cover of $\mathcal{H}$ under metric $\mu(h_1, h_2) = \mathbb{E}_{x \sim P}(h_1(x) - h_2(x))^2$ is bounded by:*

$$\mathcal{N}_P(\epsilon, \mathcal{H}, \mu) \le 2 \left( \frac{2eB}{\epsilon} \ln \frac{2eB}{\epsilon} \right)^r$$

With this lemma, we can go on to prove Theorem 2 as follows.

*Proof.* We prove by contradiction. Assume that $\sup_{\mathbf{q}, \mathbf{k}, d \le n} |w(\mathbf{q}, \mathbf{k}, d)| < a = \left( \frac{\alpha(n)}{2} \right)^{\frac{1}{2r}} \frac{\epsilon}{4e}$. Without loss of generality, we can shift all the values to range $[0, 2a]$. Then the function family $\mathcal{H} = \{w(\cdot, \cdot, d) | d \in \mathbb{N}\}$ will have cover size $\mathcal{N}_P(\epsilon, \mathcal{H}, \mu) \le 2 \left( \frac{4ea}{\epsilon} \ln \frac{4ea}{\epsilon} \right)^r < \alpha(n)$.

However, this is smaller than the number of different distances and relative weight attentions $\mathcal{H}$, which means that at least 2 functions will be close to each other $(w(\cdot, \cdot, d), w(\cdot, \cdot, d'))^2 < \epsilon$. This constitutes a contradiction with the distinguishability condition. $\qquad \square$

## D  Formal Statement and Proof of Proposition 1

The formal statement of Proposition 1 is the following:

**Proposition 2.** *(Attention Entropy Explosion) Let $w_1, w_2, \cdots, w_n \in [-B, B]$ be a sequence of attention logits. Then the entropy of the attention distribution they span is asymptotically lower bounded by $\ln n$:*

$$H \left( \left( \frac{e^{w_i}}{\sum_{j=1}^n e^{w_j}} \Big| 1 \le i \le n \right) \right) = \Omega(\ln n)$$

The entropy approaches $+\infty$ as $n$ grows larger.

*Proof.* Note that entropy on a discrete distribution is defined as $\text{Entropy}(P) = -\sum_i p_i \ln p_i$. Then the attention entropy determined by attention logits

$\{w_i | 1 \le i \le n\}$ is

$$
\begin{aligned}
&\text{Entropy(Attention)} \\
&= -\sum_i \frac{e^{w_i}}{\sum_j e^{w_j}} \ln \frac{e^{w_i}}{\sum_j e^{w_j}} \\
&= -\sum_i \frac{e^{w_i}}{\sum_j e^{w_j}} \left( w_i - \ln \sum_j e^{w_j} \right) \\
&= -\sum_i \frac{e^{w_i}}{\sum_j e^{w_j}} w_i + \ln \sum_j e^{w_j} \\
&\ge -\max_i w_i + \ln(n e^{-B}) \\
&\ge \ln n - 2B \\
&= \Omega(\ln n)
\end{aligned}
$$

$\square$

## E   More on Implicitly Encoded Positions

We also plot the token features of more layers with PCA projection to the 2D plane in Figure 7. We see that from layer 2 to higher layers, the initial few tokens occupy a distinct region with later tokens. Therefore, if these tokens are discarded by window attention during attention, the attention output, which is a weighted sum of $v_i$ vectors, will reside in a different region. This can explain why windowed attention does not work and why the first few tokens need to be kept by our $\Lambda$-shaped attention.

## F   Pseudo-Dimension Assumption on Attention Logit Functions

In Theorem 2, we assumed that the family of distance-based logit functions $\mathcal{H} = \{w(\cdot, \cdot, d) | d \in \mathbb{N}\}$ has a finite pseud-dimension. In this section, we demonstrate that most current implementations of relative positional encodings do have a finite pseudo-dimension. Let us discuss a few examples in the following:

**T5-Bias and Alibi**   It is easy to see that, the difference between any two logit functions is uniform: $w(\cdot, \cdot, d_1) - w(\cdot, \cdot, d_2) = \text{bias}(d_1) - \text{bias}(d_2)$ regardless of the input. Therefore this family cannot shatter any set larger than 2, so the pseudo-dimension is 1.

**Windowed Attention**   This operation is equivalent to limiting the family to a finite size $|\mathcal{H}| = d_{\max} + 1$, where $d_{\max}$ is the size of the window. Therefore $\dim_P(\mathcal{H}) \le d_{\max} + 1$.

**NoPE**   As there is no explicit positional encoding implemented, all distances are equivalent. The pseudo-dimension is 1.

**RoPE, CAPE, and XPos**   For RoPE, the logit function $w(\mathbf{q}, \mathbf{k}, d)$ is the weighted sum of finite fixed sinusoidal functions $\{\sin(\omega_i d), \cos(\omega_i d)\}$. The size of this set is equivalent to the feature dimension number $k$. We know that $\dim_P(\mathcal{H}_1 + \mathcal{H}_1) \le \dim_P(\mathcal{H}_1) + \dim_P(\mathcal{H}_2)$. Also, the scaling of a single function can only have a pseudo-dimension of 2. Therefore, the whole family has a bounded pseudo-dimension $\dim_P(\mathcal{H}) \le 2k$. The analysis on CAPE and XPos is similar.

## G   Computational Efficiency Evaluation

To evaluate the computational efficiency of LM-Infinite, we run the Llama-2-7B model on 100 sequences of 32k length in the ArXiv dataset. The hardware is a single A100 GPU with 80GB GPU memory. As the memory is not enough to host the whole computation graph during decoding, we use DeepSpeed (Rasley et al., 2020) with Zero3 optimization. We also have to modify the computation code to further reduce GPU memory usage and prevent out-of-memory errors. With that in mind, the vanilla Llama-2-7B model encodes with an average speed of 48.19s per sequence, while LM-Infinite encodes with an average of 15.26s per sequence, a 3.16x speedup. The vanilla Llama-2-7B model decodes with 7.34s per token, while LM-Infinite decodes with 2.70s per token, a 2.72x speedup. We also evaluate the GPU memory usage on 32k-length inputs, the statistics of which are profiled with PyTorch Profiler. The vanilla model uses 33.2Gb memory per sequence, while LM-Infinite uses 4.41Gb per sequence, a 7.53× memory saving.
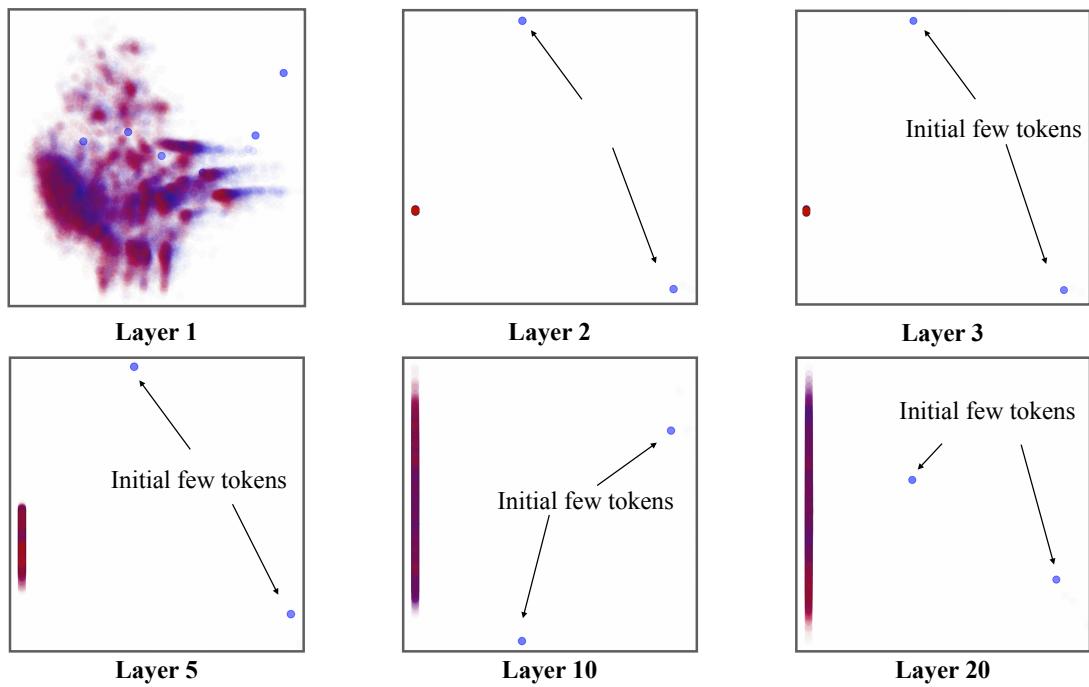
## H   Example Generation

Figure 7: In Llama, at second or higher layers, the initial few tokens encode a strong position signal and occupy a distinct feature region. Abandoning them might move the attention output vector out of the pre-training distribution.

Table 8: Example text generations on ArXiv and OpenWebText2 corpora after 8k context lengths.

| | **ArXiv Dataset** |
|---|---|
| **Context** | *Our proof of the Kottwitz formula imitates these arguments of Langlands and Rapoport of deriving the Kottwitz formula from their conjecture. As such, their theory of Galois gerbs is a major ingredient in this work. Another essential ingredient is a recent work of Kisin [@Kisin17] on the aforementioned conjecture of Langlands and Rapoport. Previously, Kottwitz [@Kottwitz92] proved the formula in PEL-type cases (of simple Lie type A or C) by a method which is based on the Honda-Tate theory. This method however cannot be applied in general Hodge-type situations, and indeed* |
| **Vanilla Llama-2** | *of this (of over-equary-her, and [( and, in the...cister '– and an of the model to by. by ..., this, by the. It, and it, 7. –(of an equist (of the. to the [[[WNE (to. and for the (((de in the (for the andistile–c.-[de (for in an inc ort, ort (betness in >with (with, based (and (> such ((c of a or for the abstract as. of *. (up. W in the, of the or an. and of a.0 in (and of (over in such [near (—of (([her (pro – [p for [with the of a: their in the in of '[alism in the, [(du of that of the ((both [equ of the [for ([of as of ... and of the of their of ((—for of the of the (t of the of the of the of his of the. of the. of the of it (D with the d. and of of the of the of the of [of of a of a, [[for for a per of a of its of a (for of the and of the* |
| **LM-Infinite** | *Kottwitz [@Kottwitz84a] had formulated the conjecture, but at the time no idea had yet been established on how to prove it. This conjecture involves some algebraic statements of automorphic forms which have an algebraic nature, but are only valid for automorphic forms (which are not holomorphic automorphic forms [@Kaplansky1966]), whose (abstract) objects are Galois extensions of $F$ and $\bar{F}$ (which in general are not necessarily finite extensions), in other words a $p$-valued Hodge-type structure (in $n$-dimension). The Langlands-Rapoport conjecture is now established [@ScholzeShin13] whose formulation is in terms of (a more explicit class of) Hodge-type structures rather than these abstract Galois objects, and the automorphic side for $p$-adic L-functions are due to Kisin [@Kisin15].* |

| | **OpenWebText2** |
|---|---|
| **Context** | *This auto-structuration of a complete whole is defined as "auto-regulation". In the individual, the latter is established by biological rhythms, biological & mental regulations and mental operations. These can be theoretically formalized. Piaget refuses to accept that "real" dialectical tensions between physical objects are the true foundations of thought and cognition (its possibility, genesis & progressive development). Piaget never fills in what reality is. He maintains no ontological view on reality-as-such, considered to be the borderline of both the developing subject and its objective world, stage after stage. The cognitive is approached as a process, for rationality grows in developmental stages, each calling for a particular cognitive structure on the side of* |
| **Vanilla Llama-2** | *the in that, of it and and ¿ " of the. of " of " of the pr of of of the of the of the [. of " – s (1 ( of s of the " of quends of it &nd of trail of the, the &, 1. – a and at-and and and the " for on the is, core by " ( as (:: a. -. ( on its to " is a, ( on, th is of all but". ".. in the, in the the the and in the for the. , and of óf the or of the :! of " A. of the of the, of the of the of [of the of the: of a th, and s of the of the the as of the, B of h. the —. the : a, and, in a: & (se and [ ( (t, d of the and....:,. of the (,. ( A. A. Pro: T. I " (b in the óf the, of of: of " of, of the of the [ of, (: - " of the , of s of all ma of the [ and of the of the of [ of the of the of the ( the & ( w [ (:, and ad and and [ " of a in. for the, of ín (* |
| **LM-Infinite** | *cognition. It develops gradually. For Piaget, the cognitive originates in experience as it is subjected to structure : i.e. as it is structured, ordered and transformed. In the first stage, the "sensorial stage", structures are formed and integrated based on perception. These structures are built by organic, biological and psychic transformations, which lead to experience. The "concrete operational stage" of intellectual development refers to a stage when the mind of the child starts to use the logic of the general rule and the rule of the general case. There is now a logical, conceptual and operational distinction of concepts. Reasoning is made explicit by applying logical operations, such as subtraction, addition and multiplication to the "mental" object as well as to the "perceived" world of reality. The childś logic can now make use of the logical operations, though for him (her) it is a conceptual understanding of abstract operations. This is Piagetś concept of "genetic development". In the "formal operational stage", logical operations are combined using logical or conceptual structures.* |