

# EFFICIENT LONG-RANGE LANGUAGE MODELING WITH SELF-SUPERVISED CAUSAL RETRIEVAL

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recently, retrieval-based language models (RLMs) have received much attention. However, most of them leverage a pre-trained retriever with fixed parameters, which may not adapt well to causal language models. In this work, we propose Grouped Cross-Attention, a novel module enabling joint pre-training of the retriever and causal LM, and apply it to long-context modeling. For a given input sequence, we split it into chunks and use the current chunk to retrieve past chunks for subsequent text generation. Our innovation allows the retriever to learn how to retrieve past chunks that better minimize the auto-regressive loss of subsequent tokens in an end-to-end manner. By integrating top- $k$  retrieval, our model can be pre-trained efficiently from scratch with context lengths up to 64K tokens. [Our experiments demonstrate that our model achieves superior performance in various tasks against strong baselines, and 100% accuracy in the needle-in-a-haystack \(NIAH\) test with a 16M context length.](#)

## 1 INTRODUCTION

Transformers (Vaswani et al., 2017), serving as the backbone of large language models (LLM), have revolutionized language modeling and demonstrated exceptional performance across a wide range of natural language processing tasks (Brown et al., 2020; Achiam et al., 2023; Touvron et al., 2023; Dubey et al., 2024). While Transformers excel in representational power, their quadratic computational complexity and increasing memory demands as input length grows pose formidable challenges for modeling long contexts. Various approaches, such as recurrent memory (Dai et al., 2019), and linear attention (Katharopoulos et al., 2020) techniques, are proposed to improve the efficiency and effectiveness of Transformers in handling extended inputs. Nevertheless, these approaches often sacrifice the random-access flexibility of attention (Mohtashami & Jaggi, 2023) during inference.

In this work, we explore long-range language modeling in Transformers from the perspective of retrieval-based language models (RLMs) (Asai et al., 2023). Typically, RLMs (Rubin & Berant, 2024; Yen et al., 2024) divide an input sequence into chunks, retrieve relevant ones from the history for the current input, and then integrate the retrieved chunks into the decoder to predict subsequent tokens. By choosing top- $k$  chunks as a “dynamic context”, RLMs overcome the efficiency challenges in long-context modeling while maintaining random-access flexibility. However, most RLMs (Lewis et al., 2020; Borgeaud et al., 2022) rely on separately pre-trained retrievers with fixed parameters, which hinders their ability to adapt to the causal LMs. Although a straightforward approach is training the retriever end-to-end to select chunks that minimize auto-regressive loss of subsequent tokens, it is rarely explored. The main challenges are twofold: firstly, while relevance scores guide chunk selection, [these scores do not participate in the next token prediction, thus unable to receive gradient backpropagation from the auto-regressive loss](#). Secondly, the large search space brought by long contexts often results in efficiency and flexibility issues for pre-training.

To tackle these challenges, we propose Grouped Cross-Attention (GCA), a novel module enabling efficient end-to-end joint optimization of the retriever and causal LM, thus the retriever can *learn to retrieve* past chunks that most effectively reduce the auto-regressive loss of subsequent tokens, which we refer to as *causal retrieval*. GCA enables the relevance scores to participate in the next token prediction in a differentiable way. Specifically, GCA can be understood as a chunk-wise analogy of token-wise self-attention. In self-attention, considering the next token prediction in causal Transformers, self-attention scores could be viewed as the relevance scores of the current token to

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

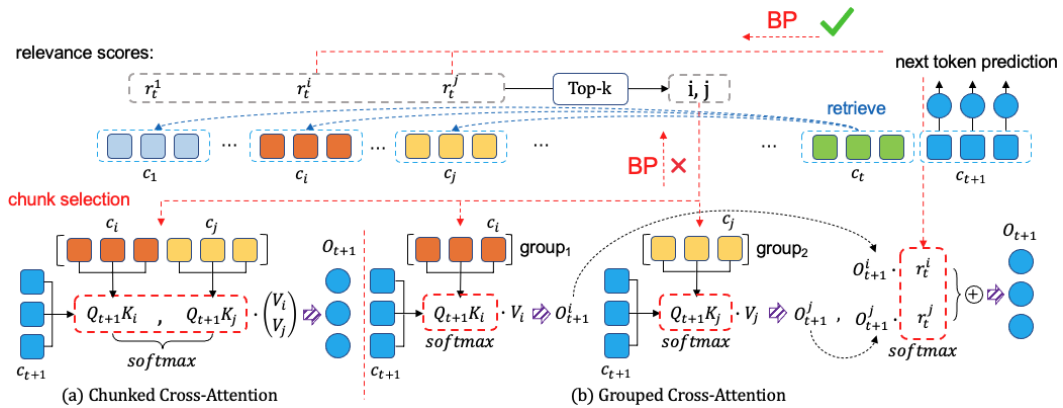


Figure 1: Comparing previous works with GCA. Consider current chunk  $c_t$  with its past chunk relevance scores  $r_t^k$ , where  $k \in \{1, \dots, t-1\}$ ,  $r_t^i$  and  $r_t^j$  are the top two. In this example, each chunk contains 3 tokens, whose query, key, and value vectors are denoted as  $Q, K, V$ . (a) In previous work, information from retrieved chunks is fused into LM decoders via Chunked Cross-Attention, in which relevance scores are merely used for chunk selection. Thus the loss can not back-propagate to the scores. (b) GCA separately applies Cross-Attention with the two chunks, yielding intermediate outputs  $O_{t+1}^i$  and  $O_{t+1}^j$ . The softmaxed relevance scores serve as weights to fuse these intermediates into LM decoders and thus can receive back-propagation from the loss.

past tokens. These scores serve as weights to fuse information gathered from past tokens to predict the next token. Analogously, we divide the input sequence into chunks and use the relevance scores between the *current* and past chunks as weights to fuse information for the *next* chunk prediction. A detailed comparison between previous works and GCA is depicted in Figure 1. By appending GCA after self-attention in Transformer layers, we introduce **Differentiable Retrieval-based Transformers (DRT)**, enabling pre-training from scratch with context lengths up to 64K. To make pre-training efficient, we sample top- $k$  past chunks according to the relevance scores for each chunk to perform GCA, along with fixed-size sliding window self-attention (Child et al., 2019), achieving linear complexity for the entire input sequence’s attention operations. During inference, we offload hidden states of past chunks to CPU memory and reload them when retrieved. It introduces additional memory-swap but largely reduces memory footprint.

In our experiments, we evaluate our model on tasks such as long-range language modeling, summarization, and the needle-in-a-haystack (NIAH) tests. The results demonstrate that DRT significantly outperforms all baselines with comparable pre-training costs and much lower inference costs. Notably, in the NIAH test, DRT trained with a 16K context length maintains nearly 100% accuracy on inputs up to 16M tokens. More interestingly, case studies on the arXiv-math dataset suggest that long-range reasoning ability emerges in DRT, which retrieves lemmas, variants, or functions defined distantly but used in the next chunk. These findings suggest that GCA has the potential to be a fundamental component in retrieval-based LMs. Overall, our main contributions are:

1. We propose a novel module called **Grouped Cross-Attention (GCA)**, which allows dense retrievers learn to retrieve guided by auto-regressive loss in an end-to-end manner efficiently.
2. Building upon GCA, we introduce **Differentiable Retrieval-based Transformers (DRT)**, which is fast and memory-efficient in both pre-training and inference on long texts, but still maintains the random-access flexibility and excellent extrapolation capability.
3. We implement a hardware-aware GCA based on FlashAttention-2 (Dao, 2024), significantly reducing the training and inference time. The code will be made publicly available.

## 2 RELATED WORKS

**Relation to RPT & Landmark Attention.** There are two long-range LMs closely related to ours. One of them is **Retrieval-Pretrained Transformer (RPT)** (Rubin & Berant, 2024). The key difference between DRT and RPT is the training approach of the retriever. During data-preparation, for each chunk, RPT picks relevant past chunks by using BM25 (Robertson & Zaragoza, 2009), concatenates them with the current chunk, and evaluates them by a *reference LM* like Pythia 1.4B (Biderman

et al., 2023). The past chunks that increase the probability of the next chunk are identified as ‘gold chunks’ to train RPT’s retriever. However, such a complex data preparation process limits scalability and flexibility in pre-training and post-training (Lee, 2024). In contrast, DRT is pre-trained end-to-end. By employing a sliding window size larger than the chunk size, it effectively uses feedback from subsequent several chunks to train the retriever. Its flexibility also allows for adaptive multi-hop retrieval. **Landmark Attention (LA)** (Mohtashami & Jaggi, 2023) is another close work. LA is pre-trained with short contexts but capable of handling long contexts during inference. It addresses long-range language modeling by modifying self-attention KV Cache. During inference, each token, at each layer, selects top- $k$  chunks based on token-to-chunk attention scores and appends their key and value vectors to the current KV cache of self-attention. The token-to-chunk attention scores are trained in an end-to-end manner with a grouped softmax technique. However, it has to perform top- $k$  chunk selection per token, per layer, which incurs significant extra costs during inference. Moreover, it fails to extrapolate on longer context length. Our method combines the chunk-retrieval and grouped softmax ideas, resolving the aforementioned issues while balancing training efficiency and inference performance.

**Long-Range Language Modeling.** Various methods have been proposed to improve long-range language modeling. One line of research is introducing memorization to Transformers via recurrence. Many works (Dai et al., 2019; Burtsev & Sapunov, 2020; Martins et al., 2022; Hutchins et al., 2022) compress past information into fixed-sized vectors. However, these methods often sacrifice the flexibility to attend to arbitrary past tokens. Meanwhile, other works focus on maintaining random-access flexibility of attention. Memorizing Transformers (Wu et al., 2022) appends retrieved past keys and values to the current attention segment via  $k$ -NN search, but they do not back-propagate gradients to them. CEPE (Yen et al., 2024) retrieves previous chunks using an independently trained dense retriever and fuses them into the decoder. During the training process, the decoder parameters are fixed, and only the encoder is adjusted. A notable distinction in our work is the end-to-end optimization of all parameters, particularly the retriever.

**Efficient Language Modeling.** Many works have been done to reduce the training and inference cost of LLM. One direction is sparse attention, which includes limiting the attention window to a small range around each token (Child et al., 2019; Zaheer et al., 2020; Beltagy et al., 2020), approximating attention matrix (Wang et al., 2020), leveraging locality-sensitive hashing (LSH) for key vectors retrieval (Kitaev et al., 2020), and hierarchical self-attention (Ren et al., 2021). However, empirically most efficient Transformers sacrifice performance for efficiency. Recently, state-space models (Gu & Dao, 2023; Dao & Gu, 2024) and RNN models (Beck et al., 2024) provide new architecture alternatives, with comparable performance to Transformers but much lower cost for inference. We argue that our core innovation GCA is flexible enough to be incorporated into these models as an additional module to obtain random-access flexibility.

**Retrieval-Augmented Language Models.** Retrieval-augmented LMs leverage a retriever to access relevant external knowledge, enhancing their generation capabilities. In some works, the retriever can be jointly trained with the LM such as REALM (Gua et al., 2020). However, its computational complexity limits its extension to causal LMs. On the other hand, in most other works (Lewis et al., 2020; Izacard & Grave, 2021; Borgeaud et al., 2022; Ivgi et al., 2023), retriever parameters are fixed, preventing optimization for retrieving information that best predicts subsequent tokens.

**Unsupervised Dense Retrieval.** In the line of research on unsupervised dense retrieval, early works leverage Inverse Cloze Task (Lee et al., 2019) to train retriever unsupervisedly, where a sentence is randomly sampled from a document and the task is to predict its surrounding context. However, this approach still lags behind BM25 on long-tail entities (Sciavolino et al., 2021). Recently, contrastive learning methods (Izacard et al., 2022; Gao & Callan, 2022) have shown improved results by creating positive and negative pairs from sentences within the same or different documents, respectively. However, these unsupervisedly trained retrievers are typically not optimized for causal LMs, so they may not guarantee to retrieve the most pertinent information.

### 3 DIFFERENTIABLE RETRIEVAL-BASED TRANSFORMER

A typical architecture of RLMs (Borgeaud et al., 2022; Yen et al., 2024; Rubin & Berant, 2024) appends Chunked Cross-Attention (CCA) after self-attention to fuse information from retrieved

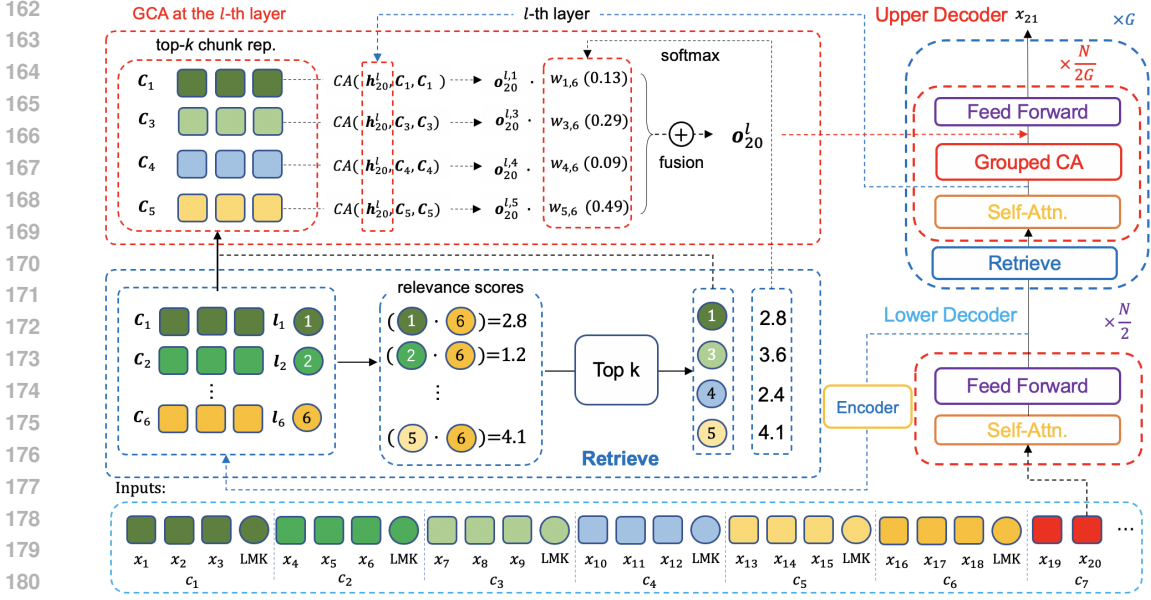


Figure 2: The illustration shows how the current chunk  $c_6$  retrieves past chunks for better next token prediction of  $x_{20}$  in the next chunk  $c_7$ . The landmark representation  $l_6$  is used to compute relevance scores with past chunks, selecting the top four. The hidden states of  $x_{20}$  at the  $l$ -th layer,  $h_{20}^l$ , perform Cross-Attention (CA) with all tokens in each separate chunk. The chunk-wise CA outputs,  $o_{20}^l$ , are fused via weighted sum, whose weights are softmaxed relevance scores.

chunks, in which a retriever is merely used to pick relevant chunks. Our approach makes retriever learnable by replacing CCA with GCA, which represents the key innovation of this work. The novelty of GCA lies in using relevance scores to fuse information from retrieved chunks for LM decoders, enabling the retriever to adaptively learn to select the best past chunks for predicting subsequent tokens, guided by the auto-regressive loss. This section details the model architecture, training, and inference.

### 3.1 MODEL ARCHITECTURE

DRT is composed of  $N$  Transformer-like layers. Similar to RETRO (Borgeaud et al., 2022), the input sequence of DRT is equally divided into chunks. Formally, given a sequence  $\mathbf{x} = [x_1, x_2, \dots, x_L]$  with  $L$  tokens, we divide the sequence into  $\frac{L}{S}$  chunks, where  $S$  is the chunk size, denoted as  $\{c_1, c_2, \dots, c_{L/S}\}$ , where  $x_i \in c_{\lceil i/S \rceil}$ . Similar to Landmark Attention, we insert a special token LMK at the end of each chunk, which summarizes the preceding content via self-attention.

**Forward Pass.** Figure 2 illustrates the forward pass of a token in DRT. DRT layers are bifurcated into upper and lower sections like in RPT (Rubin & Berant, 2024). The key differences are the introduction of GCA and the further division of the upper layers into  $G$  groups, enabling learning to retrieve on the fly and adaptive multi-hop retrieval. The lower layers comprise standard Transformer layers while each upper layer has an additional GCA module after self-attention. In the forward pass, the chunk hidden states output by the lower layers, besides being fed to the upper layers, are also fed into a bi-directional Transformer encoder, which further contextualizes the representations with inner-chunk positional encoding, yielding  $C_k \in \mathbb{R}^{S \times d}$  and  $l_k \in \mathbb{R}^d$  shared across all upper layers, where  $d$  is the hidden state dimension. At the  $g$ -th upper decoder group, chunk  $c_t$  retrieves the top- $k$  relevant past chunks for its next chunk:

$$r_t^{g,k} = \frac{\mathbf{h}_t^{g \top} \mathbf{l}_k}{\sqrt{d}}, \quad C_t^g = \text{Top-k}([r_t^{g,1}, \dots, r_t^{g,t-1}]). \quad (1)$$

Here,  $\mathbf{h}_t^g \in \mathbb{R}^d$  represents the landmark representation output by the decoder layer just before the  $g$ -th group, which accumulates all information from groups 1 to  $g-1$ . This enables  $\mathbf{h}_t^g$  to retrieve relevant information based on previously retrieved chunks thus achieving multi-hop retrieval.  $r_t^{g,k}$  represents the causal relevance score of  $c_k$  to  $c_t$ .  $C_t^g$  contains the indices of past chunks with top- $k$

relevance scores. The retrieved chunks are shared among the subsequent layers within the same group. The upper layers apply GCA to fuse retrieved information into the decoder.

**Grouped Cross-Attention.** For the  $l$ -th layer, let  $\mathbf{H}_{t+1}^l$  and  $\hat{\mathbf{H}}_{t+1}^l \in \mathbb{R}^{(S+1) \times d}$  denote the batched token representations including the landmark token in the next chunk before and after GCA. In GCA, we perform Cross-Attention (CA) separately and fuse results via relevance scores:

$$g(l) = \lceil (l - \frac{N}{2}) / \frac{N}{2G} \rceil, \quad \mathbf{O}_{t+1,k}^l = \mathbf{CA}(\mathbf{H}_{t+1}^l, \mathbf{C}_k, \mathbf{C}_k), \quad k \in \mathcal{C}_t^{g(l)},$$

$$w_t^{g(l),k} = \frac{\exp(r_t^{g(l),k})}{\sum_{k' \in \mathcal{C}_t^{g(l)}} \exp(r_t^{g(l),k'})}, \quad \mathbf{O}_{t+1}^l = \sum_k w_t^{g(l),k} \mathbf{O}_{t+1,k}^l, \quad \hat{\mathbf{H}}_{t+1}^l = \text{Norm}(\mathbf{H}_{t+1}^l + \mathbf{O}_{t+1}^l). \quad (2)$$

Here  $g(l)$  converts the layer index to the group index and  $\mathbf{C}_k \in \mathbb{R}^{S \times d}$  represents token representations of the  $k$ -th retrieved chunk.  $\mathbf{O}_{t+1,k}^l \in \mathbb{R}^{(S+1) \times d}$  represents the information that  $S+1$  tokens in chunk  $c_{t+1}$  gather from past chunk  $c_k$ .  $w_t^{g(l),k}$  is the normalized relevance score after softmax, serving as the weight of  $\mathbf{O}_{t+1,k}^l$  for information fusion. The final fused results of GCA is  $\mathbf{O}_{t+1}^l$ .

Since  $\mathbf{C}_k$  is shared across layers, we use the same  $\mathbf{K}, \mathbf{V}$  linear transformations across layers to compact model parameters and reduce memory footprint. For each head  $h$ , we have CA defined as:

$$\mathbf{CA}(\mathbf{H}_{t+1}^l, \mathbf{C}_k, \mathbf{C}_k)_h \triangleq \text{Softmax}\left(\frac{Q_h^l(\mathbf{H}_{t+1}^l)K_h(\mathbf{C}_k)^T}{\sqrt{d_h}}\right)V_h(\mathbf{C}_k) \quad (3)$$

Here,  $d_h$  is per head dimension, and  $Q_h^l, K_h, V_h$  are linear transformations per head, where  $Q_h^l$  varies across layers and  $K_h, V_h$  are layer-shared. The final CA outputs are concatenated vectors from all heads. It is worth noting that GCA is easy to integrate with FlashAttention-2, as detailed in the pseudo-code in Appendix A.2.

**GCA vs CCA.** A key distinction between GCA and CCA lies in how softmax is applied to cross-attention matrices as shown in Figure 1(a)(b). In CCA, all retrieved chunks are concatenated and softmax is directly applied to the whole attention matrix to fuse token-level information. Notably, relevance scores are entirely excluded from the process. In contrast, GCA applies softmax to each chunk’s attention matrix separately. This modification allows information to be gathered from each chunk separately. The softmaxed relevance scores then serve as soft choices (Hu et al., 2021; 2022), participating in the next token prediction thus receiving back-propagated gradients.

**Encoder-Decoder Variant.** Our model architecture could be easily modified to an encoder-decoder-based variant like RETRO (Borgeaud et al., 2022) by directly applying the encoder to chunk token embeddings instead of the lower layers’ outputs. This variant allows for retrieval from trillions of tokens with acceptable storage costs. We will elaborate on this in § 3.3.

### 3.2 TRAINING

The pre-training objective of DRT is the next token prediction, but there are certain details as discussed below. We also give a detailed time complexity analysis of training.

**Gumbel Top-k sampling.** The core idea of self-supervised retrieval is making candidate chunks compete with each other by softmaxing relevance scores as weights. The weights of the chunks contributing most to the next chunk prediction are enhanced while the weights of the rest are suppressed. To balance exploration and exploitation, we sample chunks based on relevance scores instead of always picking the top-k, enabling highly relevant chunks to be more likely chosen while still considering lower-scoring ones. A simple trick is to add Gumbel noise (Gumbel, 1954) to the raw scores before the top-k operation. Importantly, this noise doesn’t affect subsequent operations.

**Encoder-Decoder Pre-training.** To enhance the encoder’s representational ability, we train it with MLM in addition to the auto-regressive loss during the first half of pre-training. Specifically, while the decoder sees all tokens, the encoder’s tokens are partially masked, whose outputs are used in both the decoder’s GCA and computing the MLM objective. Masking causes inaccurate encoding that may disturb subsequent computation and training of the decoder, so we eventually remove the MLM training in the second half of the pretraining.

**Time Complexity.** Our approach reduces training complexity by compressing quadratic operations. Vanilla Transformers have a complexity of  $\mathcal{O}(NL^2)$  for full self-attention. In DRT, we encode chunks with  $S$  tokens into landmark representations, performing chunk-wise full attention to compute relevance scores within  $\mathcal{O}(G\frac{L^2}{S^2})$  for  $G$  groups of upper layers. By employing sliding-window attention and top- $k$  retrieval, we scale down self-attention and GCA costs to  $\mathcal{O}(G\frac{L^2}{S^2} + \frac{N}{2}LKS + NLW)$ , where  $K$  is the retrieved chunk number and  $W$  is the window size,  $K, W \ll L$ . This largely reduces the complexity but maintains the random-access flexibility.

### 3.3 INFERENCE

**Memory Offloading.** During the inference stage of Transformers, the default memory cost for KV Cache is  $\mathcal{O}(NLd)$ . To reduce GPU memory usage, we can offload past chunk representations to CPU memory. This results in a spatial complexity of the GPU memory usage  $\mathcal{O}(\frac{Ld}{S} + \frac{N}{2}Ksd + NWd)$  during inference. Here,  $\frac{Ld}{S}$  is the memory footprint of landmark representations, while the remaining terms account for the GCA and sliding-window KV cache. Although each retrieval involves gathering representations from CPU memory and transferring them to the GPU, this operation occurs  $G$  times every  $S$  tokens. Therefore, the cost of memory exchange from chunk retrievals is minimal.

**Infinitely Long Context Retrieval.** To mimic humans’ capability for random access to memories, we extend the retrievable context to all pre-trained tokens, where contextualized token representations in chunks could be regarded as memory. A straightforward approach, similar to RPT, is to store chunk-level representations and token-level hidden states as key-value pairs in a Faiss (Douze et al., 2024). However, this approach requires significant disk space. For instance, a 100 billion token corpus with 1,024-dimensional hidden states, even with Product Quantization, demands approximately 25TB. In contrast, using an encoder-decoder variant, we only store a chunk’s landmark representation and positions in the corpus without saving hidden states. When retrieving, we access the corresponding tokens by document position and re-encode them to obtain their hidden states, achieving retrieval from billions of tokens with disk cost reducing by  $S$  fold.

## 4 EXPERIMENTS

We compare DRT with prior works in long-range language modeling, wall-clock training time, inference cost, extrapolation capability, and infinite-length context retrieval. Notably, we include the closely related works RPT and Landmark Attention.

### 4.1 EXPERIMENTAL SETUP

#### 4.1.1 DATASETS

**PG19.** PG19 (Rae et al., 2020) is a language modeling benchmark widely used to evaluate long-range text understanding capabilities of models. It includes a collection of full-length books from Project Gutenberg across a wide range of genres, styles, and topics. The dataset is particularly useful for evaluating models’ abilities to capture dependencies and context over long sequences of text.

**ArXiv-math.** ArXiv-math is a corpus consisting of mathematical papers from arXiv. Characterized by a sustained coherence over long distance, the corpus requires models’ capability of correctly referring to long-range history information and using long-range context effectively for predictions. We use the preprocessed corpus and data splits from Azerbayev et al. (2023).

**MiniPile.** MiniPile (Kaddour, 2023) is a 6GB subset of the deduplicated 825GB *The Pile* (Gao et al., 2021) corpus, which covers sub-datasets from webpages, dialogue, books, science and code.

#### 4.1.2 MODELS

**DRT<sub>retrieval</sub> $\times G$ .** A DRT consists of 12 Transformer decoder layers, divided into 6 lower and 6 upper layers, with upper layers further split into  $G$  groups. The sliding window size is set to  $W=512$ , the chunk size is set to  $S=64$ , and 8 chunks are retrieved for GCA, resulting in an attention field of

512 ( $8 \times 64$ ). We implement hardware-aware GCA based on Triton (Tillet et al., 2019)<sup>1</sup>. As we employ various parameter settings across different experiments, further details are provided in Appendix A.1. **DRT<sub>enc-dec</sub>**: the variant of DRT introduced in § 3.1.

**Base LM.** Our base LM is based on the implementation of TinyLlama (Zhang et al., 2024) combined with Flash Attention2 (Dao, 2024) enabling ALiBi (Press et al., 2022) and sliding window attention (Child et al., 2019). We compare models against the baseline across various configurations. One configuration involves 12 layers with a sliding window of 512 tokens, aligning with the DRT sliding window size. Another configuration of 12 layers with a 768-token sliding window ensures the same attention field coverage, as  $12 \times 768 = 12 \times 512 + 6 \times 512$ (GCA). The strongest baseline, with 14 layers and a 658-token sliding window, has a parameter count comparable to our DRT while maintaining a similar total attention field across all 12 layers, calculated as  $658 \times 14 \approx 12 \times 512 + 6 \times 512$ .

**Retrieval-Pretrained Transformer (RPT).** Since the official implementation is in JAX and the code for distilling the retriever is not released, we reimplement RPT in PyTorch and replace the retriever with Contriever (Izacard et al., 2022). [Elaborations can be found in Appendix A.6.](#)

**Landmark Attn.** We use the official Llama-like implementation<sup>2</sup> of Landmark Attention. Similar to Base LM, we extend the length of the self-attention range from 512 to 768 to ensure it shares the same attention field as DRT.

**Block-Recurrent TFM.** Since the official implementation of Block-Recurrent Transformer is also based on JAX, we utilized a PyTorch implementation<sup>3</sup> to ensure all baselines are running with the same framework.

**Ablations.** *w/o Triton*: A naively implemented version of GCA without Triton. *w/o enc*: The decoder-only version of DRT, which uses the hidden states of the middle layer of the decoder stack to retrieve history chunks. Differently put, this model can be attained by simply removing the encoder from DRT<sub>enc-dec</sub>. *w/o mlm*: The architecture is exactly the same as DRT<sub>enc-dec</sub>, while the only difference lies in the training process by eliminating the masked language modeling part. *w/o gumbel top-k*: The architecture is exactly the same as DRT, while the only difference lies in the training process by eliminating the gumbel noise when selecting the top-k chunks.

## 4.2 LONG-RANGE LANGUAGE MODELING

In this section, we evaluate DRT against baselines in long-range language modeling on PG19 and arXiv-math, and report their respective perplexities. All models are pre-trained with the same attention field and a 16K context by default, except for baselines that cannot efficiently pre-train on long contexts. To ensure fairness, we adjusted these baselines. Detailed hyper-parameters are provided in Appendix A.1.

**Results.** From Table 1, we have several observations. **Firstly**, DRT outperforms all baselines where the evaluation length exceeds 16K. While DRT performs retrieval for  $G$  times every 64 tokens, LA performs retrieval at every token and in every layer, potentially offering better random access flexibility. However, DRT still surpasses LA on longer inputs. This is likely because LA follows a “train short, test long” approach due to the need for full attention during pre-training, whereas DRT can be directly pre-trained on long input sequences. Thanks to GCA, our key innovation, DRT can randomly access distant contexts during pre-training with the same attention field as in the baselines, allowing it to better utilize long-range information during pre-training with negligible extra training costs. **Secondly**, in terms of parameter efficiency, we compared against Base LM with two additional layers in the decoder stack. It can be observed that on shorter evaluation lengths, the baseline has a slight advantage. However, with a longer context, our model consistently leads. This experiment suggests precisely retrieving semantic knowledge from a long context may be more beneficial for improving the language model than increasing model parameters. **Thirdly**, multi-hop retrieval yields positive gains with a low marginal cost. It allows upper layers to access more diverse

<sup>1</sup><https://github.com/triton-lang/triton/blob/main/python/tutorials/06-fused-attention.py>

<sup>2</sup><https://github.com/epfml/landmark-attention>

<sup>3</sup><https://github.com/lucidrains/block-recurrent-transformer-pytorch>

Model	Time	#Param. dec/enc	k	attn. win.	PG19↓		ArXiv-math↓	
					valid	test	valid	test
Train length=16K, eval length = 2K								
BaseLM	1×	124M	-	512	15.00	14.10	3.31	3.31
(Sliding window+Alibi)	1.03×	124M	-	768	14.86	13.96	3.24	3.24
+2 layers	1.15×	138M	-	658	14.71	13.83	<b>3.06</b>	<b>3.06</b>
RPT <sub>contriever</sub> (our impl.)	2.5×	133M/14M	8	512	<u>14.81</u>	<u>13.92</u>	<u>3.24</u>	<u>3.24</u>
Landmark Attn.	1.5×	124M	4	768	<b>14.41</b>	<b>13.40</b>	3.17	3.16
Block Recurrent TFM	2×	155M	-	768	15.99	15.00	3.33	3.32
DRT <sub>enc-dec</sub>	1.22×	133M/14M	4	512	<u>14.90</u>	<u>14.02</u>	<u>3.24</u>	<u>3.24</u>
DRT <sub>retrieval × 1</sub>	1.22×	133M/14M	8	512	<u>14.65</u>	<u>13.78</u>	<u>3.24</u>	<u>3.24</u>
DRT <sub>retrieval × 2</sub>	1.24×	133M/14M	8	512	<u>14.56</u>	<u>13.69</u>	<u>3.22</u>	<u>3.22</u>
Train length=16K, eval length = 16k								
BaseLM	1×	124M	-	512	14.55	13.68	3.06	3.06
(Sliding window+Alibi)	1.03×	124M	-	768	14.36	13.49	2.95	2.95
+2 layers	1.15×	138M	-	658	14.23	13.37	2.95	2.94
RPT <sub>contriever</sub> (our impl.)	2.5×	133M/14M	8	512	<u>14.39</u>	<u>13.52</u>	<u>2.93</u>	<u>2.92</u>
Landmark Attn.	1.5×	124M	4	768	14.10	<u>13.21</u>	3.02	3.02
Block Recurrent TFM	2×	155M	-	768	15.59	14.60	3.14	3.14
DRT <sub>enc-dec</sub>	1.22×	133M/14M	4	512	<u>14.42</u>	<u>13.53</u>	<u>2.93</u>	<u>2.93</u>
DRT <sub>retrieval × 1</sub>	1.22×	133M/14M	8	512	<u>14.05</u>	<u>13.21</u>	<u>2.89</u>	<u>2.89</u>
DRT <sub>retrieval × 2</sub>	1.24×	133M/14M	8	512	<b>14.02</b>	<b>13.18</b>	<b>2.85</b>	<b>2.85</b>
Train length=16K, eval length = 32k								
BaseLM	1×	124M	-	512	14.50	13.64	3.05	3.04
(Sliding window+Alibi)	1.03×	124M	-	768	14.30	13.46	2.93	2.92
+2 layers	1.15×	138M	-	658	14.18	13.34	2.93	2.92
RPT <sub>contriever</sub> (our impl.)	2.5×	133M/14M	8	512	<u>14.35</u>	<u>13.49</u>	<u>2.91</u>	<u>2.91</u>
Landmark Attn.	1.5×	124M	4	768	14.19	13.33	3.07	3.07
Block Recurrent TFM	2×	155M	-	768	15.61	14.56	3.13	3.12
DRT <sub>enc-dec</sub>	1.22×	133M/14M	4	512	<u>14.38</u>	<u>13.52</u>	<u>2.91</u>	<u>2.91</u>
DRT <sub>retrieval × 1</sub>	1.22×	133M/14M	8	512	<u>14.01</u>	<u>13.19</u>	<u>2.85</u>	<u>2.85</u>
DRT <sub>retrieval × 2</sub>	1.24×	133M/14M	8	512	<b>13.98</b>	<b>13.16</b>	<b>2.81</b>	<b>2.81</b>
Ablation studies								
-w/o Triton	1.45×	133M/14M	8	512	—	—	—	—
-w/o enc.	-eval len=16k		8	512	<u>14.31</u>	<u>13.44</u>	<u>2.97</u>	<u>2.97</u>
-w/o MLM	-eval len=16k		4	512	<u>14.43</u>	<u>13.57</u>	<u>3.07</u>	<u>3.06</u>
-w/o gumbel top-k	-eval len=16k		8	512	<u>14.36</u>	<u>13.46</u>	<u>2.90</u>	<u>2.90</u>

Table 1: Perplexity for all datasets. We highlight the best results in **bold** and underline the second best.

chunks and enables further retrieval based on previous retrieval results. **Finally**, ablation studies show that all the training techniques we add bring positive improvements. In conclusion, the above results fully demonstrate that the GCA module can indeed bring effective gains in modeling long texts, and it is more advantageous compared to other baselines.

### 4.3 DOWNSTREAM TASK EVALUATION AND EFFICIENCY ANALYSIS

In this section, we fine-tune all baselines and evaluate them against downstream tasks including summarization (Nallapati et al., 2016; Narayan et al., 2018), single NIAH test, and multi-hop NIAH proposed by Hsieh et al. (2024). The details for the downstream tasks are described in Appendix A.4. Then we analyze the inference cost, the relationship between training time and context length, and the extrapolation capability of DRT. In the inference cost analysis, We skip RPT because it has a similar cost to DRT. In the extrapolation experiments, we utilize the pre-trained models described in § 4.2 to assess their performance with extended context lengths.

**Results.** From Table 2, we observe that DRT significantly outperforms all baselines in the summarization tasks, validating its capability to effectively utilize long contexts. Notably, in the single NIAH test, DRT maintains 100% performance even with a context length of up to 16 million tokens, demonstrating its strong length generalization ability in long-context scenarios. Furthermore, in the 2-hop NIAH test, DRT<sub>retrieval × 2</sub> performs comparably to Landmark Attention at context lengths below 16K tokens while successfully extrapolating to longer context lengths beyond 64K tokens. Additionally, DRT<sub>retrieval × 2</sub> significantly outperforms DRT<sub>retrieval × 1</sub>, confirming our hypothesis that conducting causal retrieval every  $G$ -layer contributes to scenarios requiring multi-hop retrievals.



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444

Models	Retrieval	Single NIAH $\uparrow$							
		1K	2K	4K	8K	16K	32K	64K	16M
Base LM (+2 layers)	-	60.43	29.54	15.37	8.30	3.89	2.13	0.0	-
Block Recurrent TFM	-	66.80	30.61	13.96	7.60	6.01	2.13	4.29	-
RPT <sub>Contriever</sub>	fixed	42.13	18.45	11.66	6.71	4.24	1.42	2.86	-
Landmark Attn.	adaptive	<b>99.98</b>	99.08	<b>99.82</b>	97.74	97.88	96.45	0.00	-
DRT <sub>retrieval<math>\times</math>1</sub>	adaptive	98.23	99.12	98.50	98.76	98.59	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>
DRT <sub>retrieval<math>\times</math>2</sub>	adaptive	99.69	<b>99.56</b>	99.65	<b>99.47</b>	<b>99.65</b>	99.99	<b>100.00</b>	<b>100.00</b>

Models	XSum $\uparrow$			CNN/DailyMail $\uparrow$			2-hop NIAH with noises $\uparrow$					
	R-1	R-2	R-L	R-1	R-2	R-L	1K	4K	16K	64K	256K	4M
BaseLM	29.43	8.04	23.26	32.38	12.57	22.61	3.60	1.15	0.71	0.0	0.0	-
+2 layers	29.74	8.26	23.48	34.14	14.09	23.60	16.60	5.83	1.06	0.0	0.0	-
Landmark Attn.	27.98	6.96	21.99	34.06	13.80	23.77	<b>90.82</b>	<b>88.35</b>	<b>86.41</b>	0.0	0.0	-
DRT <sub>retrieval<math>\times</math>1</sub>	30.30	8.59	23.92	36.27	15.88	25.08	41.07	33.39	39.93	38.57	35.29	34.29
DRT <sub>retrieval<math>\times</math>2</sub>	<b>30.39</b>	<b>8.64</b>	<b>23.98</b>	<b>36.39</b>	<b>15.96</b>	<b>25.15</b>	88.52	84.45	86.21	<b>81.43</b>	<b>94.11</b>	<b>79.41</b>

Table 2: The performances of various models on summarization tasks and NIAH tests.

445  
446  
447  
448  
449  
450  
451  
452

	Prompt #tokens	Generated #tokens	w/ cpu offload		w/o cpu offload	
			time/token $\downarrow$	mem. cost $\downarrow$	time/token $\downarrow$	mem. cost $\downarrow$
Landmark Attn. / Base LM	16K	128	160.9 $\times$	1.98 $\times$	4.16 $\times$	32 $\times$
	48K	128	163.1 $\times$	2.98 $\times$	4.25 $\times$	96 $\times$
Block Recurrent TFM / Base LM	16K	128	-	-	2.85 $\times$	2 $\times$
	48K	128	-	-	2.85 $\times$	2 $\times$
DRT <sub>retrieval<math>\times</math>1</sub> / Base LM	16K	128	1.25 $\times$	1.54 $\times$	1.06 $\times$	4.08 $\times$
	48K	128	1.27 $\times$	1.62 $\times$	1.08 $\times$	9.41 $\times$

Table 3: The inference time per token and memory footprint ratio compared to the Base LM (12 layers with a 512 sliding window), with lower values indicating better performance.

453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467

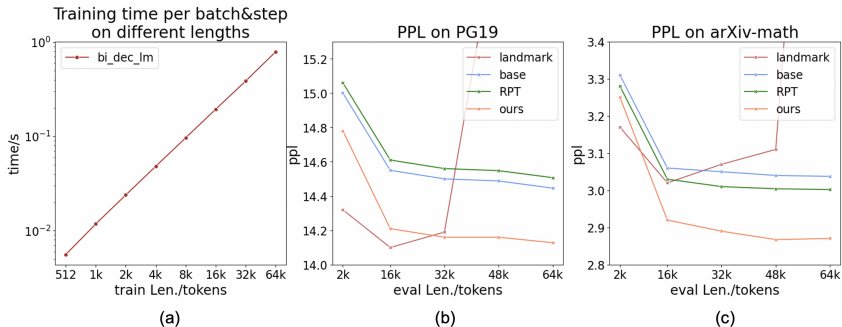


Figure 3: Training speed, extrapolation ability

468  
469  
470  
471  
472

Table 3 shows DRT significantly outperforms Landmark Attn in terms of memory footprint and inference speed. The main overhead for Landmark Attn arises from modifications to the self-attention KV cache and gathering tensors from memory offloaded to the CPU. In DRT, thanks to the chunk-wise retrieval, we perform retrieval only once every 64 tokens, which is  $1/(12 \times 64)$  of the corresponding operation in Landmark Attn.

473  
474  
475  
476  
477  
478  
479  
480

From Figure 3(a), it can be observed that the total training time increases linearly with the increment of the sequence length. In the extrapolation experiments, both BaseLM and DRT perform well as shown in Figure 3(b)(c). However, Landmark Attn fails to extrapolate at longer eval lengths. We believe a possible explanation is that a longer context increases the probability of retrieving irrelevant distant chunks, which stems from its limitations in pre-training directly on long contexts. An overly short pre-training window prevents access to longer-range information, making the model unable to learn to reduce the attention scores of long-range noise chunks. DRT benefits from being able to pre-train directly on long contexts, alleviating this issue to a certain extent.

481  
482  
483  
484  
485

#### 4.4 RETRIEVAL FROM SIMULATED INFINITELY LONG CONTEXT

We wonder whether a densely pre-trained retriever can generalize from a limited to an unlimited context. To verify this, we simulate an infinitely long context by retrieving from all pre-trained tokens. The details could be found in Appendix A.5.

Model	Corpus Retrieval	Self Retrieval	MiniPile↓	Settings:
Base LM (attn. win. 512)	No	No	12.68	<i>Self-Retrieval:</i>
DRT w/ random retrieval	Yes	No	12.68	Indicates if retrieval of past 48K tokens is enabled.
DRT w/ Contriever	Yes	No	12.65	
DRT w/ Contriever	No	Yes	12.25	<i>Corpus-Retrieval:</i>
DRT	Yes	No	12.67	Indicates if retrieval from the pre-training corpus is enabled.
DRT	Yes	Yes	12.30	
DRT	No	Yes	12.18	

Table 4: Valid set perplexity for MiniPile under different settings.

**Results.** From Table 4, we observe that the perplexity of DRT slightly rises instead of declines when we extend the context to infinity. Retrieving information from a limited-length context yielded the best results among all methods. A possible explanation is that retrieval from the vast amount of chunks from the corpus may yield results with similar representations but semantically irrelevant content. Notably, only Contriever benefits from corpus-retrieval when self-retrieval is disabled. The key distinction between Contriever and DRT’s inherent retriever is that Contriever incorporates random negative sampling during training. However, unlike Contriever, our negative samples are drawn from a fixed-size context, meaning the negative sample candidates are fixed. Contriever, on the other hand, performs random negative sampling via in-batch sampling, allowing for a more extensive negative sample space. This insight could potentially contribute to retrieval from trillion tokens in future works.

#### 4.5 CASE STUDIES

... We start first with the following lemma which is useful to establish the Quillen equivalence.  $\begin{matrix} \text{\ref{lem-reflect-equiv}} \end{matrix}$  Let  $\mathbf{M}$  be a symmetric monoidal model category that is combinatorial and left proper. Assume that the transferred (natural) model structure on  $\mathit{com}$  exists. Let  $\sigma : \mathbf{C} \rightarrow \mathbf{D}$  be a morphism between usual commutative...  
 ... weak equivalence between co-Segal categories is just a level-wise weak equivalence.  $\begin{matrix} \text{\ref{prop-eta-kx-loc-equiv}} \end{matrix}$  For any  $\mathbf{F} \in \mathit{coms}$ , the canonical map  $\mathbf{F} \rightarrow |\mathbf{F}|$  is an equivalence in  $\mathit{comsepc}$  i.e, it’s a  $kb(I)$ -local equivalence in  $\mathit{comsep}$  (whence in  $\mathit{comse}$ ).  $\end{matrix}$ ...  
 ... Thanks to  $\text{\ref{prop-eta-kx-loc-equiv}}$ , we know that  $\eta : \mathbf{F} \rightarrow |\mathbf{F}|$  is always a  $kb(I)$ -local equivalence. Then by 3-for-2 we see that  $\sigma$  is a  $kb(I)$ -local equivalence if and only if  $ol(\sigma)$  is. Now thanks to  $\text{\ref{lem-reflect-equiv}}$  we know that ...

Figure 4: In the case above, Retrieved 2nd best chunk describes a proposition which appears in the current chunk. retrieved best chunk and retrieved 3rd best chunk are adjacent chunks which introduce the lemma used in the next chunk.

By analyzing DRT’s retrieval results on the arXiv-math dataset, we find some intriguing cases. A case is given in Figure 4. When retrieving past chunks, the results not only include the definition of prepositions referenced in the current chunk but also lemmas to be used in the next chunk. This validates the idea of causal retrieval, allowing us to not only retrieve semantically similar content but also information that better predicts the next chunk. More case studies can be found in Appendix A.3.

## 5 CONCLUSION & FUTURE WORKS

In this study, we successfully optimize the retriever module with the auto-regressive LM objective in an end-to-end manner. The core innovation lies in the Grouped Cross-Attention (GCA), which makes relevance scores learnable by using them to fuse information retrieved by the current chunk for next chunk prediction. Combined with Gumbel top-k sampling, this approach enables the pre-training of LMs on context lengths extending up to 64K tokens.

In future work, we will explore self-supervised causal retrieval from vast amounts of tokens outside the context. Meanwhile, we will combine structured representations (Hu et al., 2024b;a) to achieve multi-granular retrieval.

## REFERENCES

- 540  
541  
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Akari Asai, Sewon Min, Zexuan Zhong, and Danqi Chen. Retrieval-based language models and  
546 applications. In Yun-Nung (Vivian) Chen, Margot Margot, and Siva Reddy (eds.), *Proceed-*  
547 *ings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 6:*  
548 *Tutorial Abstracts)*, pp. 41–46, Toronto, Canada, July 2023. Association for Computational Lin-  
549 guistics. doi: 10.18653/v1/2023.acl-tutorials.6. URL <https://aclanthology.org/2023.acl-tutorials.6>.
- 551 Zhangir Azerbayev, Edward Ayers, and Bartosz Piotrowski. Proof-pile: A pre-training dataset of  
552 mathematical text, 2023. URL [https://huggingface.co/datasets/hoskinson-center/](https://huggingface.co/datasets/hoskinson-center/proof-pile)  
553 [proof-pile](https://huggingface.co/datasets/hoskinson-center/proof-pile).
- 555 Maximilian Beck, Korbinian Pöppel, Markus Spanring, Andreas Auer, Oleksandra Prudnikova,  
556 Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. xlstm: Extended  
557 long short-term memory. *CoRR*, abs/2405.04517, 2024. doi: 10.48550/ARXIV.2405.04517. URL  
558 <https://doi.org/10.48550/arXiv.2405.04517>.
- 559 Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer.  
560 *CoRR*, abs/2004.05150, 2020. URL <https://arxiv.org/abs/2004.05150>.
- 561 Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien,  
562 Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff,  
563 Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large  
564 language models across training and scaling. In Andreas Krause, Emma Brunskill, Kyunghyun  
565 Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference*  
566 *on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of  
567 *Proceedings of Machine Learning Research*, pp. 2397–2430. PMLR, 2023. URL [https://](https://proceedings.mlr.press/v202/biderman23a.html)  
568 [proceedings.mlr.press/v202/biderman23a.html](https://proceedings.mlr.press/v202/biderman23a.html).
- 569 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Mil-  
570 lican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego  
571 de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren  
572 Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol  
573 Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. Im-  
574 proving language models by retrieving from trillions of tokens. In Kamalika Chaudhuri, Ste-  
575 fanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International*  
576 *Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*,  
577 volume 162 of *Proceedings of Machine Learning Research*, pp. 2206–2240. PMLR, 2022. URL  
578 <https://proceedings.mlr.press/v162/borgeaud22a.html>.
- 579 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-  
580 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-  
581 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,  
582 Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler,  
583 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCand-  
584 lish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot  
585 learners. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan,  
586 and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual*  
587 *Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-*  
588 *12, 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html)  
589 [1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html).
- 590 Mikhail S. Burtsev and Grigory V. Sapunov. Memory transformer. *CoRR*, abs/2006.11527, 2020.  
591 URL <https://arxiv.org/abs/2006.11527>.
- 592  
593 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse  
transformers. *CoRR*, abs/1904.10509, 2019. URL <http://arxiv.org/abs/1904.10509>.

- 594 Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc Viet Le, and Ruslan Salakhut-  
595 dinov. Transformer-xl: Attentive language models beyond a fixed-length context. In Anna Ko-  
596 rihonen, David R. Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Conference of the*  
597 *Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019,*  
598 *Volume 1: Long Papers*, pp. 2978–2988. Association for Computational Linguistics, 2019. doi:  
599 10.18653/V1/P19-1285. URL <https://doi.org/10.18653/v1/p19-1285>.
- 600  
601 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In  
602 *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Aus-*  
603 *tria, May 7-11, 2024*. OpenReview.net, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=mZn2Xyh9Ec)  
604 [mZn2Xyh9Ec](https://openreview.net/forum?id=mZn2Xyh9Ec).
- 605 Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through  
606 structured state space duality. In *Forty-first International Conference on Machine Learning, ICML*  
607 *2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL [https://openreview.](https://openreview.net/forum?id=ztn8FCR1td)  
608 [net/forum?id=ztn8FCR1td](https://openreview.net/forum?id=ztn8FCR1td).
- 609  
610 Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-  
611 Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. *CoRR*,  
612 abs/2401.08281, 2024. doi: 10.48550/ARXIV.2401.08281. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2401.08281)  
613 [48550/arXiv.2401.08281](https://doi.org/10.48550/arXiv.2401.08281).
- 614 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
615 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
616 *arXiv preprint arXiv:2407.21783*, 2024.
- 617  
618 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason  
619 Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile:  
620 An 800gb dataset of diverse text for language modeling. *CoRR*, abs/2101.00027, 2021. URL  
621 <https://arxiv.org/abs/2101.00027>.
- 622 Luyu Gao and Jamie Callan. Unsupervised corpus aware language model pre-training for dense  
623 passage retrieval. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Pro-*  
624 *ceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Vol-*  
625 *ume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pp. 2843–2853. Asso-  
626 ciation for Computational Linguistics, 2022. doi: 10.18653/V1/2022.ACL-LONG.203. URL  
627 <https://doi.org/10.18653/v1/2022.acl-long.203>.
- 628  
629 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces.  
630 *CoRR*, abs/2312.00752, 2023. doi: 10.48550/ARXIV.2312.00752. URL [https://doi.org/](https://doi.org/10.48550/arXiv.2312.00752)  
631 [10.48550/arXiv.2312.00752](https://doi.org/10.48550/arXiv.2312.00752).
- 632 Emil Julius Gumbel. *Statistical theory of extreme values and some practical applications: a series*  
633 *of lectures*, volume 33. US Government Printing Office, 1954.
- 634  
635 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented  
636 language model pre-training. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th*  
637 *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*  
638 *Research*, pp. 3929–3938. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v119/guu20a.html)  
639 [v119/guu20a.html](https://proceedings.mlr.press/v119/guu20a.html).
- 640  
641 Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, and  
642 Boris Ginsburg. RULER: What’s the real context size of your long-context language models? In  
643 *First Conference on Language Modeling*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=kIoBbc76Sy)  
644 [kIoBbc76Sy](https://openreview.net/forum?id=kIoBbc76Sy).
- 645 Xiang Hu, Haitao Mi, Zujie Wen, Yafang Wang, Yi Su, Jing Zheng, and Gerard de Melo. R2D2:  
646 recursive transformer based on differentiable tree for interpretable hierarchical language model-  
647 ing. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the*  
*59th Annual Meeting of the Association for Computational Linguistics and the 11th International*

- 648 *Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Pa-*  
649 *pers), Virtual Event, August 1-6, 2021*, pp. 4897–4908. Association for Computational Linguistics, 2021. doi: 10.18653/V1/2021.ACL-LONG.379. URL [https://doi.org/10.18653/v1/](https://doi.org/10.18653/v1/2021.acl-long.379)  
650 [2021.acl-long.379](https://doi.org/10.18653/v1/2021.acl-long.379).  
651
- 652 Xiang Hu, Haitao Mi, Liang Li, and Gerard de Melo. Fast-r2d2: A pretrained recursive neural  
653 network based on pruned CKY for grammar induction and text representation. In Yoav Gold-  
654 berg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Em-*  
655 *pirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emi-*  
656 *rates, December 7-11, 2022*, pp. 2809–2821. Association for Computational Linguistics, 2022.  
657 doi: 10.18653/V1/2022.EMNLP-MAIN.181. URL [https://doi.org/10.18653/v1/2022.](https://doi.org/10.18653/v1/2022.emnlp-main.181)  
658 [emnlp-main.181](https://doi.org/10.18653/v1/2022.emnlp-main.181).  
659
- 660 Xiang Hu, Pengyu Ji, Qingyang Zhu, Wei Wu, and Kewei Tu. Generative pretrained structured  
661 transformers: Unsupervised syntactic language models at scale. In Lun-Wei Ku, Andre Mar-  
662 tins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association*  
663 *for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, Au-*  
664 *gust 11-16, 2024*, pp. 2640–2657. Association for Computational Linguistics, 2024a. URL  
665 <https://aclanthology.org/2024.acl-long.145>.
- 666 Xiang Hu, Qingyang Zhu, Kewei Tu, and Wei Wu. Augmenting transformers with recursively  
667 composed multi-grained representations. In *The Twelfth International Conference on Learning*  
668 *Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024b. URL  
669 <https://openreview.net/forum?id=u859gX7ADC>.
- 670 DeLesley Hutchins, Imanol Schlag, Yuhuai Wu, Ethan Dyer, and Behnam Neyshabur. Block-  
671 recurrent transformers. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho,  
672 and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on*  
673 *Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November*  
674 *28 - December 9, 2022*. URL [http://papers.nips.cc/paper\\_files/paper/2022/](http://papers.nips.cc/paper_files/paper/2022/hash/d6e0bbb9fc3f4c10950052ec2359355c-Abstract-Conference.html)  
675 [hash/d6e0bbb9fc3f4c10950052ec2359355c-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/d6e0bbb9fc3f4c10950052ec2359355c-Abstract-Conference.html).
- 676 Maor Ivgi, Uri Shaham, and Jonathan Berant. Efficient long-text understanding with short-text  
677 models. *Trans. Assoc. Comput. Linguistics*, 11:284–299, 2023. doi: 10.1162/TACL\_A\_00547.  
678 URL [https://doi.org/10.1162/tacl\\_a\\_00547](https://doi.org/10.1162/tacl_a_00547).  
679
- 680 Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for  
681 open domain question answering. In Paola Merlo, Jörg Tiedemann, and Reut Tsarfaty (eds.),  
682 *Proceedings of the 16th Conference of the European Chapter of the Association for Computa-*  
683 *tional Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021*, pp. 874–880. As-  
684 sociation for Computational Linguistics, 2021. doi: 10.18653/V1/2021.EACL-MAIN.74. URL  
685 <https://doi.org/10.18653/v1/2021.eacl-main.74>.
- 686 Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand  
687 Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learn-  
688 ing. *Trans. Mach. Learn. Res.*, 2022, 2022. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=jKN1pXi7b0)  
689 [jKN1pXi7b0](https://openreview.net/forum?id=jKN1pXi7b0).
- 690 Jean Kaddour. The minipile challenge for data-efficient language models. *CoRR*, abs/2304.08442,  
691 2023. doi: 10.48550/ARXIV.2304.08442. URL [https://doi.org/10.48550/arXiv.2304.](https://doi.org/10.48550/arXiv.2304.08442)  
692 [08442](https://doi.org/10.48550/arXiv.2304.08442).  
693
- 694 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are  
695 rnns: Fast autoregressive transformers with linear attention. In *Proceedings of the 37th Inter-*  
696 *national Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, vol-  
697 *ume 119 of Proceedings of Machine Learning Research*, pp. 5156–5165. PMLR, 2020. URL  
698 <http://proceedings.mlr.press/v119/katharopoulos20a.html>.
- 699 Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In  
700 *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia,*  
701 *April 26-30, 2020*. OpenReview.net, 2020. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=rkgNKkHtvB)  
[rkgNKkHtvB](https://openreview.net/forum?id=rkgNKkHtvB).

- 702 Jieh-Sheng Lee. Instructpatentgpt: Training patent language models to follow instructions with  
703 human feedback. *CoRR*, abs/2406.16897, 2024. doi: 10.48550/ARXIV.2406.16897. URL  
704 <https://doi.org/10.48550/arXiv.2406.16897>.  
705
- 706 Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised open  
707 domain question answering. In Anna Korhonen, David R. Traum, and Lluís Màrquez (eds.),  
708 *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019,*  
709 *Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pp. 6086–6096. Association for  
710 Computational Linguistics, 2019. doi: 10.18653/V1/P19-1612. URL [https://doi.org/10.](https://doi.org/10.18653/v1/p19-1612)  
711 [18653/v1/p19-1612](https://doi.org/10.18653/v1/p19-1612).
- 712 Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin,  
713 Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian  
714 Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP  
715 tasks. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan,  
716 and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual*  
717 *Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-*  
718 *12, 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html)  
719 [6b493230205f780e1bc26945df7481e5-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html).
- 720 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International*  
721 *Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019.*  
722 OpenReview.net, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- 723 Pedro Henrique Martins, Zita Marinho, and André F. T. Martins.  $\infty$ -former: Infinite mem-  
724 ory transformer. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Pro-*  
725 *ceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Vol-*  
726 *ume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pp. 5468–5485. Asso-  
727 ciation for Computational Linguistics, 2022. doi: 10.18653/V1/2022.ACL-LONG.375. URL  
728 <https://doi.org/10.18653/v1/2022.acl-long.375>.
- 729 Amirkeivan Mohtashami and Martin Jaggi. Random-access infinite context length for transform-  
730 ers. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey  
731 Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on*  
732 *Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December*  
733 *10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/](http://papers.nips.cc/paper_files/paper/2023/hash/ab05dc8bf36a9f66edbff6992ec86f56-Abstract-Conference.html)  
734 [ab05dc8bf36a9f66edbff6992ec86f56-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/ab05dc8bf36a9f66edbff6992ec86f56-Abstract-Conference.html).
- 735 Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Gulçehre, and Bing Xiang. Abstrac-  
736 tive text summarization using sequence-to-sequence RNNs and beyond. In Stefan Riezler and  
737 Yoav Goldberg (eds.), *Proceedings of the 20th SIGNLL Conference on Computational Natural*  
738 *Language Learning*, pp. 280–290, Berlin, Germany, August 2016. Association for Computational  
739 Linguistics. doi: 10.18653/v1/K16-1028. URL <https://aclanthology.org/K16-1028>.
- 740 Shashi Narayan, Shay B. Cohen, and Mirella Lapata. Don’t give me the details, just the sum-  
741 mary! topic-aware convolutional neural networks for extreme summarization. In Ellen Riloff,  
742 David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii (eds.), *Proceedings of the 2018 Confer-*  
743 *ence on Empirical Methods in Natural Language Processing*, pp. 1797–1807, Brussels, Bel-  
744 gium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/  
745 D18-1206. URL <https://aclanthology.org/D18-1206>.
- 746 Ofir Press, Noah A. Smith, and Mike Lewis. Train short, test long: Attention with linear bi-  
747 ases enables input length extrapolation. In *The Tenth International Conference on Learning*  
748 *Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022. URL  
749 <https://openreview.net/forum?id=R8sQPpGCv0>.
- 750 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language  
751 models are unsupervised multitask learners. 2019. URL [https://api.semanticscholar.](https://api.semanticscholar.org/CorpusID:160025533)  
752 [org/CorpusID:160025533](https://api.semanticscholar.org/CorpusID:160025533).  
753
- 754 Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap.  
755 Compressive transformers for long-range sequence modelling. In *International Conference on*  
*Learning Representations*, 2020. URL <https://openreview.net/forum?id=SylKikSYDH>.

- 756 Hongyu Ren, Hanjun Dai, Zihang Dai, Mengjiao Yang, Jure Leskovec, Dale Schuurmans, and  
757 Bo Dai. Combiner: Full attention transformer with sparse computation cost. In Marc’Aurelio  
758 Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan  
759 (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neu-  
760 ral Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*,  
761 pp. 22470–22482, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/  
762 bd4a6d0563e0604510989eb8f9ff71f5-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/bd4a6d0563e0604510989eb8f9ff71f5-Abstract.html).
- 763 Stephen E. Robertson and Hugo Zaragoza. The probabilistic relevance framework: BM25 and  
764 beyond. *Found. Trends Inf. Retr.*, 3(4):333–389, 2009. doi: 10.1561/1500000019. URL <https://doi.org/10.1561/1500000019>.  
765 //doi.org/10.1561/1500000019.
- 766 Ohad Rubin and Jonathan Berant. Retrieval-pretrained transformer: Long-range language modeling  
767 with self-retrieval, 2024. URL <https://arxiv.org/abs/2306.13421>.
- 769 Christopher Sciavolino, Zexuan Zhong, Jinhyuk Lee, and Danqi Chen. Simple entity-centric ques-  
770 tions challenge dense retrievers. In Marie-Francine Moens, Lucia Specia, and Scott Wen-tau  
771 Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language  
772 Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November,  
773 2021*, pp. 6138–6148. Association for Computational Linguistics, 2021. doi: 10.18653/V1/2021.  
774 EMNLP-MAIN.496. URL <https://doi.org/10.18653/v1/2021.emnlp-main.496>.
- 775 Philippe Tillet, Hsiang-Tsung Kung, and David D. Cox. Triton: an intermediate language and com-  
776 piler for tiled neural network computations. In Tim Mattson, Abdullah Muzahid, and Armando  
777 Solar-Lezama (eds.), *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine  
778 Learning and Programming Languages, MAPL@PLDI 2019, Phoenix, AZ, USA, June 22, 2019*,  
779 pp. 10–19. ACM, 2019. doi: 10.1145/3315508.3329973. URL [https://doi.org/10.1145/  
780 3315508.3329973](https://doi.org/10.1145/3315508.3329973).
- 781 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
782 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
783 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 784 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
785 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike  
786 von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Ro-  
787 man Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Confer-  
788 ence on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA,  
789 USA*, pp. 5998–6008, 2017. URL [https://proceedings.neurips.cc/paper/2017/hash/  
790 3f5ee243547dee91fbd053c1c4a845aa-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html).
- 791 Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention  
792 with linear complexity. *CoRR*, abs/2006.04768, 2020. URL [https://arxiv.org/abs/2006.  
793 04768](https://arxiv.org/abs/2006.04768).
- 794 Yuhuai Wu, Markus Norman Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing trans-  
795 formers. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual  
796 Event, April 25-29, 2022*. OpenReview.net, 2022. URL [https://openreview.net/forum?  
797 id=TrjbxzRcnf-](https://openreview.net/forum?id=TrjbxzRcnf-).
- 799 Howard Yen, Tianyu Gao, and Danqi Chen. Long-context language modeling with parallel context  
800 encoding. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd  
801 Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL  
802 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 2588–2610. Association for Computational  
803 Linguistics, 2024. URL <https://aclanthology.org/2024.acl-long.142>.
- 804 Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago  
805 Ontañón, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: Trans-  
806 formers for longer sequences. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-  
807 Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems  
808 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, De-  
809 cember 6-12, 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/  
hash/c8512d142a2d849725f31a9a7a361ab9-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/c8512d142a2d849725f31a9a7a361ab9-Abstract.html).

810 Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small  
811 language model, 2024.  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863



## A APPENDIX

### A.1 HYPER-PARAMETERS

**Long-Range Language Modeling.** We employ a Llama-like architecture (Touvron et al., 2023) featuring a 12-layer, decoder-only transformer with 12 heads per layer (64 dimensions each), an embedding dimension of 768, and an FFN size of 2048. Training utilizes the AdamW optimizer (Loshchilov & Hutter, 2019) with  $\beta_1 = 0.9$  and  $\beta_2 = 0.95$ , and a weight decay factor of 0.001. We used base learning rate  $2 \times 10^{-3}$  for all our experiments with a warmup stage that was 2% of the whole training and applied a cosine scheduler with final learning rate being  $4 \times 10^{-4}$ . We used GPT-2’s (Radford et al., 2019) tokenizer. We used mixed-precision training with bfloat16 over at 8 Nvidia A100 GPUs. We train all models with an effective batch size of  $2^{19}$  tokens for 60K steps resulting in a total training budget of 32.2 billion tokens. We train Base LM, RPT and DRT on each dataset with a context length of 16K tokens. Due to Landmark Attention doesn’t support sliding-window attention, the model is pre-trained with full self-attention with a context length of 768. Due to Block Recurrent Transformer cannot be fully paralleled, which takes  $5 \times$  wall-clock training time with 16K context length, we pre-train it with a context length of 4K.

### A.2 HARDWARE-AWARE GCA PSUEDO-CODE

---

#### Algorithm 1 FLASHGCA forward pass

---

**Require:** Matrices  $\mathbf{Q} \in \mathbb{R}^{N_q \times d}$ ,  $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{K \times N_{kv} \times d}$  in HBM, vector  $\mathbf{w} \in \mathbb{R}^k$  in HBM, block sizes  $B_c, B_r$ .

- 1: Divide  $\mathbf{Q}$  into  $T_r = \lceil \frac{N_q}{B_r} \rceil$  blocks  $\mathbf{Q}_1, \dots, \mathbf{Q}_{T_r}$  of size  $B_r \times d$  each, and divide  $\mathbf{K}, \mathbf{V}$  in to  $K \times T_c$  blocks where  $T_c = \lceil \frac{N_{kv}}{B_c} \rceil$   $\mathbf{K}_{1,1}, \dots, \mathbf{K}_{K,T_c}$  and  $\mathbf{V}_{1,1}, \dots, \mathbf{V}_{K,T_c}$ , of size  $B_c \times d$  each.
  - 2: Divide the output  $\mathbf{O} \in \mathbb{R}^{N_q \times d}$  into  $T_r$  blocks  $\mathbf{O}_1, \dots, \mathbf{O}_{T_r}$  of size  $B_r \times d$  each, and divide the logsumexp  $L \in \mathbb{R}^{N_q \times K}$  into  $T_r \times K$  blocks  $L_{1,1}, \dots, L_{T_r,K}$  of size  $B_r$  each.
  - 3: Divide the output  $\mathbf{O}' \in \mathbb{R}^{K \times N_q \times d}$  into  $T_r$  blocks  $\mathbf{O}'_{1,1}, \dots, \mathbf{O}'_{K,T_r}$  of size  $K \times B_r \times d$  each.
  - 4: **for**  $1 \leq i \leq T_r$  **do**
  - 5:   Load  $\mathbf{Q}_i$  from HBM to on-chip SRAM.
  - 6:   Load  $\mathbf{w}_k$  from HBM to on-chip SRAM.
  - 7:   **for**  $1 \leq k \leq K$  **do**
  - 8:     On chip, initialize  $\mathbf{O}_i^{(0)} = (0)_{B_r \times d} \in \mathbb{R}^{B_r \times d}$ ,  $\ell_i^{(0)} = (0)_{B_r} \in \mathbb{R}^{B_r}$ ,  $m_i^{(0)} = (-\infty)_{B_r} \in \mathbb{R}^{B_r}$ .
  - 9:     **for**  $1 \leq j \leq T_c$  **do**
  - 10:      Load  $\mathbf{K}_{k,j}, \mathbf{V}_{k,j}$  from HBM to on-chip SRAM.
  - 11:      On chip, compute  $\mathbf{S}_i^{(j)} = \mathbf{Q}_i \mathbf{K}_{k,j}^T \in \mathbb{R}^{B_r \times B_c}$ .
  - 12:      On chip, compute  $m_i^{(j)} = \max(m_i^{(j-1)}, \text{rowmax}(\mathbf{S}_i^{(j)})) \in \mathbb{R}^{B_r}$ ,  $\tilde{\mathbf{P}}_i^{(j)} = \exp(\mathbf{S}_i^{(j)} - m_i^{(j)}) \in \mathbb{R}^{B_r \times B_c}$  (pointwise),  $\ell_i^{(j)} = e^{m_i^{(j-1)} - m_i^{(j)}} \ell_i^{(j-1)} + \text{rowsum}(\tilde{\mathbf{P}}_i^{(j)}) \in \mathbb{R}^{B_r}$ .
  - 13:      On chip, compute  $\mathbf{O}_i^{(j)} = \text{diag}(e^{m_i^{(j-1)} - m_i^{(j)}})^{-1} \mathbf{O}_i^{(j-1)} + \tilde{\mathbf{P}}_i^{(j)} \mathbf{V}_{k,j}$ .
  - 14:     **end for**
  - 15:     On chip, compute  $\mathbf{O}'_{i,k} = \text{diag}(\ell_i^{(T_c)})^{-1} \mathbf{O}_i^{(T_c)}$ .
  - 16:     On chip, compute  $\mathbf{O}_i \leftarrow \mathbf{O}_i + \mathbf{w}_k \mathbf{O}'_{i,k}$ .
  - 17:     Write  $\mathbf{O}_i$  to HBM.
  - 18:     On chip, compute  $L_{i,k} = m_i^{(T_c)} + \log(\ell_i^{(T_c)})$ .
  - 19:     Write  $L_{i,k}$  to HBM.
  - 20:   **end for**
  - 21:   Write  $\mathbf{O}_i$  to HBM as the  $i$ -th block of  $\mathbf{O}$ .
  - 22: **end for**
  - 23: Return the output  $\mathbf{O}$  and the logsumexp  $L$ .
-

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

---

**Algorithm 2** FLASHGCA Backward Pass
 

---

**Require:** Matrices  $\mathbf{Q}, \mathbf{O}, \mathbf{dO} \in \mathbb{R}^{N_q \times d}$ ,  $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{K \times N_{kv} \times d}$ ,  $L \in \mathbb{R}^{N_q \times K}$ ,  $\mathbf{O}' \in \mathbb{R}^{K \times N_q \times d}$  in HBM, vector  $\mathbf{w} \in \mathbb{R}^K$  in HBM, block sizes  $B_c, B_r$ .

- 1: Divide  $\mathbf{Q}$  into  $T_r = \lceil \frac{N}{B_r} \rceil$  blocks  $\mathbf{Q}_1, \dots, \mathbf{Q}_{T_r}$  of size  $B_r \times d$  each, and divide  $\mathbf{K}, \mathbf{V}$  in to  $K \times T_c$ , where  $T_c = \lceil \frac{N}{B_c} \rceil$  blocks  $\mathbf{K}_{1,1}, \dots, \mathbf{K}_{K,T_c}$  and  $\mathbf{V}_{1,1}, \dots, \mathbf{V}_{K,T_c}$ , of size  $B_c \times d$  each.
  - 2: Divide  $\mathbf{O}$  into  $T_r$  blocks  $\mathbf{O}_1, \dots, \mathbf{O}_{T_r}$  of size  $B_r \times d$  each, divide  $\mathbf{dO}$  into  $T_r$  blocks  $\mathbf{dO}_1, \dots, \mathbf{dO}_{T_r}$  of size  $B_r \times d$  each, and divide  $L$  into  $T_r \times K$  blocks  $L_{1,1}, \dots, L_{T_r,K}$  of size  $B_r$  each.
  - 3: Initialize  $\mathbf{dQ} = (0)_{N_q \times d}$  in HBM and divide it into  $T_r$  blocks  $\mathbf{dQ}_1, \dots, \mathbf{dQ}_{T_r}$  of size  $B_r \times d$  each. Divide  $\mathbf{dK}, \mathbf{dV} \in \mathbb{R}^{K \times N_{kv} \times d}$  in to  $K \times T_c$  blocks  $\mathbf{dK}_{1,1}, \dots, \mathbf{dK}_{K,T_c}$  and  $\mathbf{dV}_{1,1}, \dots, \mathbf{dV}_{K,T_c}$ , of size  $B_c \times d$  each. Initialize  $\mathbf{dW} = (0)_{T_r \times K}$  in HBM.
  - 4: Compute  $D = \text{rowsum}(\mathbf{dO} \circ \mathbf{O}') \in \mathbb{R}^{N_q \times K}$  (pointwise multiply), write  $D$  to HBM and divide it into  $T_r$  blocks  $D_1, \dots, D_{T_r}$  of size  $B_r$  each.
  - 5: **for**  $1 \leq k \leq K$  **do**
  - 6:   Load  $\mathbf{w}_k$  from HBM to on-chip SRAM.
  - 7:   **for**  $1 \leq j \leq T_c$  **do**
  - 8:     Load  $\mathbf{K}_{k,j}, \mathbf{V}_{k,j}$  from HBM to on-chip SRAM.
  - 9:     Initialize  $\mathbf{dK}_{k,j} = (0)_{B_c \times d}$ ,  $\mathbf{dV}_{k,j} = (0)_{B_c \times d}$ ,  $\mathbf{dW}_{k,j} = (0)$  on SRAM.
  - 10:    **for**  $1 \leq i \leq T_r$  **do**
  - 11:     Load  $\mathbf{Q}_i, \mathbf{dO}_i, \mathbf{dQ}_i, D_i$  from HBM to on-chip SRAM.
  - 12:     Load  $L_{i,k}$  from HBM to on-chip SRAM.
  - 13:     On chip, compute  $\mathbf{S}_i^{(j)} = \mathbf{Q}_i \mathbf{K}_{k,j}^T \in \mathbb{R}^{B_r \times B_c}$ .
  - 14:     On chip, compute  $\mathbf{P}_i^{(j)} = \exp(\mathbf{S}_i^{(j)} - L_{i,k}) \in \mathbb{R}^{B_r \times B_c}$ .
  - 15:     On chip, compute  $\mathbf{dV}_{k,j} \leftarrow \mathbf{dV}_{k,j} + (\mathbf{w}_k \mathbf{P}_i^{(j)})^\top \mathbf{dO}_i \in \mathbb{R}^{B_c \times d}$ .
  - 16:     On chip, compute  $\mathbf{dP}_i^{(j)} = \mathbf{dO}_i \mathbf{V}_j^\top \in \mathbb{R}^{B_r \times B_c}$ .
  - 17:     On chip, compute  $\mathbf{dW}_{i,k} = \text{rowsum}(\mathbf{P}_i^{(j)} \circ \mathbf{dP}_i^{(j)})$ .
  - 18:     On chip, compute  $\mathbf{dS}_i^{(j)} = \mathbf{w}_k \mathbf{P}_i^{(j)} \circ (\mathbf{dP}_i^{(j)} - D_{i,k}) \in \mathbb{R}^{B_r \times B_c}$ .
  - 19:     Write  $\mathbf{dW}_{i,k}$  to HBM.
  - 20:     Load  $\mathbf{dQ}_i$  from HBM to SRAM, then on chip, update  $\mathbf{dQ}_i \leftarrow \mathbf{dQ}_i + \mathbf{dS}_i^{(j)} \mathbf{K}_j \in \mathbb{R}^{B_r \times d}$ , and write back to HBM.
  - 21:     On chip, compute  $\mathbf{dK}_{k,j} \leftarrow \mathbf{dK}_{k,j} + \mathbf{dS}_i^{(k,j)\top} \mathbf{Q}_i \in \mathbb{R}^{B_c \times d}$ .
  - 22:    **end for**
  - 23:    Write  $\mathbf{dK}_{k,j}, \mathbf{dV}_{k,j}$  to HBM.
  - 24:   **end for**
  - 25: **end for**
  - 26:  $\mathbf{dW} = \mathbf{dW}.\text{sum}(\text{dim} = 0)$
  - 27: Return  $\mathbf{dQ}, \mathbf{dK}, \mathbf{dV}, \mathbf{dW}$ .
-

## A.3 MORE CASE STUDIES

...An alternate approach is given below in Corollary \ref{cor:infnty}. Along the way we obtain more information about the eigenfunctions, which leads directly to an explicit formula for  $u_m(x; \infty)$ , see \eqref{conjsum2} and \eqref{Bkmexplicit}. As  $\sigma$  increases, the derivatives of  $u_m(x; \sigma)$  remain bounded, and so to ensure that the interior condition in \eqref{deltabc} continues to hold, the values  $u_m(x_k; \sigma)$  converges to infinity...

...must converge to 0 as  $\sigma$  converges to infinity. Our first corollary of Theorem \ref{thm:main} is that these values converge to 0 at the same rate for each node  $x_k$ . \begin{cor} \label{cor:nodes}. Up to an overall normalization factor, for each  $\sigma \geq 0$  and  $1 \leq m \leq n-1$ , the values of the eigenfunctions...

...To obtain the limiting eigenfunctions, which we denote by \label{def:uminfnty}  $\nu_m(x; \infty) = \lim_{\sigma \rightarrow \infty} u_m(x; \sigma)$ , one can use the fact that  $\gamma_m(\sigma) \rightarrow n\pi$  for  $1 \leq m \leq n$  to obtain...

...the eigenvalues  $la_m(\sigma)$  for  $1 \leq m \leq n$  all converge to  $la_n = n^2\pi^2$  as  $\sigma$  tends to infinity. (Note that this is consistent with our implicit expression for the eigenvalues  $la_m(\sigma)$  from Theorem \ref{thm:main}.) From Corollary \ref{cor:nodes}, this ensures that  $u_m(x_k; \sigma)$  converges to zero as  $\sigma$  tends to infinity. This means that  $u_m(x; \infty)$  (defined in \eqref{def:uminfnty}) is proportional...

Figure 5: In the case above, retrieved top-1 chunk introduces the definition used in the target chunk, while the adjacent retrieved 3rd best chunk and retrieved 4th best chunk both cover the same variants as those appear in the target chunk. Retrieved 2nd best chunk contains the theorem and corollary used in the query chunk.

... denote the projection  $\pi : \mathbb{R}^{pvc} \rightarrow \mathbb{C}^{pvc}$ . \begin{lemma} \label{lemma:mu circ pi=2n} The complex and real moment maps for  $G^{\mathbb{C}}$  are related by  $\mu^* \circ \pi = 2n$  \end{lemma} \begin{proof} Many of our computations...

...  $\pi(\omega([v])) = \omega(\pi[v])$  where  $\omega(p)$  denotes the  $\omega$ -limit set of the negative gradient flow starting from  $p$ . \end{prop} \begin{proof} Applying Lemma \ref{lemma:mu circ pi=2n} we have  $4 < grad||n||^2[v], w_{[v]} > = 4$  ...

Figure 6: In the case above, retrieved top-1 chunk introduces the lemma used in the target chunk.

... They are not the same: see Section \ref{sec:15}. To establish Theorem \ref{thm:ABn} it suffices to prove it for  $B(n)$ ; the estimate for  $A(n)$  then follows from the linear relation  $A(n) = \log G_n + B(n)$  (from \eqref{eqn:GABx}) combined with the asymptotic estimate for  $G(n)$  in \eqref{eqn:logG-asymp}. The main contribution in the sum  $B(n)$  comes from those primes  $p$  having  $p > \sqrt{n}$ , whose key property ...

... that exponential sum methods yield alternative unconditional estimates for  $A(n, x)$ ,  $B(n, x)$  and  $\log G(n, x)$ , which are nontrivial when  $x = o(n)$ , and apply for  $x > \sqrt{n}$ . These estimates improve on the estimates of our main theorems for certain ranges ...

... exponents  $\nu_p(G_n)$  as a difference of quantities given by statistics of the base  $p$  radix expansion of integers up to  $n$  (see Theorem \ref{thm:explicit}). Summing over  $p \leq x$  yields a formula  $\log G(n, x) = A(n, x) - B(n, x)$  involving nonnegative arithmetic functions  $A(n, x)$  and  $B(n, x)$  ...

... The implied constant in the  $O$ -notation does not depend on  $\alpha$ . \end{thm} The limit function  $f_B(\alpha)$  is pictured in Figure \ref{fig:B2}. The function lies strictly above the diagonal line  $\beta = (1 - \gamma)\alpha$ ; note that in \eqref{eqn:GABx} in its relation to  $\log G(n, x)$  it appears with a negative sign, consistent with  $f_G(\alpha)$  ...

Figure 7: In the case above, retrieved top-1 chunk and retrieved 3rd best chunk are adjacent, which mentions the same equation as target chunk. retrieved 2nd best chunk and retrieved 4th best chunk both mention  $\log G(n, x)$ , which also appears in target chunk.

#### A.4 THE DETAILS FOR THE NIAH TEST

In all evaluations conducted for the NIAH tests, we fine-tune all models using checkpoints derived from PG19, employing the same set of synthetic data. The number of fine-tuning steps is set to one-tenth of the total steps used during pre-training, while all other hyperparameters are kept constant. Examples of the synthetic data utilized for each task are presented in the table 5. Specifically, we pad the input tokens to ensure that the landmark token can be inserted before “is” in the question.

Task	Example
Single NIAH	(essays)...
	The passkey is: {tokens}.
2-hop NIAH with noises	...
	What is the passkey? The passkey is {tokens}.
	(essays)...
	DEF {tokens_5}->{tokens_6} ...
	DEF {tokens_2}->{tokens_3} ...
DEF {tokens_4}->{tokens_5} ...	
DEF {tokens_1}->{tokens_2} ...	
...	The path from {tokens_1} is: {tokens_2}, {tokens_3}

Table 5: Task examples for the two NIAH tests.

#### A.5 THE DETAILS FOR INFINITELY LONG CONTEXT RETRIEVAL

We pre-train  $DRT_{enc-dec}$  with 200M parameters on MiniPile and store all landmark representations in a Faiss as described in § 3.3 to emulate an infinite context. The trained model has an embedding dimension of 1,024, and an FFN size of 2,816. We train DRT on MiniPile for 20 epochs with 384K tokens per batch. Specifically, we prepare DRT with different settings.  $DRT_{w/random\ retrieval}$  uses a randomly generated vector for retrieval.  $DRT_{w/Contriever}$  utilizes a pre-trained retriever with fixed parameters to select top- $k$  relevance chunks, with information still fused via GCA.

#### A.6 DISCUSSIONS ABOUT THE $RPT_{CONTRIVER}$ BASELINE

In the original RPT, a reference LM is used to pre-prepare target chunks for each chunk, as discussed in the related works. Compared to using Contriever as the retrieval module, the original method offers stronger causal retrieval capabilities. However, since the code for retriever distillation in the original RPT is not released and the approach is costly and less flexible, we opt to use Contriever instead.  $RPT_{contriever}$  can be considered a fusion of RETRO and RPT. It retrieves past chunks in a manner similar to RETRO and integrates the retrieved information in the style of RPT. The retrieval process involves dividing every 64 tokens into a chunk, encoding them with Contriever to obtain chunk representations, and then retrieving past chunks based on cosine similarity with the current chunk.