

# CacheFormer: High Attention-based Segment Caching

## Abstract

Efficiently handling long contexts in transformer-based language models with low perplexity is an active area of research. Although, numerous approaches have been recently presented like Linformer, Longformer, Performer, Structured state space models (SSMs) etc., yet it remains an unresolved problem. All these models strive to reduce the quadratic time complexity of the attention mechanism to approximate linear time complexity while minimizing the loss in quality due to the effective compression of the long context. Inspired by the cache and virtual memory concepts in computer architecture, we improve the work presented in Long-Short Transformer (Transformer-LS) that implements a sliding window for the short attention and compressed contextual segments for the long attention. Our enhancements include augmenting the architecture with attention on dynamically retrieved uncompressed context segments that indicate high attention at the compressed level. Similar to the cache and virtual memory principle in computers, where in case of a cache or page miss, not only the needed data is retrieved from the random-access memory or the hard disk, but the nearby following data is also obtained. On a similar note, we too retrieve the nearby segments in uncompressed form when a high attention occurs at the compressed level. We further enhance the long-short transformer by augmenting the long attention with compressed overlapping segments to reduce the loss in quality due to segment fragmentation that occurs in sequences with long context. Our results indicate significant improvements over the base line of the long-short transformer in terms of perplexity on the popular benchmarks.

## 1 Introduction

Deep Convolutional Neural Networks (CNNs) were fundamental in revolutionizing the field of computer vision. Similarly, the pioneering induction of the Transformer (Vaswani et al., 2017) architecture in Natural Language Processing (Singh and Mahmood, 2021) has resulted in the AI revolution with Large

Language Models (LLMs) such as ChatGPT (J. Achiam et al., 2023), Bard (G. Team et al., 2023), Llama (H. Touvron et al., 2023) among others have yielded impressive performances. The Transformer uses a simple similarity computation in the form of an inner product on the learnt positional encoded embeddings of a sequence of  $n$  input words. If the matrix  $Q$  and  $K$  contain rows representing embedding of each word ( $1 \times d$ ), then  $A = \text{softmax}(QK^T)$  referred to as the “attention”, contains the dot product similarity of each input word with every other word in the input sequence. If there are  $n$  words being input, referred to as the context, then  $Q, K \in \mathbb{R}^{n \times d}$ , and  $A \in \mathbb{R}^{n \times n}$ . Like parallel feature maps in a CNN, each layer in the Transformer divides the attention calculation into parallel heads. The output from a Transformer layer has the same dimensionality as input and is obtained by a simple matrix computation of  $(A \times V) \in \mathbb{R}^{d \times n}$  where  $V \in \mathbb{R}^{n \times d}$  is similar to  $K$  and contains rows of learnt position encoded embeddings of input words. For language models, where text generation is carried out based on a given context, the attention matrix is masked in a triangular fashion so that future tokens are not visible in the training process. Multiple layers of Transformer blocks are used before feeding the result of the last layer to a classification head. Because attention computation in each head is  $O(n^2)$ , for long contexts, this becomes a computational bottleneck. Many approaches have been proposed in the last few years to reduce the quadratic time complexity of attention to either linear or sub quadratic complexity. Some of the notable works include (Dai Z et al., 2019), (Wang et al., 2020), (Beltagy et al., 2020), (Kitaev et al., 2020), (Choromanski et al., 2021), (Hawthorne et al., 2022), (H. Ji et al., 2022), (Martins et al., 2022) among others. We provide a brief background in the above-mentioned approaches used in reducing the attention complexity. Then we elaborate on the

93 Long-Short Transformer that we will further  
94 enhance in this work.

## 95 **2 Background and Related Work**

96 An important earlier work in handling long  
97 contexts was presented by (Dai Z et al., 2019).  
98 The authors divided the context into segments  
99 and used segment level recurrence and a  
100 corresponding positional encoding to allow it to  
101 handle longer contexts. It achieved impressive  
102 results on the perplexity and BPC at that time.  
103 (Wang et al., 2020) accomplished  $O(n)$   
104 complexity through linear self-attention. The  
105 authors demonstrate that the attention is  
106 typically low rank, and thus can be  
107 approximated by a low rank matrix. Here, the Q  
108 and V matrices  $\in \mathbb{R}^{n \times d}$  are projected to lower  
109 dimension matrices  $\in \mathbb{R}^{k \times d}$  where  $k < n$ . Thus  
110 attention  $A = QK^T \in \mathbb{R}^{n \times k}$ . The output  
111  $(A \times V) \in \mathbb{R}^{n \times d}$ , i.e., same as the original  
112 transformer. Since  $k$  is fixed, the attention  
113 complexity is  $O(n)$ .

114 Although (Wang et al., 2020) reduced the  
115 attention complexity significantly, especially if  
116  $k \ll n$ , note that, it cannot be effectively  
117 used in autoregressive training and generation,  
118 as the projection of Q compresses the  
119 information along the context, making the  
120 masking of attention for future tokens invalid.  
121 However, for classification problems where  
122 masking of attention is not needed, their  
123 architecture is effective in reducing complexity.

124  
125 Another approach introduced by (Beltagy et al.,  
126 2020) used sparse attention patterns instead of  
127 the full dense attention. The authors proposed  
128 sliding window attention, where tokens  
129 attended only to the nearby past, a dilated  
130 sliding window, and a mix of global and sliding  
131 window attention where some tokens attend to  
132 all tokens while others only attend to nearby  
133 tokens. For autoregressive modeling (Beltagy et  
134 al., 2020) used dilated sliding window  
135 attention. Another notable work in reducing the  
136 attention complexity was performed by (Kitaev  
137 et al., 2020). The authors key idea was to use  
138 locality sensitive hashing which reduces the  
139 attention complexity to  $O(n \log n)$ . Note that  
140 because of the hashing process, the architecture  
141 is not suited for autoregressive modeling.

142 A different approach to reduce the attention  
143 complexity was taken by (Choromanski et al.,  
144 2021) where the attention is decomposed as a  
145 product of non-linear functions of original  
146 query and key matrices referred to as random  
147 features. This allows the attention to be  
148 encoded more efficiently via the transformer  
149 query and key matrices. Further efficient  
150 handling of long contexts accomplished by  
151 (Hawthorne et al., 2022) divided the input  
152 sequence into smaller key/value and query  
153 components. These components underwent  
154 cross attention in the first layer with a latent  $\in$   
155  $\mathbb{R}^{l \times d}$  where  $l$  is the size chosen in splitting the  
156 input sequence into the query part. The  
157 remaining layers operate on the  $l \times d$  size  
158 instead of the usual  $n \times d$  size as in a standard  
159 transformer. Although this cross attention on  
160 the partitioned input sequence results in  
161 efficient handling of long sequences, because  
162 of the reduced query size, the equivalent effect  
163 is more like a sliding window attention.

164 More recently, a different approach to handling  
165 long contexts was proposed via structure state  
166 space models. The work by (A. Gu et al., 2022)  
167 proposes the Structured State Space Sequence  
168 model (S4) based on a new parameterization  
169 that can be computed much more efficiently. A  
170 variation of the state space approach proposed  
171 by (X. Ma et al., 2023) uses single-head gated  
172 attention mechanism equipped with exponential  
173 moving average to incorporate inductive bias of  
174 position-aware local dependencies into the  
175 position-agnostic attention mechanism. They  
176 also present its variation with linear time  
177 complexity for handling long sequences.  
178 Further progression on the state space models  
179 yielded better results (D. Y. Fu et al., 2023),  
180 (Gu, Albert and Tri Dao., 2023) who achieved  
181 a very low perplexity score. Most recently  
182 (Maximilian et al., 2024) introduced  
183 exponential gating and parallelization in  
184 LSTMs to achieve extended memory. Some of  
185 the model sizes consisted of several billion  
186 parameters. We outperform the smaller version  
187 of these models with similar size as ours on the  
188 perplexity metric as shown in Table 2.

189 An interesting concept in handling long  
190 sequences was presented by (C. Zhu et al.,  
191 2021). Here a sliding window approach is used  
192 in handling near term attention, while a set of

193 compressed segments for the entire past context  
 194 is used as long-term attention. Both short and  
 195 long attention are combined in the overall  
 196 attention. The slight drawback of the approach  
 197 is that the longer context is effectively used in  
 198 compressed form and thus may lose some key  
 199 contextual information in being able to generate  
 200 the output in an autoregressive environment.  
 201 We address this problem by further augmenting  
 202 the long-short attention by using uncompressed  
 203 highly attentive segments. Since long short  
 204 attention divides the context into equal size  
 205 segments before projecting each segment to a  
 206 smaller size, there is potential for a loss in  
 207 information due to segment fragmentation. We  
 208 also improve this aspect by using overlapping  
 209 segments and augment this to the existing long-  
 210 short model. Thus, our enhanced long-short  
 211 architecture involves four components in the  
 212 overall attention, a sliding window attention,  
 213 long attention based on compressed segments,  
 214 long attention based on overlapping segments,  
 215 and uncompressed segmented attention for few  
 216 high attentive segments beyond the sliding  
 217 window part. We describe the details of our  
 218 design in the section 3. For completeness, we  
 219 summarize the composition of a Transformer,  
 220 followed by the ideas of long-short  
 221 Transformer, that we build upon in our work.

## 2.1 Canonical Transformer

223 In normal multi-headed attention, if  $Q, K, V \in$   
 224  $\mathbb{R}^{n \times d}$  are the query, key and value  
 225 transformations of the input embeddings with  
 226 sequence length of  $n$  and embedding dimension  
 227 of  $d$ , then the scaled dot-product attention in the  
 228  $i$ -th Head  $H_i \in \mathbb{R}^{n \times d_k}$  is given as:

$$229 H_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) =$$

$$230 \text{Softmax}\left[\frac{QW_i^Q(KW_i^K)^T}{\sqrt{d_k}}\right] VW_i^V = A_i VW_i^V \quad (1)$$

231 Where  $d_k = d/h$  is the dimension of each  
 232 head. The output in each transformer layer is  
 233 obtained by catenation of the output of all  
 234 heads and transformed further via this  
 235 projection matrix.

$$236 W^o \in \mathbb{R}^{d \times d} \text{ as } \text{Layer}_j =$$

$$237 \text{Concat}(H_0, H_1, \dots, H_{h-1})W^o \quad (2)$$

238

239 After feeding the embedding of a sequence of  
 240 one hot encoded words,  $x$  (with position  
 241 encoding  $PE$  added) through  $p$  transformer  
 242 layers, a classification layer is used at the output  
 243 of the last layer to decide the output produced by  
 244 the transformer. For autoregressive text  
 245 generation, the classification layer's final output  
 246 is equal to the size of the dictionary of unique  
 247 words in the corpus.

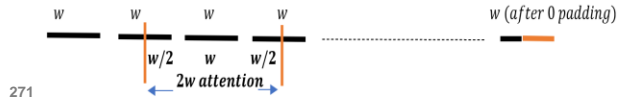
$$248$$

$$249 \text{out} = \text{classifier}[\text{layer}_{p-1}(\text{layer}_{p-2}(\dots \text{layer}_0$$

$$250 (\text{embedding}(x) + PE(x)))] \quad (3)$$

## 2.2 Long Short Transformer

253 (C. Zhu et al., 2021) aggregated the local  
 254 attention around a smaller window (sliding  
 255 window), with a projection of the full sequence  
 256 attention to a smaller size, so that we can  
 257 efficiently handle long sequences without the  
 258 quadratic attention complexity. For short  
 259 attention, the approach here is to use a segment  
 260 level sliding window attention, where the input  
 261 sequence is divided into disjoint segments with  
 262 length  $w$  (e.g.,  $w=128$  and sequence length is  
 263 1024). For non-autoregressive applications, all  
 264 tokens within a segment attend to all tokens  
 265 within its home segment, as well as  $w/2$   
 266 consecutive tokens on the left and right side of  
 267 its home segment (zero-padding when  
 268 necessary), resulting in an attention span over a  
 269 total of  $2w$  key-value pairs. This is depicted in  
 270 Figure 1.



271

272 Figure 1. Segment-based Sliding Window Attention

273

274 For each query  $Q_t$  at the position  $t$  within the  $i$ -th  
 275 head, the  $2w$  key-value pairs within its window  
 276 are:  $\tilde{K}_t, \tilde{V}_t \in \mathbb{R}^{2w \times d}$ . The short attention  $\bar{A}_{S_i} \in$   
 277  $\mathbb{R}^{2w \times d_k}$  is then given by the following equation:

$$278$$

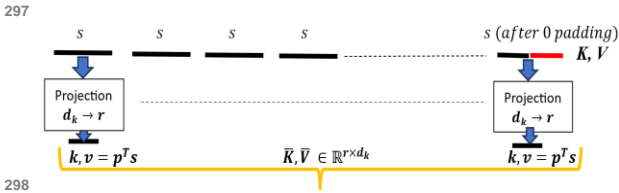
$$279 \bar{A}_{S_i} = \text{softmax}\left[\frac{QW_i^Q \tilde{K}_t^T}{\sqrt{d_k}}\right] \quad (4)$$

280

281 Execution wise the segment-level sliding  
 282 window attention (referred to as short attention)  
 283 is more time efficient than the per-token sliding  
 284 window attention where each token attends to  
 285 itself and  $w$  tokens to its left and right, and its  
 286 memory consumption scales linearly with

287 sequence length. For auto-regressive  
 288 applications, the future tokens in the current  
 289 segment are masked, and only the previous  
 290 segment is used.

291 For long attention, the key and value  
 292 transformations for the input sequence are first  
 293 divided into segments of fixed size  $s$ , and then  
 294 projected to a smaller dimension  $r$ , where the  
 295 projection  $P_{l_i} \in \mathbb{R}^{n \times r}$ . Figure 2 depicts this  
 296 process.



299 Figure 2. Segmented Long Attention with Compressed  
 300 Segments

302 Mathematically, the long attention  $\bar{A}_{l_i}$  (in each  
 303 head  $i$ ) as followed by the long-short Transformer  
 304 can be described as

$$305 P_{l_i} = \text{Softmax}(KW_i^P), \bar{K}_{l_i} = P_{l_i}^T KW_i^K, \bar{V}_{l_i} = 306 P_{l_i}^T VW_i^V \quad (5)$$

$$307 \bar{A}_{l_i} = \text{softmax} \left[ \frac{QW_i^Q \bar{K}_{l_i}^T}{\sqrt{d_k}} \right] \quad (6)$$

309 The output of in the  $i^{th}$  head is:

$$310 \bar{H}_i = \bar{A}_{l_i} (P_{l_i}^T VW_i^V) \quad (7)$$

312 Note that the long attention is effectively done on  
 313 a compressed form of  $K$  and  $V$ , as the projection  
 314 causes the input sequence of size  $n$  to be  
 315 compressed to size  $r$ . This results in full attention  
 316 to now be replaced with the implicit product of  
 317 two low-rank matrices  $\bar{P}_{l_i}^T \in \mathbb{R}^{r \times n}$  and  $QW_i^Q \in$   
 318  $\mathbb{R}^{n \times d}$ , and thus the computational complexity is  
 319 of long attention is reduced from  $O(n^2)$  to  
 320  $O(rn)$ .

321 Long-Short Transformer integrates the short and  
 322 long attentions into a single attention. While the  
 323 short attention can attend to most recent input,  
 324 the long attention is in compressed form. Further,  
 325 the long attention is based on segmentation of the  
 326 input sequence that may suffer from segment  
 327 fragmentation as the information in each segment  
 328 is compressed via the projection mechanism. We  
 329 improve upon these shortcomings and present  
 330 our enhanced long short architecture in the

331 following section. Our contributions can be  
 332 summarized as:

- 333 1. Improving the segment fragmentation of the  
 334 projection mechanism followed in long  
 335 attention by adding projections of segments  
 336 that have an  $s/2$  overlap where  $s$  is the  
 337 segment size, with the existing segment-based  
 338 projection mechanism.
- 339 2. Since the long attention is based on an  
 340 effective compression of the input sequence,  
 341 we develop an innovative uncompressed  
 342 attention mechanism where some of the  
 343 highly attentive segments are used  
 344 dynamically in uncompressed form.
- 345 3. The existing short and long attention are  
 346 combined with our two enhancements in an  
 347 effective manner to result in an architecture  
 348 that can efficiently handle long attentions  
 349 without causing much loss of attention  
 350 information.

### 3. Enhanced Long-Short Transformer

353 The long-term attention in the existing Long-  
 354 Short Transformer is done at a compressed level  
 355 (projection to  $r$  causes an effective compression  
 356 of the input context). Therefore, one of our  
 357 enhancements is to augment the long attention  
 358 with an attention that is based on a subset of  
 359 highly attentive uncompressed segments.

#### 3.1 Enhanced Long Attention with Segment Caching

362 The subset of segments that are selected for  
 363 attention at the uncompressed level is completely  
 364 dynamic and obtained by the vector magnitude  
 365 of the compressed segment-wise attention. In  
 366 simple words, we examine the segment-wise  
 367 long attention  $\bar{A}_{l_i}$  as given by Equation 6. Since  
 368  $\bar{A}_{l_i} \in \mathbb{R}^{n \times r}$ , and if there are  $n_s$  segments, then  
 369 each row in  $\bar{A}_{l_i}$  contains a set of row vectors of  
 370 size  $r/n_s$ , as denoted by segmented attention  
 371  $\bar{A}_{seg_i}$  in Equation 8. Magnitude of each vector  
 372  $\vec{a}_{i,j} \in \mathbb{R}^{1 \times r/n_s}$  in Equation 8, indicates the  
 373 attention of word  $i$  to the  $j^{th}$  segment in the long  
 374 attention.

$$375 \bar{A}_{seg_i} = \begin{bmatrix} \vec{a}_{1,1} & \vec{a}_{1,2} & \dots & \dots & \vec{a}_{1,n_s} \\ \vdots & \vdots & \dots & \dots & \vdots \\ \vec{a}_{n,1} & \vec{a}_{n,2} & \dots & \dots & \vec{a}_{n,n_s} \end{bmatrix} \quad (8)$$

376 For execution efficiency, we average the  
 377 segment attention vectors in  $p$  consecutive rows  
 378 resulting in a segment attention matrix  
 379  $A_{segavg_i} \in \mathbb{R}^{m \times n_s}$  where  $m = n/p$ . Then we  
 380 choose top  $k$  segments by magnitude of each  
 381 vector in each row of the segment attention  
 382 matrix  $A_{segavg_i}$

$$383 \quad A_{segavg_i} = \begin{bmatrix} \overrightarrow{\text{topk}\left(\sum_{t=1}^p \bar{A}_{S_i}[t,:]/p\right)} \\ \overrightarrow{\text{topk}\left(\sum_{t=p+1}^{2p} \bar{A}_{S_i}[t,:]/p\right)} \\ \dots \\ \overrightarrow{\text{topk}\left(\sum_{t=n-p}^n \bar{A}_{S_i}[t,:]/p\right)} \end{bmatrix} \quad (9)$$

384 Note that each entry in the segment attention  
 385 matrix,  $A_{segavg_i}[i,j]$ , indicates the segment  
 386 number that has high attention to the sequence of  
 387  $p$  words (positioned from  $(i-1)xp$  to  $ixp$ ) in  
 388 the input context. Rather than using these  
 389 attentive segments in compressed form, we  
 390 extract them from the segmented  $K$  and  $V$   
 391 matrices before doing any compression on these.  
 392 Similar to how in cache memory design (in  
 393 computer architecture), in case of a cache miss,  
 394 we not only retrieve the needed data from the  
 395 RAM, but also bring a few consecutive following  
 396 words, as there is high probability that these may  
 397 be needed in the near future. In case of segments  
 398 that we determine most attentive (by the top  $k$   
 399 order), we also retrieve  $u$  consecutive segments.  
 400 To clarify our approach, if the sequence length is  
 401  $n = 1024$ , and long attention segment size = 16,  
 402 then there will be 64 segments in the  
 403 uncompressed  $K$  and  $V$  matrices. If the projection  
 404 size  $r = 256$  (ratio of  $1024/256=4$ ), then each  
 405 segment of size 16 will be compressed to size of  
 406 4, resulting in long attention matrix  $\bar{A}_i$  of size  
 407  $1024 \times (64 \times 4)$  i.e.,  $1024 \times 256$ . If we choose to  
 408 average  $p=32$  consecutive rows in  $\bar{A}_i$ , and take  
 409 the magnitude of each of the  $1 \times 4$  vectors in each  
 410 row (corresponding to the 64 segments), then the  
 411 segment attention matrix  $A_{segavg_i}$  will be  $32 \times 64$ .  
 412 Taking the index of top  $k$  entries in each row of  
 413  $A_{segavg_i}$  will give us the index of most attentive

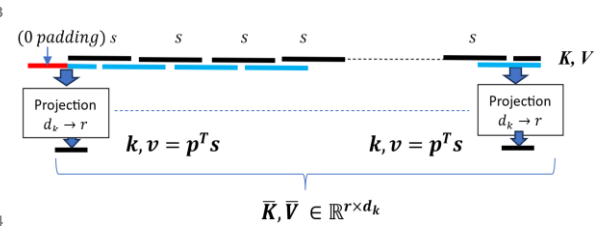
414  $k$  segments to the corresponding set of 32 words  
 415 in the input sequence. Assembling these top  $k$   
 416 attentive segments, and one segment before and  
 417 one segment after the attentive segment (if  $u=3$ ),  
 418 will result in 15 segments per row. If  $k=5$  is  
 419 chosen in  $topk$  and  $u=3$  which indicates using of  
 420  $u-1$  many nearby segments for each attentive  
 421 segment. Thus, the cache  $K, V$  matrices  $K_c, V_c \in$   
 422  $\mathbb{R}^{(n/p) \times (k \times u)}$  (e.g.,  $32 \times (15 \times 16) = 32 \times 240$  in our  
 423 example) contain the most attentive 15 segments  
 424 in uncompressed form. From the most attentive  
 425  $kxu$  segments in  $K_c$ , we can obtain the cache  
 426 attention  $\bar{A}_{c_i} \in \mathbb{R}^{n \times (k \times u)}$  as,

$$427 \quad \bar{A}_{c_i} = \text{softmax}\left(QW_i^Q \begin{bmatrix} K_c \\ K_c \\ \vdots \\ K_c \end{bmatrix}^T\right) / \sqrt{d_k} \quad (10)$$

428 Note that we stack the  $K_c$   $p$  times to match the  
 429 dimensionality with  $Q$ .

### 430 3.2 Enhanced Long Attention with 431 Overlapping Segments

432 In addition to the original long attention in the  
 433 Long-Short Transformer that uses the projections  
 434 on each segment, we augment the existing long  
 435 attention by using overlapping segments (with  
 436 50% overlap in augmented long attention) as  
 437 shown in Figure 3. The motivation behind the  
 438 overlap is to reduce the effect of segment  
 439 fragmentation in long attention. Zero padding in  
 440 the beginning segment is added to ensure the  
 441 same dimensionality for the overlapped long  
 442 segment attention.



443  
 444 **Figure 3.** Overlapping Segmented Long Attention  
 445 with Compressed Segments

447 The overlapped long segment attention  $\bar{A}_{o_i} \in$   
 448  $\mathbb{R}^{n \times r}$  similar to Equation 5 is given below.

$$449 \quad P_{o_i} = \text{Softmax}(KW_i^{P_o}), \bar{K}_{o_i} = P_{o_i}^T K W_i^K \\ 450 \quad \bar{V}_{o_i} = P_{o_i}^T V W_i^V \quad (11)$$

$$\bar{A}_{o_i} = \text{Softmax} \left[ \frac{QW_i^Q \bar{K}_{o_i}^T}{\sqrt{d_k}} \right] \quad (12)$$

$$\bar{H}_{o_i} = \bar{A}_{o_i} (P_{o_i}^T V W_i^V) \quad (13)$$

### 3.3 Aggregated Long-Short Attention

The final attention in our enhanced architecture is obtained by aggregating the four attentions described earlier, i.e., the short attention  $\bar{A}_{s_i} \in \mathbb{R}^{n \times 2w}$  that uses segment-wise sliding window, the segment based compressed long attention  $\bar{A}_{l_i} \in \mathbb{R}^{n \times r}$  as proposed by (C. Zhu et al., 2021). Our cache attention  $\bar{A}_{c_i} \in \mathbb{R}^{n \times (k \times u \times s)}$  is based on uncompressed high attention segments, and overlapping segment-based compressed attention,  $\bar{A}_{o_i} \in \mathbb{R}^{n \times r}$ . We add the two long and overlapping attentions,  $\bar{A}_{l_i}$  and  $\bar{A}_{o_i}$ . Thus, the final enhanced attention  $A_{e_i} \in \mathbb{R}^{n \times f}$  is:

$$A_{e_i} = [\bar{A}_{s_i} \parallel (\bar{A}_{l_i} + \bar{A}_{o_i}) \parallel \bar{A}_{c_i}] \quad (14)$$

where  $\parallel$  indicates the catenation of different attentions, and  $f = 2w + r + (k \times u \times s)$ ,  $w$  is the window size in short i.e., sliding window attention,  $r$  is the projection size in compressing the long attention,  $k$  is the *top k* factor in retrieving high attention top  $k$  segments,  $u-1$  is the number of neighboring segments to retrieve for cache attention,  $s$  is the segment size in long attention. For example, for *top k* of 5 and  $u = 3$ , segment size in short attention,  $w = 128$ , segment size in long attention = 16,  $r = 512$ , for an input sequence length of 2048, the size of our combined attention matrix is 2048x762.

## 4. Results

We use the long-short transformer (C. Zhu et al., 2021) as the baseline architecture. Instead of focusing on the absolute best results for perplexity and BPC, which often are achieved through extremely refined training schedules and large model sizes, we focus on the improvements over the baseline. Therefore, the results we show are more accurate reflection of the architectural improvements of our design. The baseline architecture is also programmed by us, and the enhancements we propose are programmed in the same implementation and can be selectively

turned on or off to see the contribution of each enhancement. We also use similar training schedules for the different architectures being compared. Table 1 shows the perplexity results for wikitext-103 dataset. It uses sequence length of 1024, short attention segment size of 128, long attention segment size of 16, compression of the long sequence by a factor of 4, i.e.,  $r=256$ , and different values of  $k$  in top k cache attention, and neighboring segments retrieval  $u$  of 1 or 3 (which indicates the segment before the attentive segment, and the one after it is also retrieved).

Model	Model Size	Perplexity
Long-Short Baseline	122.52 million	23.74
Enhanced Long-Short ( $k=3, u=1$ )	122.52 million	23.31
Enhanced Long-Short ( $k=5, u=1$ )	122.52 million	22.75
Enhanced Long-Short ( $k=7, u=1$ )	122.52 million	21.32
Enhanced Long-Short ( $k=5, u=3$ )	122.52 million	21.26

**Table 1.** Perplexity results Comparing the Baseline and our Enhanced Architecture

Note that our enhanced architecture does not cause any increase in the number of model parameters over the baseline long short Transformer. The models used for results in Table 1 have 12 layers, 12 heads, and an embedding size of 768 (for all architectural variations). For a sequence length of 1024 (which is same as used in GPT-2), using 7 segments ( $k=7, u=1$ ) yielded considerable improvement in perplexity. Increasing  $k$  beyond 7 did not seem to considerably reduce perplexity further. Since we have two major enhancements of cache attention and overlapping segment-based attention over the baseline, Table 2 shows an ablation study of the effects of each architectural improvement. Figure 4 depicts the 64 attention vectors for each segment (from compressed long attention, after averaging  $p=256$  rows) corresponding to the 64 segments during the beginning of training. The highest top  $k$  magnitude vectors then determine the segment to use in uncompressed form for our cache attention.

535 Table 3 shows the BPC results on the enwik-8  
 536 benchmark. The 23 million model uses 8 layers,  
 537 8 heads and embedding size of 512. The 34.88  
 538 million models used 12 layers. It is interesting to  
 539 note that the relative improvement in BPC by our  
 540 enhanced architecture is less pronounced as  
 541 compared to the perplexity improvements. This  
 542 could be attributed to the fact that majority of  
 543 improvements are attributed to cache attention  
 544 which uses a few highly attentive uncompressed  
 545 segments in long attention.

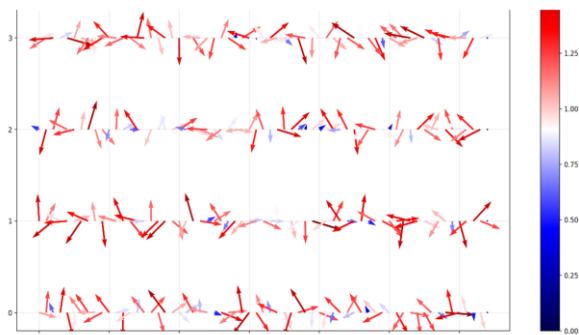
546

Architecture	Model Size (Millions)	Perplexity
Long-Short (Baseline-Ours)	122.52	23.74
Transformer-XL (Standard)	151	24
$\infty$ -former	160	24.22
LaMemo	151	23.77
H3 (Hungry Hungry Hippos)	125	23.7
Llama	125	23.16
Mamba	125	22.49
xLSTM[7:1]	125	<u>21.47</u>
Enhanced Long Short with overlapping segments only	122.52	23.47
Enhanced Long Short with cache attention only ( $k=7, u=1$ )	122.52	21.67
Enhanced Long Short with overlapping segments and cache attention ( $k=7, u=1$ )	122.52	<b>21.32</b>

547 **Table 2.** Ablation Study of Architectural  
 548 Enhancements

549  
 550 While this benefits the perplexity which is a  
 551 measure of the model’s prediction capability, but  
 552 BPC not as much, as BPC is more of a  
 553 compression efficiency measure of the model.

554  
 555



557 **Figure 4.** Attention Vectors from Compressed Long  
 558 Attention

Model	Model Size	BPC
Long-Short Baseline	23 million	1.192
Enhanced Long-Short ( $k=7, u=1$ )	23 million	1.188
Long-Short Baseline	34.88 million	1.173
Enhanced Long-Short ( $k=7, u=1$ )	34.88 million	1.167

559 **Table 3.** Comparison of BPC on the enwik-8  
 560 Benchmark

## 561 5. Discussion

562  
 563 Since the uncompressed segments to be used in  
 564 our cache attention design are dynamically  
 565 decided based on the input sequence, the  
 566 execution time increases as more segments (i.e.,  
 567 higher  $k$ ) are used. When we use, sequence  
 568 length of 1024, compression  $r = 256, k = 7, u =$   
 569 1, short attention segment size of 128, then the  
 570 size of aggregated attention (short, long, cache,  
 571 overlapping) is 1024x624.

572 Since our cache attention mechanism as  
 573 explained in section 3.1 is completely dynamic,  
 574 and uses the most attentive segments in  
 575 uncompressed form, we average the attention  
 576 vectors over  $p$  rows (to improve efficiency of  
 577 execution) as given by Equation 9.

578 If we use a sequence length of 1024, and average  
 579 over 256 rows, then the segments determined by  
 580 our cache attention mechanism part way through  
 581 the training of the model appears as shown in  
 582 Table 4. Note that to implement the  
 583 autoregressive behavior, the input sequence  
 584 cannot attend to a future segment. Our  
 585 implementation guarantees that the input  
 586 sequence can only attend to a previous segment.  
 587 For example, when attending to words 768-1023  
 588 in the input sequence, the maximum segment that  
 589 the cache attention can use is 47 (if the long  
 590 segment size is 16, then there are 64 segments in  
 591 the 1024 size sequence).

592 One of the important recent papers in handling  
 593 long contexts has indicated that current language  
 594 models do not robustly make use of information  
 595 in long input contexts (N. F. Liu et al., 2023).  
 596 They studied different models and concluded that  
 597 “performance is often highest when relevant  
 598 information occurs at the beginning or at the end  
 599 of the input context, and significantly degrades  
 600 when models must access relevant information in  
 601 the middle of long contexts.”

Input Sequence	Top k Attentive Segments ( $k=7, u=1$ )	Comments
0 – 255 words	[-1, -1, -1, -1, -1, -1, -1]	No cache segments are used to prevent future token leakage
256-511 words	[ 7, 8, 11, 12, 13, 14, <b>15</b> ]	Maximum segment allowed = <b>15</b>
512-767 words	[ 7, 8, 27, 28, 29, 30, <b>31</b> ]	Maximum segment allowed = <b>31</b>
768-1023 words	[ 8, 29, 32, 35, 37, 44, <b>47</b> ]	Maximum segment allowed = <b>47</b>

603 **Table 4.** Most Attentive Segments Used by our Cache  
604 Attention Part way in Training.

605  
606 Note that our cache attention model addresses  
607 this aspect nicely in the sense it uses attentive  
608 segments dynamically regardless they are needed  
609 in the beginning or the middle of input context.  
610 For example, see the last row in Table 4 which  
611 indicates the highest attentive segments that are  
612 used. Segments 32, 35, 37 are relatively in the  
613 middle of the input context. When we determine  
614 the most attentive segment to use in our cache  
615 attention, if the neighboring segment parameter  
616 count  $u > 1$ , then as we look at the segment index  
617 of the next or previous index, a duplicate may  
618 occur as the next segment may already be one of  
619 the high attentive segments. Similarly, if the high  
620 attentive segments belong to a future segment,  
621 we replace them by one of the allowed segments.  
622 Since information segmentation should not  
623 occur, the segment we select to be added is the  
624 one that is contiguous to an existing high  
625 attention segment.

## 626 6 Conclusions

627 Handling long contexts in an efficient manner  
628 without loss of performance is an important area  
629 of research in language models. Although many  
630 approaches have been recently proposed to  
631 address this problem, we present a new  
632 innovative solution that is motivated by the  
633 cache and virtual memory concepts in computer  
634 architecture. In such designs, if there is a cache  
635 or page miss, the needed data is retrieved from  
636 the disk or RAM. We handle long contexts by

637 diving them into small segments. By the  
638 magnitude of the compressed attention vectors,  
639 we determine the most attentive segments, and  
640 then use these in uncompressed form. Similar to  
641 the cache memory design, we also use  
642 consecutive segments near to the high attention  
643 segments to improve the language model  
644 predictive performance. Our results on the  
645 perplexity indicate significant improvement over  
646 the baseline architecture that uses short and long  
647 compressed attention. For the BPC, the cache  
648 attention mechanism does not show remarkable  
649 improvement on the baseline. We conjecture that  
650 the BPC that favors compression capability is not  
651 benefited by the relevant segment usage that our  
652 model provides which is helpful in model  
653 prediction capability. Another advantage of our  
654 approach is that the use of high attention  
655 segments is dynamic and depends on the input  
656 sequence. Thus, if the model needs to use  
657 information in the middle or anywhere in the  
658 input context, it is provided in uncompressed  
659 form via the high attention determination on the  
660 compressed segments.

## 661 7 Limitations

662 The only shortcoming of our approach we feel is  
663 that the dynamic segment attention is relatively  
664 slow during training. We partially overcome this  
665 by initially pretraining the model without  
666 dynamic attention, and then fine tune it on our  
667 cached attention. Our future work involves in  
668 applying the cache attention to reduce the model  
669 complexity of large language models and to  
670 create a hierarchical cache design such that very  
671 long contexts can be efficiently handled.

672 Further, our model sizes and datasets were  
673 constrained by computational resources available  
674 to us. We used GPU RTX 4090 and therefore  
675 could not use larger datasets such as PG-19 and  
676 run larger models with larger embedding size,  
677 layers, and heads.



## 8 References

- 684 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J.,  
685 Jones, L., Gomez, A. N., Kaiser, L., and  
686 Polosukhin, I. "Attention is all you need".  
687 Proceedings of Neural Information Processing  
688 Systems (NeurIPS), 2017
- 689 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
690 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
691 Diogo Almeida, Janko Altenschmidt, Sam  
692 Altman, Shyamal Anadkat, et al., "Gpt-4  
693 technical report," arXiv:2303.08774, 2023
- 694 G. Team, R. Anil, S. Borgeaud, Y. Wu, J.-B. Alayrac,  
695 J. Yu, R. Soricut, J. Sc (H. Touvron et al.,  
696 2023)halkwyk, A. M. Dai, A. Hauth et al.,  
697 "Gemini: a family of highly capable multimodal  
698 models," arXiv:2312.11805, 2023
- 699 Hugo Touvron, Thibaut Lavril, Gautier Izacard,  
700 Xavier Martinet, Marie-Anne Lachaux, Timothée  
701 Lacroix, Baptiste Rozière, Naman Goyal, Eric  
702 Hambro, Faisal Azhar, et al. "Llama: Open and  
703 efficient foundation language models,"  
704 arXiv:2302.13971, 2023.
- 705 Hugo Touvron, Louis Martin, Kevin Stone, Peter  
706 Albert, Amjad Almahairi, Yasmine Babaei,  
707 Nikolay Bashlykov, Soumya Batra, Prajjwal  
708 Bhargava, Shruti Bhosale, et al., "Llama 2: Open  
709 foundation and fine-tuned chat models,"  
710 arXiv:2307.09288, 2023.
- 711 Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q. V.,  
712 and Salakhutdinov, R. "Transformer-XL:  
713 Attentive language models beyond a fixed-length  
714 context," Proceedings of the Annual Meetings of  
715 the Association for Computational Linguistics  
716 (ACL), 2019.
- 717 Sinong Wang, Belinda Z. Li, Madian Khabsa, Han  
718 Fang, Hao Ma, "Linformer: Self-Attention with  
719 Linear Complexity," arXiv:2006.04768, 2020.
- 720 Iz Beltagy, Matthew E Peters, and Arman Cohan,  
721 "Longformer: The long-document transformer,"  
722 arXiv:2004.05150, 2020.
- 723 Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya,  
724 "Reformer: The efficient transformer," ICLR,  
725 2020.
- 726 Krzysztof Choromanski, Valerii Likhoshesterov, David  
727 Dohan, Xingyou Song, Andreea Gane, Tamas  
728 Sarlos, Peter Hawkins, Jared Davis, Afroz  
729 Mohiuddin, Lukasz Kaiser, et al., "Rethinking  
730 attention with performers," ICLR, 2021
- 731 Curtis Hawthorne, Andrew Jaegle, Cătălina  
732 Cangea, Sebastian Borgeaud et al., "General-  
733 purpose, long-context autoregressive modeling  
734 with Perceiver AR," ICML 2022
- 735 Albert Gu, Karan Goel, and Christopher Ré,  
736 "Efficiently Modeling Long Sequences with  
737 Structured State Spaces," arXiv:2111.00396v3,  
738 2022
- 739 Xuezhe Ma, Chunting Zhou, Xiang Kong, et.al.,  
740 "Mega: Moving Average Equipped Gated  
741 Attention," arXiv:2209.10655v3, 2023.
- 742 Daniel Y. Fu, Tri Dao, Khaled K. Saab, Armin W.  
743 Thomas, Atri Rudra, Christopher Ré, "Hungry  
744 Hungry Hippos: Towards Language Modeling  
745 with State Space Models," arXiv:2212.14052v3,  
746 2023
- 747 Chen Zhu, Wei Ping, Chaowei Xiao, Mohammad  
748 Shoeyb, Tom Goldstein, Anima Anandkumar,  
749 and Bryan Catanzaro, "Long-Short Transformer:  
750 Efficient Transformers for Language and  
751 Vision," 35th Conference on Neural Information  
752 Processing Systems (NeurIPS 2021).
- 753 Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin  
754 Paranjape, Michele Bevilacqua, Fabio Petroni,  
755 Percy Liang, "Lost in the Middle: How  
756 Language Models Use Long Contexts,"  
757 arXiv:2307.03172v3 [cs.CL] 20 Nov 2023.
- 758 Singh, Sushant and Ausif Mahmood. "The NLP  
759 Cookbook: Modern Recipes for Transformer  
760 Based Deep Learning Architectures." *IEEE  
761 Access* (2021)
- 762 Ji, Haozhe, Rongsheng Zhang, Zhenyu Yang,  
763 Zhipeng Hu and Minlie Huang. "LaMemo:  
764 Language Modeling with Look-Ahead  
765 Memory." *North American Chapter of the  
766 Association for Computational  
767 Linguistics* (2022).
- 768 Martins, Pedro Henrique, Zita Marinho and André F.  
769 T. Martins. "∞-former: Infinite Memory  
770 Transformer-former: Infinite Memory  
771 Transformer." *Proceedings of the 60th Annual  
772 Meeting of the Association for Computational  
773 Linguistics (Volume 1: Long Papers)* (2022)
- 774 Gu, Albert and Tri Dao. "Mamba: Linear-Time  
775 Sequence Modeling with Selective State  
776 Spaces." *ArXiv abs/2312.00752* (2023)
- 777 Beck, Maximilian, Korbinian Poppel, Markus  
778 Spanring, Andreas Auer, Oleksandra  
779 Prudnikova, Michael K Kopp, Günter  
780 Klambauer, Johannes Brandstetter and Sepp

781 Hochreiter. “xLSTM: Extended Long Short-  
782 Term Memory.” (2024).

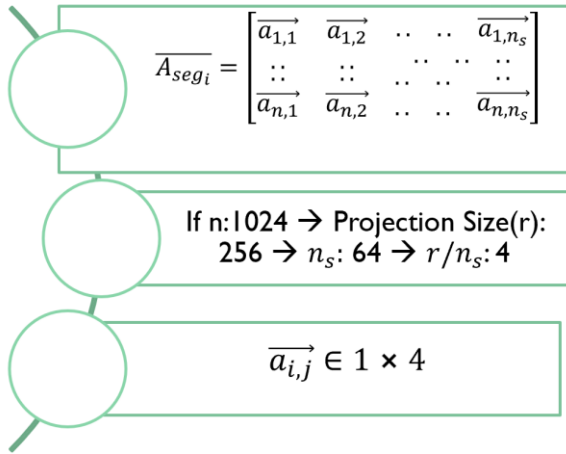
783

## 784 Appendix

### 785 **A. Further Details on our Enhanced** 786 **Caching Transformer**

787 In our caching protocol we compress and  
788 dynamically retrieve the most relevant  
789 compressed segments for any given input. Based  
790 on the design constraints an appropriate amount  
791 of input sequence compression is performed.  
792 Thereafter the sequence is split into the desired  
793 segments and we choose the most similar  
794 segments for each query and retrieve them in the  
795 original uncompressed form. It ensures only the  
796 most relevant information is being picked. This  
797 not only helps in reducing the context size but it  
798 also enables in preserving key information. This  
799 enhanced caching attention technique is  
800 explained in greater detail in the subsequent  
801 sections.

#### 802 **A.1 Enhanced Caching Attention**



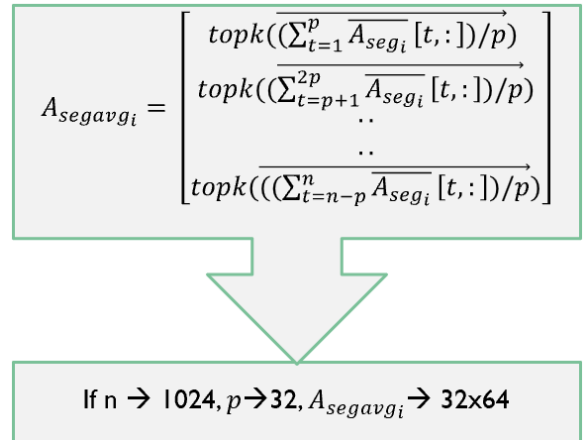
804 **Figure 5.** Downsized Compression of Attention  
805 Matrix along  $K_c, V_c$

806 Consider the length of the input sequence to be  
807 1024 tokens that need to be compressed and  
808 down projected to 256 tokens. Here we choose to  
809 divide the row into  $(n_s)$  64 segments. This will  
810 yield to a compression ratio  $(r/n_s)$  of 4. The  
811 attention matrix will be of size  $\overline{A}_{seg_i} \in \mathbb{R}^{n \times r}$ .  
812 Therefore for  $n_s$  segments, each row in  $\overline{A}_{seg_i}$   
813 will consist of row vectors with size  $r/n_s$ .

814 Further, the magnitude of the vector  $\overline{a}_{i,j} \in$   
815  $\mathbb{R}^{1 \times r/n_s}$  will represent the attention of the  $i^{th}$   
816 word token to the  $j^{th}$  compressed segment in the  
817 long attention as shown in Figure 5. Thereafter,  
818 we compute the root mean square for each of the  
819  $(1 \times 4)$  sized attention vectors  $\overline{a}_{i,j}$ , hence the  
820 dimension across each row is downsized from  
821 256 to 64. We use this size for the subsequent  
822 attention processing steps as demonstrated in the  
823 following section.

#### 824 **A.2 Averaging in Segment Caching**

825 Attention computation and top-k segment  
826 retrieval across all 1024 rows turned out to be  
827 computationally cumbersome and time intensive.  
828 Therefore, to achieve execution efficiency, we  
829 averaged all 1024 input vectors across  $p$   
830 consecutive rows for the previous attention  
831 matrix  $\overline{A}_{seg_i} \in \mathbb{R}^{n \times r}$  where  $p$  is a  
832 hyperparameter.



833

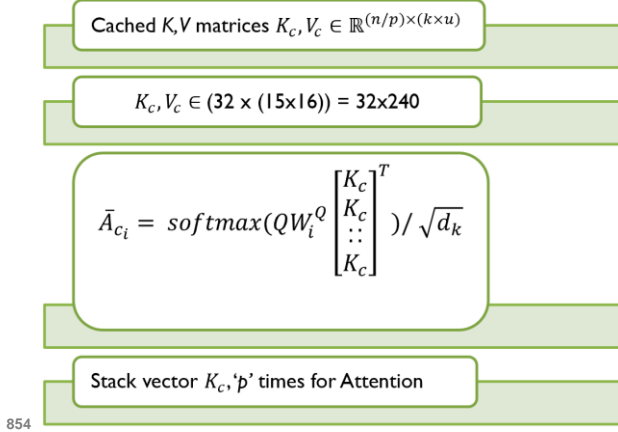
834 **Figure 6.** Averaged Compression of  
835 Attention Matrix along the Input Length

836 This segment attention matrix is further reshaped  
837 and compressed into  $A_{seg_avg_i} \in \mathbb{R}^{m \times n_s}$ , where  
838  $m = n/p = 32$  as shown in Figure 6. This  
839 implementation was key for our model to  
840 achieve superior results outperforming other  
841 popular language models of similar size as  
842 mentioned in Table 2 and resulted in a faster run  
843 time as well.

#### 844 **A.3 Top-k Retrieval in Segment Caching**

845 Post the compression and averaging, the top  $k$   
846 most similar segments were chosen to be  
847 retrieved by the order of the attention magnitude

848 between the modified input and key/value  
 849 matrices. These segments were picked  
 850 corresponding to each row  $m$ , which is an  
 851 averaged input sequence of 32 consecutive words  
 852 (averaged down from 1024) from the segment  
 853 attention matrix  $A_{segavg_i}$ .

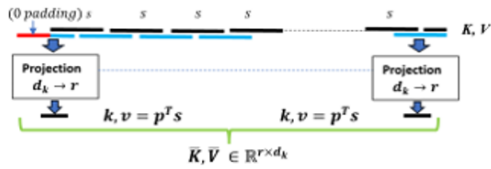


855 **Figure 7.** Enhanced Attention Matrix after  
 856 top-k retrieval

857 The hyperparameter  $k$  is chosen based on the  
 858 performance needs and based on that value along  
 859 with the  $k^{th}$  segment, we also extract one  
 860 segment before and after the  $k^{th}$  attentive  
 861 segment.

862 Therefore, we define  $u$  as the hyperparameter  
 863 that regulates the number of adjacent segments  
 864 around  $k$  that need to be retrieved from the  
 865 sequence. For instance, with  $k = 5$  and  $u = 3$   
 866 will result in a total of 15 uncompressed  
 867 extracted segments of length 16 from each row  
 868 as shown in Figure 7.

#### 869 A.4 Overlapping Segments in Long Attention



$$P_{o_i} = \text{Softmax}(KW_i^{P_o}), \bar{K}_{o_i} = P_{o_i}^T KW_i^K, \bar{V}_{o_i} = P_{o_i}^T VW_i^V$$

$$\bar{A}_{o_i} = \text{Softmax} \left[ \frac{QW_i^Q \bar{K}_{o_i}^T}{\sqrt{d_k}} \right]$$

$$\bar{H}_{o_i} = \bar{A}_{o_i} (P_{o_i}^T VW_i^V)$$

870

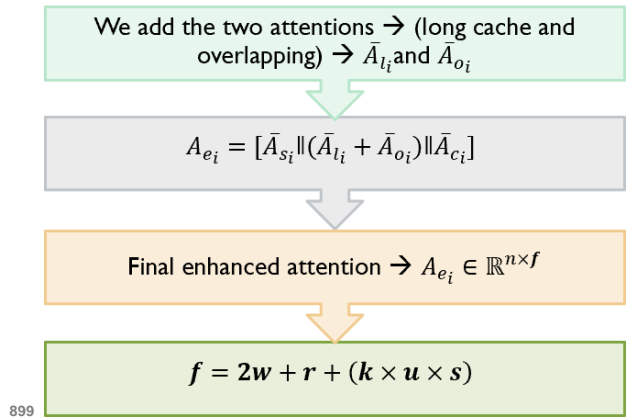
871 **Figure 8.** Long Attention with Overlapping Segments

872 As discussed earlier that the segmentation of  
 873 input into chunks leads to fragmentation of long-  
 874 term information. This becomes a challenge in  
 875 building long term dependency. This issue hasn't  
 876 been addressed in prior Transformer based  
 877 language models. Therefore, we augment the  
 878 long attention with segments with a 50% overlap  
 879 to maintain the continuity of data as shown in  
 880 Figure 8. The model is trained with the  
 881 overlapping data as the query that needs to learn  
 882 the original chunks as key and values.

#### 883 A.5 Aggregated Enhanced Long Short 884 Attention

885 Thereafter we add the overlapping attention  $\bar{A}_{o_i}$   
 886 to the long cache attention  $\bar{A}_{l_i}$  who have similar  
 887 shapes. The sliding window (short) attention  $\bar{A}_{s_i}$   
 888 and our caching attention  $\bar{A}_{c_i}$  are concatenated to  
 889 the above summed attention as pictorially  
 890 demonstrated in Figure 9.

891 Here  $\parallel$  indicates the catenation of different  
 892 attentions,  $w$  is the window size in short i.e.,  
 893 sliding window attention,  $r$  is the projection  
 894 size in compressing the long attention,  $k$  is  
 895 the *top k* factor in retrieving high attention  
 896 top  $k$  segments,  $s$  is the segment size in long  
 897 attention,  $u$  determines the number of segments  
 898 to be retrieved adjacent to the top  $k^{th}$  one.



899 **Figure 9.** Complexity of the Enhanced Attention

900 Finally, Figure 10 shows the four attention  
 901 mechanisms that are simultaneously aggregated  
 902 and successfully inducted in our model  
 903 architecture.  
 904

905

## Enhanced architecture obtained by aggregating the four attentions

Short attention  $\bar{A}_{s_t} \in \mathbb{R}^{n \times 2w}$ , uses segment-wise sliding window

Segment based compressed long attention  $\bar{A}_{l_t} \in \mathbb{R}^{n \times r}$  as proposed in Transformer LS

Our cache attention  $\rightarrow \bar{A}_{c_t} \in \mathbb{R}^{n \times (k \times u \times s)} \rightarrow$   
Retrieval of **uncompressed** high attention segments

Overlapping segment-based compressed attention  $\rightarrow$   
 $\bar{A}_{o_t} \in \mathbb{R}^{n \times r}$

906

907 **Figure 10.** Aggregated Enhanced Attention

908

909

910

911

912

913

914

915

916

917

918

919