

DAPE V2: PROCESS ATTENTION SCORE AS FEATURE MAP FOR LENGTH EXTRAPOLATION

Anonymous authors

Paper under double-blind review

ABSTRACT

The attention mechanism is a fundamental component of the Transformer model, contributing to interactions among distinct tokens. In general, the attention scores are determined simply by the key-query products. However, this work’s occasional trial (combining DAPE and NoPE) of including additional MLPs on attention scores without position encoding indicates that the classical key-query multiplication may limit the performance of Transformers. In this work, we conceptualize attention as a feature map and apply the convolution operator (for neighboring attention scores across different heads) to mimic the processing methods in computer vision. Specifically, **the main contribution of this paper is identifying and interpreting the Transformer length extrapolation problem as a result of the limited expressiveness of the naive query and key dot product, and we successfully translate the length extrapolation issue into a well-understood feature map processing problem.** The novel insight, which can be adapted to various attention-related models, reveals that the current Transformer architecture has the potential for further evolution. Extensive experiments demonstrate that treating attention as a feature map and applying convolution as a processing method significantly enhances Transformer performance.

1 INTRODUCTION

Transformer-based models (Vaswani et al., 2017) have delivered exceptional performances across widespread applications, including language processing (Zhang et al., 2020; Guo et al., 2022; Ainslie et al., 2023), computer vision (Alexey, 2020; Touvron et al., 2021; Liu et al., 2021a; Chen et al., 2024; Peebles & Xie, 2023), quantitative research (Zhou et al., 2024b; Liu et al., 2021b; Wu et al., 2023), and scientific machine learning (Taylor et al., 2022; Geneva & Zabarar, 2022). However, the quadratic cost of the key-query multiplication for processing a sequence raised much concern about the modern architecture of Transformers especially for long context inputs. To address the issue of storage and computation efficiency, recent research delves into developing more efficient architectures, such as sparse structural attention (Xiao et al., 2024d; Zhu et al., 2024), adaptive key selection (Xiao et al., 2024a; Fountas et al., 2024), and hybrid models (Lieber et al., 2024). While these adaptations enhance efficiency, they often involve tradeoffs with model effectiveness.

At the same time, there is another voice advocating for refining the model design for tackling complex tasks, rather than prioritizing efficiency. Positional encoding is one of the key components of the attention mechanism. Although the widely recognized decoder-based Transformer can implicitly incorporate the positional information of tokens, growing evidence both theoretically and empirically shows that the well-designed explicit positional encoding significantly enhances the model performances, especially in long-context tasks (Su et al., 2024b; Press et al., 2021; Zhao et al., 2023). In practice, Transformers depend on positional encoding to explicitly incorporate positional information, enabling the model to make meaningful token predictions. **Without these encodings, token generation would lack the necessary contextual order.** The well-recognized RoPE (Su et al., 2024b), which is adopted in LLaMA (Touvron et al., 2023), distinguishes the token order by rotating with different angles depending on the token position. However, it demonstrated a notable performance degradation, failing entirely when the input length is double that of the training length (Peng et al., 2023b; Chen et al., 2023a; Ding et al., 2024b). The undesirable performance degradation is also observed for other positional encoding methods, e.g., ALiBi (Press et al., 2021) and Kerple (Chi et al., 2022). FIRE (Li et al., 2023c) alleviates the long-context extrapolation by learnable posi-

054 tional encodings, trying to capture the suitable positional representation by MLPs. Recently, the
 055 data-adaptive positional encoding method, namely DAPE (Zheng et al., 2024), which adjusts dy-
 056 namically with context, enhances the length generalization by incorporating the attention scores and
 057 positional information with a more complex mechanism.

058 In this paper, we propose that precise attention scores are crucial for improving Transformer length
 059 extrapolation, and we introduce a new perspective on the attention mechanisms. Traditionally, at-
 060 tention scores are computed through the dot product of the query and key vectors. As illustrated
 061 in Figure 1, further processing these attention scores using a neural network—a general case of
 062 DAPE (Zheng et al., 2024)—can significantly enhance the length generalization of Transformers,
 063 even in the absence of positional encoding (NoPE). Therefore, we suggest treating attention scores as
 064 feature maps. By conceptualizing attention as an image feature map (with dimensions $[B, C, W, H]$
 065 for batch size, channel size, width, and height), we can achieve more accurate attention scores by ap-
 066 plying techniques used in image processing. In this work, we employ different kernel sizes (such as
 067 1×3) to process attention, finding that the perplexity (ppl) of attention decreases significantly—from
 068 over 600 to just above 100—when trained on a sequence length of 128 and evaluated on a length of
 069 8192.

070 In summary, our contributions are as follows:

- 071 1. We highlight that the coarse attention mechanism, which is the direct result of the query
 072 and key dot product, limits the Transformer’s ability to extrapolate to longer sequences.
 073 However, Transformers can achieve good length extrapolation performance with careful
 074 processing of attention scores.
 075
- 076 2. Besides developing better position encoding (Vaswani et al., 2017) or position interpola-
 077 tion (Chen et al., 2023b) for length extrapolation, we propose the third direction: by treating
 078 attention scores as feature maps and refining them using image processing techniques like
 079 convolution, we can enhance the Transformer’s extrapolation capabilities.
- 080 3. We conducted extensive experiments on language tasks to support our claims and believe
 081 that these insights can significantly improve the Transformer’s performance in length ex-
 082 trapolation.
 083

084 2 RELATED WORKS

085 **Absolute Positional Encoding** Absolute positional encoding (APE), introduced by Vaswani et al.
 086 (2017), enables Transformers to incorporate positional information. Specifically, at the first layer,
 087 each position i is assigned a real-valued encoding $e_i \in \mathbb{R}^d$, which can be either learnable or a fixed
 088 sinusoidal encoding (Vaswani et al., 2017; Kiyono et al., 2021; Likhomanenko et al., 2021; Wang
 089 et al., 2020; Liu et al., 2020), and this encoding is then added to the input sequence. Although
 090 this approach is straightforward, Transformers relying on APE tend to struggle with generalizing to
 091 longer sequences (Press et al., 2021).
 092

093 **Relative Positional Encoding** Relative positional encoding (RPE) offers an alternative for em-
 094 bedding positional information (Shaw et al., 2018; Raffel et al., 2020; Press et al., 2021). A widely
 095 used RPE method in large language models is rotary positional encoding (RoPE)(Su et al., 2024b;
 096 Chowdhery et al., 2023; Touvron et al., 2023). To address length extrapolation challenges(Press
 097 et al., 2021; Kazemnejad et al., 2024), positional interpolation (PI) has been introduced (Chen et al.,
 098 2023b) to extend the context window. Building on this approach, models like LongLora (Chen et al.,
 099 2023c), LongRope (Ding et al., 2024b), YaRN (Peng et al., 2023b), and CLEX (Chen et al., 2023a)
 100 have emerged. Another notable direction involves additive positional encoding. For most additive
 101 RPE techniques, the computation of pre-softmax attention logits can be expressed using the formula:
 102 $A_{\text{RPE}}(\mathbf{X}) = \mathbf{X} \mathbf{W}_Q (\mathbf{X} \mathbf{W}_K)^\top + \mathbf{B}$, where the bias matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ is derived from the
 103 positional encoding function $b : \mathbb{N}^2 \rightarrow \mathbb{R}$, with the (i, j) -th entry of \mathbf{B} defined as $b(i, j)$. Different
 104 parameterizations of b give rise to various RPE variants. Methods supporting arbitrary sequence
 105 lengths include T5’s RPE (Raffel et al., 2020), ALiBi (Press et al., 2021), Kerple (Chi et al., 2022),
 106 Sandwich (Chi et al., 2023a), and FIRE (Li et al., 2023c). Recently, DAPE (Zheng et al., 2024) has
 107 been introduced, employing MLPs to dynamically adjust bias values based on the input data.

Data-Adaptive Related Positional Encoding. Transformer-XL (Dai et al., 2019) introduced the use of learnable query and key biases for adaptive positional encodings. Data-Adaptive Positional Encoding (DAPE)(Zheng et al., 2024) extends this idea by leveraging MLPs to adjust positional encodings based on attention over the head dimension for length extrapolation, ensuring different input data receive unique positional encodings. Contextual Positional Encoding(Golovneva et al., 2024) further refines this by conditioning position increments on specific tokens, as determined by the model, allowing positions to adapt based on context.”

3 METHOD

In this section, we first review the previously developed Data-Adaptive Positional Encoding method (DAPE), which incorporates attention scores and positional information through MLPs. As a proof-of-concept, our occasional trial on DAPE without the positional information (as shown in Figure 1) suggests that regarding attention as a feature map and processing it with classical operators (e.g., convolution) can enhance the Transformers’ behavior. As discussed in some previous works the perplexity scores come mostly from the associative recall (i.e., copy) tasks. In addition, we theoretically show by construction that the proposed method can explicitly realize the associative recall task, in contrast to the implicit conduct through positional encoding in standard Transformers. **The two key differences between DAPE (Zheng et al., 2024) and this work are: 1) Insight:** DAPE attributes length extrapolation performance gains to adaptive position encoding, while this work finds DAPE could still improve performance without position encoding so that we take a broader view, explaining that the Transformer’s length extrapolation ability is limited by the expressiveness of the naive query-key dot product, which can be enhanced using image processing techniques; **2) Performance:** As shown in Figure 1, DAPE is designed for additive RPE and may underperform with non-additive RPE (e.g., RoPE), whereas this work suggests that increasing kernel size (e.g., with DAPE_{1×3}) may improve RoPE’s performance. The DAPE_{1×3} implementation is shown in Appendix L.

3.1 ADDITIVE RELATIVE POSITIONAL ENCODING

For most additive relative positional encoding (ARPE) methods, the computation of pre-softmax attention logits can be unified under the following formula:

$$\mathbf{A}_{\text{ARPE}}(\mathbf{X}) = \mathbf{X}\mathbf{W}_Q(\mathbf{X}\mathbf{W}_K)^\top + \mathbf{B}, \quad (1)$$

where the bias matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ is induced by the position encoding function $b : \mathbb{N}^2 \rightarrow \mathbb{R}$ and the (i, j) -th entry of \mathbf{B} is defined as $b(i, j)$. Various formulations and parameterizations of b give rise to different variants of RPE. Examples of additive RPE include: (1) ALiBi: $b(i, j) = -r|i - j|$, with the scalar $r > 0$ as a hyper-parameter; (2) Kerple: $b(i, j) = -r_1 \log(1 + r_2|i - j|)$ with r_1 and r_2 are two learnable parameters; (3) FIRE: $b(i, j) = f_\theta \left(\frac{\psi(i-j)}{\psi(\max\{L, i\})} \right)$, where the positional encoding function f_θ parameterized by θ is learned from data and ψ is a transformation function aimed at assigning more model capacity to local positions.

Data-Adaptive Position Encoding (DAPE) The DAPE rewrite the Equation 1 as the following:

$$\mathbf{A}_{\text{DAPE}}(\mathbf{X}) = \mathbf{X}\mathbf{W}_Q(\mathbf{X}\mathbf{W}_K)^\top + f(\mathbf{X}\mathbf{W}_Q(\mathbf{X}\mathbf{W}_K)^\top, \mathbf{B}). \quad (2)$$

Here, $f : \mathbb{R}^{T \times T} \times \mathbb{R}^{T \times T} \rightarrow \mathbb{R}^{T \times T}$ is an element-wise function and T is the sequence length. Another variant of DAPE is with residual, which is the following:

$$\mathbf{A}_{\text{DAPE}}(\mathbf{X}) = \mathbf{X}\mathbf{W}_Q(\mathbf{X}\mathbf{W}_K)^\top + \mathbf{B} + f(\mathbf{X}\mathbf{W}_Q(\mathbf{X}\mathbf{W}_K)^\top, \mathbf{B}). \quad (3)$$

In practice, DAPE (Zheng et al., 2024) utilizes a two-layer *LeakyReLU* MLP with hidden dimension D_{DAPE} (default value is 32) to parameterize $f(\cdot)$ due to its universal approximability (Leshno et al., 1993). All parameters are learned directly from the data during the training process. This architecture allows $f(\cdot)$ to dynamically adjust positional embeddings based on the input sequence data, ensuring that the encoding method is both adaptive and dependent on the input data.

3.2 SPECIAL CASE OF DAPE: BIAS IS ZERO

DAPE was originally designed to dynamically adjust the positional encoding by incorporating input data information. Generally, any additive positional encoding method that includes positional information can be represented as the matrix \mathbf{B} in the DAPE model, as outlined in Equation 2. Notably,

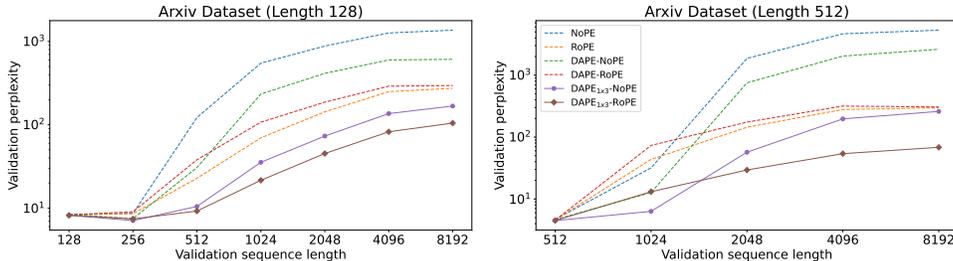


Figure 1: The result of DAPE (Zheng et al., 2024) (equivalent to kernel 1×1 in our explanation) and DAPE_{1x3} (kernel 1×3 by this work), with baseline NoPE and RoPE. The model is trained with length 128 and length 512 respectively. The DAPE_{1x3} denotes that we use $H \times 1 \times 3$ convolutions kernel size on the attention score with shape $[B, H, T, T]$. We find that DAPE can even improve the performance of NoPE (without biased position encoding), suggesting that the explanation in Zheng et al. (2024), which attributes the improvement to adaptive position encoding, may have a more general underlying cause.

No Positional Encoding (NoPE) (Kazemnejad et al., 2024) is a special case of additive RPE that assigns zero value to the matrix B . The mathematical formulation of DAPE equipped with NoPE is given by:

$$A_{\text{DAPE}}(\mathbf{X}) = \mathbf{X} \mathbf{W}_Q (\mathbf{X} \mathbf{W}_K)^\top + f(\mathbf{X} \mathbf{W}_Q (\mathbf{X} \mathbf{W}_K)^\top). \quad (4)$$

The DAPE Zheng et al. (2024) is designed for additive RPE but not trying NoPE or RoPE, and we present the results of DAPE-NoPE and DAPE-RoPE in the following.

The result of DAPE-NoPE Compared with the standard Transformer architecture, DAPE-NoPE introduces additional MLPs post the key-query multiplication and prior to the softmax operator. As shown in Figure 1, experimental evidence suggests that DAPE with NoPE significantly outperforms the basic NoPE, prompting a reconsideration of the behaviors of standard Transformers. The additional MLPs (i.e., denoted as $f(\cdot)$ in Equation 4) facilitate information sharing across attention heads and complicate the attention calculation with nonlinear transformation beyond the simple key-query multiplication. This leads to a critical question: *Is the current Transformer architecture, particularly the attention mechanism, sufficiently expressive for real-world language tasks?* Although numerous studies aim to enhance efficiency by reducing computation and storage in standard Transformers, these often come at the cost of effectiveness, potentially hindering the evolution of next-generation Transformer models. Motivated by these insights and observations, we enhance the Transformer’s expressiveness and behavior by regarding attention as a feature map and applying convolutional operations, akin to those used in computer vision.

The result of DAPE-RoPE. Building on the hypothesis that DAPE enhances Transformer performance by processing pre-softmax scores with MLPs, we explore its applicability to non-additive positional encoding methods, specifically RoPE (Su et al., 2024b). In the DAPE-RoPE configuration, DAPE-RoPE first computes the classic attention scores of key-query multiplication with RoPE, which are then refined using the MLPs described in Equation 4. The visualized results of the validation perplexity for DAPE-RoPE and other positional encoding methods are presented in Figure 1. The results indicate that DAPE-RoPE may degrade the performance, while DAPE_{1x3}-RoPE (with kernel size 1×3 , proposed by this work) not only improves overall performance but also excels in length extrapolation tasks, particularly at larger sequence lengths. This finding substantiates the effectiveness of DAPE_{1x3}-RoPE, confirming its superior performance compared to standard RoPE, attributing to the additionally introduced convolution operations to the attention scores.

3.3 DAPE V2: PROCESS ATTENTION SCORES AS FEATURE MAPS

As discussed above, improving Transformer performance necessitates refining the processing of attention score computation beyond the conventional key-query multiplication. We propose regarding the pre-softmax attention scores as feature maps (4-dimensional tensors) and applying convolutional operators, which may could additionally involve position information with zero padding and higher expressiveness (Kayhan & Gemert, 2020) but MLP does not involve additional position information because there is no zero padding. This approach facilitates enhanced communication across neighboring tokens and heads, drawing parallels to popular techniques used in computer vision. This

novel method aims to leverage the spatial relationships within tokens, potentially unlocking new aspects of model capabilities.

Rethink the DAPE formulation. In DAPE (Zheng et al., 2024), MLPs are utilized to process and integrate attention and biases. Notably, these MLP operations can be equated to convolution operations with 1×1 kernel (Krizhevsky et al., 2012; Simonyan & Zisserman, 2014; He et al., 2016), a stride of one, and no padding. Consequently, we can reformulate the DAPE in Equation 3 as the following:

$$A_{\text{DAPE}}(\mathbf{X}) = \mathbf{X}\mathbf{W}_Q(\mathbf{X}\mathbf{W}_K)^\top + \mathbf{B} + \text{Conv}(\text{tril}((\mathbf{X}\mathbf{W}_Q(\mathbf{X}\mathbf{W}_K)^\top, \mathbf{B})). \quad (5)$$

where \mathbf{X} is the input embedding, $\mathbf{X}\mathbf{W}_Q$ gives the query embedding and the $\mathbf{X}\mathbf{W}_K$ gives the key embedding. Under such formulation, DAPE employs convolution operation to process the pre-softmax attention scores of key-query multiplication. The `tril(\cdot)` returns the lower triangular part of the matrix and the other elements of the result tensor out are set to 0. The resulting attention tensor has a shape of $[B, H, T, T]$, where the four dimensions correspond to the batch size, number of heads, and the context length for both the query and key. This mirrors the structure of an image feature tensor with shape $[B, C, H, W]$, where the dimensions represent the batch size, number of channels, image height, and image width, respectively. This structural similarity underscores the feasibility of considering attention scores as a tensor of feature mappings, where popular and effective convolution operations can be leveraged for refined processing.

Process attention with more powerful convolution operation. In computer vision, the limitations of 1×1 kernels for processing image features are well-recognized. To improve upon the attention scores processed by these kernels (e.g., DAPE), we introduce $1 \times k$ kernels with a stride of 1 and padding of $k - 1$. This approach allows for wider and deeper convolution across key dimensions and heads without information leakage, as we ensure the attention scores remain lower-triangular. This mechanism is visualized in Appendix K. The use of $1 \times k$ kernels suggests a targeted convolution along the key dimensions across heads. In general, while extending this to include the query dimensions as a standard kernel is theoretically possible, it would significantly increase computational demands. Our forthcoming analysis demonstrates that Transformers modified with $1 \times k$ convolution are adept at associative recall tasks (i.e., the copy task), validating the benefits of integrating convolution in attention calculation. We left as a future work investigating the performances and the soundness of general convolution kernels, such as square sizes. **The key contribution of this work is providing a novel insight that suggests applying convolution operations and processing attention as feature maps to improve Transformers’ performances.**

Realizing associate recall tasks through convolution. As pointed out in some previous works (Arora et al., 2024), the perplexity scores of Transformers mostly result from the performances on associate recall tasks (i.e., the copy tasks). Numerous studies have explored the mechanism of associative recall within Transformers, both from theoretical perspectives and experimental validations (Arora et al., 2024; Bietti et al., 2024; Golovneva et al., 2024). Here, we theoretically prove that the proposed model can realize the associative recall tasks. Notably, this capability is achieved independently of positional encodings, marking a significant advancement in the flexibility and applicability of the proposed architecture. By integrating convolutional operations, we enable the model to handle associative tasks more effectively, leveraging spatial relationships inherent in the data, similar to methods used in image processing. To explain the associative recall mechanism, (Bietti et al., 2024) proved that the first layer of the Transformer is responsible for the previous token mechanism through the positional encoding. More specifically, given a sequence of input tokens $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$ with corresponding orthogonal positional encoding vectors $[\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N]$, the first layer primarily facilitates the copying of the previous token to the current token (e.g., $\mathbf{x}_i + \mathbf{W}_V^1 \mathbf{x}_{i-1}$, where \mathbf{W}_V^1 is the value matrix at the first layer of the Transformer). The input tokens are combined with positional encodings $\mathbf{x}_i + \mathbf{p}_i$ and the key-query weight matrix is defined as $\mathbf{W}_K^{1\top} \mathbf{W}_Q^1 = \sum_{i=1}^N \mathbf{p}_{i-1} \mathbf{p}_i^\top$. The orthogonality of positional encoding vectors and the special choices of the key-query matrix ensure that attention scores predominantly focus on the previous token. In contrast to this implicit mechanism in standard Transformers, our proposed method leverages a convolution operation to explicitly realize associative recall. This approach not only simplifies the process but also enhances its effectiveness by directly manipulating the spatial relationships within tokens and attention scores. Consider a scenario where the word ‘‘Hakuna’’ is

consistently followed by “Matata” within a lengthy paragraph. Without the loss of generality, we assume that x_1 and x_2 represent the tokens of “Hakuna” and “Matata” respectively, and $x_N = x_1$ implies that the N-th token in the sequence is “Hakuna”. Then we expect that the Transformer can predict and output the next token x_{N+1} as “Matata”. For simplicity, we consider a one-head Transformer without positional encoding. We employ a convolution operation with a kernel size of 1×2 and weights $[-1, 1]$. Note that the convolution is linear and processing the attention scores along the key dimensions is effectively equivalent to applying convolutions directly to the key vectors themselves. Consequently, the key vector of x_2 can be expressed as $W_K^1(x_2 - x_1)$ and the query vector for x_N admits $W_Q^1 x_N$. By configuring the matrix $W_K^{1\top} W_Q^1$ to be $-I$, the attention mechanism after the convolution predominantly allocates the attention values of x_N to the token x_2 . This ensures that the token values of x_2 are effectively copied to x_N , resulting in the model outputting “Matata” following “Hakuna”.

Proposition 1. *Transformers incorporating convolution operations can perform associative recall tasks without the need for positional encoding.*

Comparisons with hybrid models of convolution and Transformers. Recent developments in hybrid architectures have seen the integration of convolutional and Transformer models to capitalize on the strengths of both. For instance, Fu et al. (2022) introduced the FlashConv layer, which combines the efficiency of State Space Models (SSMs) with the capabilities of attention-based models. Similarly, Arora et al. (2024) developed a gated convolution layer, noted for its effectiveness in addressing associative recall tasks. These models typically stack convolution layers directly with standard Transformer layers, resulting in modifications to the token values through convolution. In contrast, our model adopts a distinctive approach by applying convolution along the key dimension during the computation of attention scores. This method preserves the original token values while still leveraging the convolution’s benefits for processing attention.

4 EXPERIMENT

Baselines. We evaluate the proposed $\text{DAPE}_{1 \times 3}$ against several well-established baselines, including NoPE (Kazemnejad et al., 2024), RoPE (Su et al., 2024b), T5’s Bias (Raffel et al., 2020), ALiBi (Press et al., 2021), Kerple (Chi et al., 2022), FIRE (Li et al., 2023c), CoPE (Golovneva et al., 2024), and DAPE (Zheng et al., 2024). As our kernels are applied across all heads, we simplify by omitting the kernel size description at the head dimension. For example, $\text{DAPE}_{1 \times 3}$ indicates the use of a $H \times 1 \times 3$ convolution kernel size on the attention scores, with a shape of $[B, H, T, T]$.

Datasets. Our analysis is based on training language models using the Arxiv and Books3 datasets, commonly employed benchmarks for assessing model performance (Press et al., 2021; Chi et al., 2022; Li et al., 2023c; Ding et al., 2024b). We begin our evaluation by processing entire sequences and comparing the zero-shot perplexity of the last 256 tokens across various input lengths. In addition to perplexity, we also leverage downstream datasets with randomized positional encoding (Rouss et al., 2023) to further assess $\text{DAPE}_{1 \times 3}$.

Experiment settings. Initially, we compare $\text{DAPE}_{1 \times 3}$ with other baselines at training lengths of 128, 512, and 1024, using 125M decoder-only Transformers (Brown et al., 2020), with model configurations detailed in Appendix I. Subsequently, we evaluate the performance of different training lengths using the same number of training tokens but with larger model sizes (350M and 2.7B). We also explore the impact of the convolutional hidden dimension D_{DAPE} , the effect of information leakage, and the influence of varying kernel sizes. Additionally, we examine the computational efficiency of $\text{DAPE}_{1 \times 3}$, focusing on processing times. Lastly, we evaluate $\text{DAPE}_{1 \times 3}$ on algorithmic reasoning datasets using accuracy metrics. Compared to DAPE (Zheng et al., 2024), $\text{DAPE}_{1 \times 3}$ demonstrates a more pronounced attention sink (Xiao et al., 2024d), as visualized in Appendix K.

4.1 COMPARE WITH BASELINES

$\text{DAPE}_{1 \times 3}$ -Kerple improves performance within training length, proving its ability to process the entire sequence. According to Figure 2, the proposed $\text{DAPE}_{1 \times 3}$ -Kerple demonstrates superior performance across various training and evaluation lengths. Specifically, $\text{DAPE}_{1 \times 3}$ -Kerple achieves

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

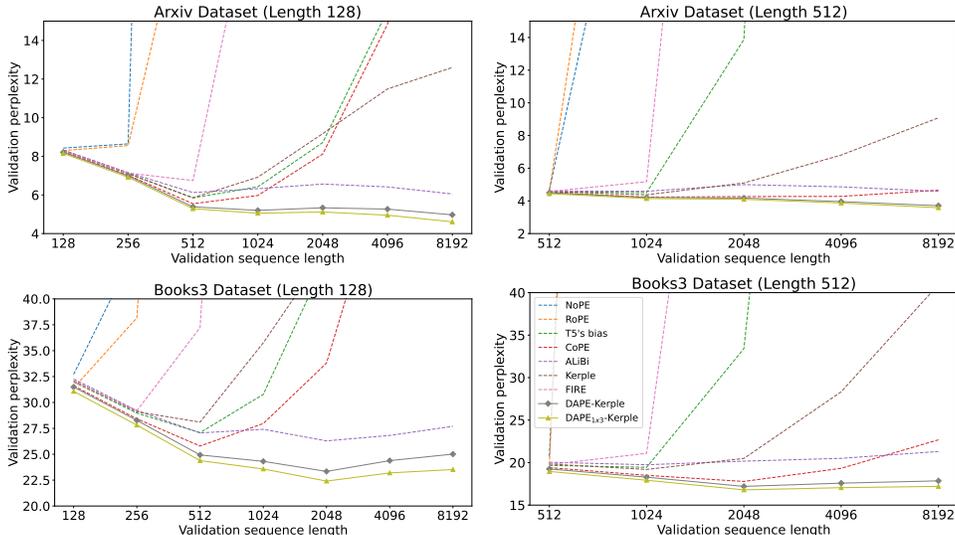


Figure 2: **Comparisons with baselines:** performance with training lengths 128 and 512 on Arxiv and Books3 datasets.

the best performance where the training length is 128 or 512 and the evaluation length ranges from 128 to 8192. This performance consistency is observed across both the arXiv and Books datasets. For instance, on the arXiv dataset with a training length of 512, DAPE_{1×3}-Kerple achieves a perplexity score of 4.44. This score surpasses those of other methods, such as DAPE-Kerple with a perplexity of 4.49, CoPE with 4.51, Kerple with 4.57, and RoPE with 4.57. These results indicate that DAPE_{1×3}-Kerple has a more robust modeling capability within the training length compared to the other methods evaluated. The Appendix B also presents the performance of different methods with training length 1024. The improvements are not only significant but also consistent, reinforcing the efficacy of the DAPE_{1×3}-Kerple approach in handling various training lengths effectively.

DAPE_{1×3}-Kerple improves performance beyond training length. The advantages of DAPE_{1×3}-Kerple extend beyond the training length. When the training length is set to 128 and the evaluation length is extended to 8192, DAPE_{1×3}-Kerple achieves a perplexity score of 4.60 on the arXiv dataset and 23.52 on the Books3 dataset. These scores are significantly better than those achieved by DAPE-Kerple, which records perplexity scores of 4.97 and 25.01 on the arXiv and Books3 datasets, respectively. Similarly, CoPE performs poorly with perplexity scores of 29.86 on the arXiv dataset and 90.66 on the Books3 dataset under the same conditions. Furthermore, when the training duration is increased to 512, DAPE_{1×3}-Kerple continues to deliver the best performance, further validating its superior generalization capabilities. These findings highlight the scalability and robustness of DAPE_{1×3}-Kerple, which is attributed to the introduced convolution operator, making it a promising approach for diverse data scenarios and lengths.

4.2 PERFORMANCE WITH SAME TRAINING TOKENS AND DIFFERENT TRAINING LENGTH

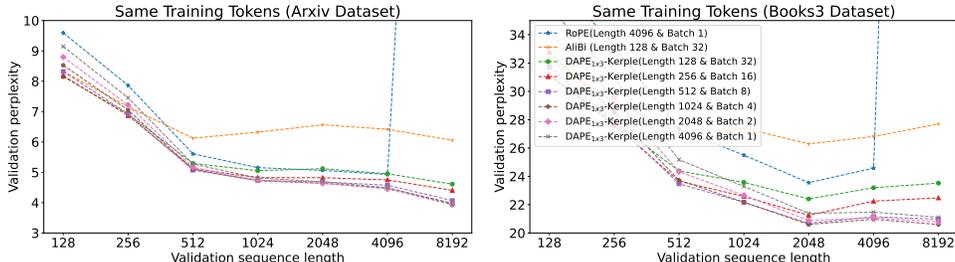


Figure 3: **The performance with same training tokens and different training length.** With the same training tokens, DAPE_{1×3} with training length 512 could even achieve better performance than RoPE with training length 4096.

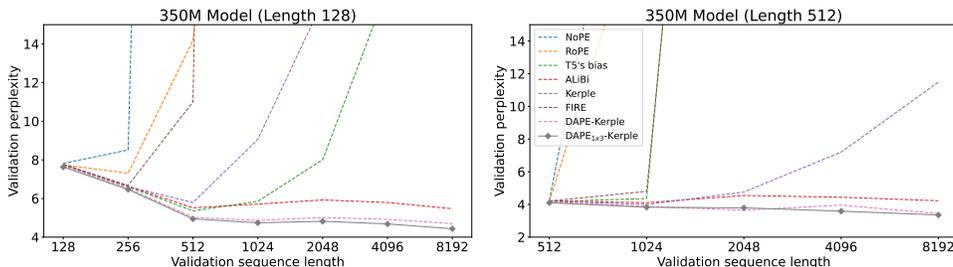
378 **Compared to RoPE, with the same training tokens, DAPE_{1×3}-Kerple with a training length**
 379 **of 128 achieves performance comparable to RoPE with a training length of 4096, for varying**
 380 **evaluation length.** As shown in Figure 3, for DAPE_{1×3}-Kerple trained with a length of 128, it
 381 achieves a perplexity (ppl) of 8.15 at an evaluation length of 128 and 4.95 at an evaluation length of
 382 4096 on the arXiv dataset. In comparison, RoPE trained with a length of 4096 achieves a ppl of 9.59
 383 at an evaluation length of 128 and 4.92 at an evaluation length of 4096. Similarly, on the Books3
 384 dataset, DAPE_{1×3}-Kerple trained with a length of 128 achieves a ppl of 31.07 at an evaluation length
 385 of 128 and 23.19 at an evaluation length of 4096, while RoPE trained with a length of 4096 achieves
 386 38.36 and 24.58, respectively. This suggests the superiority of the proposed DAPE_{1×3} with the
 387 introduced convolution operators among heads and neighboring tokens.

388 **With the same training tokens, compared to DAPE_{1×3} with longer training lengths, DAPE_{1×3}**
 389 **with shorter training lengths can achieve comparable performance, indicating that DAPE_{1×3}**
 390 **enhances the model’s understanding of text structure.** On the arXiv dataset, DAPE_{1×3}-Kerple
 391 with training lengths of 512 demonstrates performance close to that of training with a length of 4096
 392 when the evaluation length is 4096. Moreover, the performance curves for training lengths of 1024,
 393 and 2048 are almost identical. This trend is also observed with the Books3 dataset. These results
 394 indicate that DAPE_{1×3}-Kerple effectively helps the model comprehend text structure, enabling it to
 395 extend to longer lengths.
 396

397 **Transformers may overfit their training length: training on longer sequences may decrease**
 398 **performance when testing on shorter sequences.** On the arXiv dataset, DAPE_{1×3}-Kerple with
 399 a training length of 128 achieves the best performance when the evaluation length is 128. Similarly,
 400 DAPE_{1×3}-Kerple with training lengths of 256, 512, 1024, and 2048 achieves the best performance
 401 at evaluation lengths of 256, 512, 1024, and 2048, respectively. Also, on evaluation 128, the RoPE
 402 with training length 4096 and batch size 1 also achieves worse performance than the RoPE with
 403 training length 128 and batch size 32. This suggests that training on longer sequences may worsen
 404 a Transformer’s performance at shorter sequence lengths.

405 **DAPE_{1×3} can reduce the training time cost via larger batch size and shorter training length,**
 406 **achieving comparable performance compared to trained on longer length.** The cost of
 407 DAPE_{1×3} is $\mathcal{O}(B \cdot (h \cdot d \cdot T^2 + h \cdot D_{\text{DAPE}} \cdot T^2))$, where B, h, d, T and D_{DAPE} are the batch
 408 size, attention hidden dimension, attention head number, sequence length and DAPE hidden dimen-
 409 sion. By reducing the training length from T to $\frac{T}{K}$ and increasing the batch size from B to $B \cdot K$
 410 with the same training tokens, the cost becomes $\mathcal{O}(B \cdot K \cdot (h \cdot d \cdot (\frac{T}{K})^2 + h \cdot D_{\text{DAPE}} \cdot (\frac{T}{K})^2))$, which
 411 simplifies to $\mathcal{O}(\frac{B \cdot (h \cdot d \cdot T^2 + h \cdot D_{\text{DAPE}} \cdot T^2)}{K})$. For example, when the training length is 128 and the batch
 412 size is 32, the time cost of one step is 40.30ms. The time cost of length 256 (batch 16), length 512
 413 (batch 8), length 1024 (batch 4), and length 2048 (batch 2) are 42.61ms, 50.38ms, 79.36ms, and
 414 120.14ms. This reduction demonstrates the potential for significant training time savings.
 415

416 **4.3 THE EFFECT OF LARGER MODEL SIZE**



426 Figure 4: **The Effect of Larger Model Size 350M.** We show the results with training length 128 and training
 427 length 512 on Arxiv dataset.
 428

429 **DAPE_{1×3} performs well with larger model sizes, such as 350M and 2.7B.** As illustrated in
 430 Figure 4, the proposed DAPE_{1×3} shows superior performance at varying evaluation lengths with a
 431 model size of 350M. For a training length of 128, DAPE_{1×3}-Kerple achieves a perplexity (ppl) of
 7.63 at an evaluation length of 128 and 4.43 at an evaluation length of 8192, compared to DAPE’s

7.69 and 4.69, respectively. Similarly, for a training length of 512, DAPE_{1×3}-Kerple achieves a ppl of 4.10 at an evaluation length of 128 and 3.35 at an evaluation length of 8192, whereas DAPE achieves 4.14 and 3.44, respectively. We also present the 2.7B model size result in Appendix C. Therefore, the proposed DAPE_{1×3} demonstrates excellent performance with larger model sizes, showing the potential of including the proposed processing techniques in existing large language models.

4.4 THE EFFECT OF DAPE_{1×3}

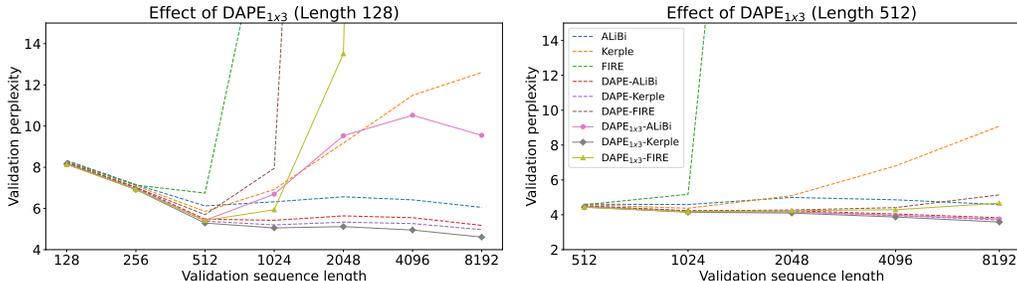


Figure 5: **The effect of DAPE_{1×3}.** Whatever the baseline is ALiBi, Kerple or FIRE, the proposed DAPE_{1×3} can all improve their performance. The Figure 1 also proves that the proposed DAPE_{1×3} is effective for NoPE and RoPE.

For Additive Positional Encoding, DAPE_{1×3} enhances performance within and beyond the training length. As demonstrated in Figure 5, for varying additive positional encoding such as ALiBi, Kerple, and FIRE, their incorporations with DAPE_{1×3} (i.e., DAPE_{1×3}-ALiBi, DAPE_{1×3}-Kerple, and DAPE_{1×3}-FIRE) consistently improve performance, **while DAPE_{1×3}-ALiBi may needs longer training length to achieve better performance than DAPE-ALiBi.** Furthermore, regardless of the specific additive positional encoding used, the proposed DAPE_{1×3} (configured with a kernel size of 1 × 3) outperforms the standard DAPE method (which employs a kernel size of 1 × 1). Also, as shown in Figure 1, DAPE_{1×3} improves the performance of NoPE, both within and beyond the training length. These results highlight the robustness and scalability of DAPE_{1×3}, suggesting its broad applicability in enhancing additive positional encoding frameworks.

For Non-Additive Positional Encoding, DAPE_{1×3} also improves performance within and beyond the training length. As illustrated in Figure 1, DAPE_{1×3} enhances the performance of RoPE, both within and beyond the training length. In contrast, naive DAPE reduces the performance of RoPE, with training lengths of 128 and 512. This indicates that the proposed DAPE_{1×3} is a versatile and widely applicable method with the potential to be applied to various position encoding techniques on the language modeling task.

4.5 THE PERFORMANCE OF DAPE_{1×3} WITH INFORMATION LEAKAGE

The DAPE_{1×3} can utilize attention data, which is supported by almost zero loss (perplexity is 1) under information leakage. To prevent the information leakage, we use the `torch.tril` before DAPE_{1×3} to make the attention score lower-triangular matrix. For the cheating version, we do not use the `torch.tril`. As shown in Figure 6, whatever DAPE_{1×3}-ALiBi, DAPE_{1×3}-Kerple or DAPE_{1×3}-FIRE, their cheating version can all achieve about zero loss within evaluation length 1024. Furthermore, the DAPE_{1×3}-Kerple can even achieve zero loss when the evaluation length is extended to 8096. This suggest that the proposed DAPE_{1×3} can really realize and utilize the information of attention score.

4.6 COMPARE DAPE AND DAPE_{1×3} WITH APPROXIMATE COMPUTATIONAL COST

DAPE_{1×3} achieves even better performance at a lower computational cost. As shown in Appendix E, when the training length is set to 128, DAPE_{1×3}-Kerple with D_{DAPE} as 10 achieves a perplexity (ppl) of 8.16 at an evaluation length of 128 and 4.74 at an evaluation length of 8192. This performance is notably better than that of DAPE-Kerple with D_{DAPE} as 64, which achieves

perplexities of 8.21 and 4.87, respectively. Moreover, when the training length is extended to 512 and the evaluation length is smaller or equal to 4096, DAPE_{1×3}-Kerple with D_{DAPE} as 10 continues to surpass the performance of DAPE-Kerple with D_{DAPE} as 64. Also, DAPE_{1×3}-Kerple with D_{DAPE} as 21 always achieves better performance than DAPE-Kerple with D_{DAPE} as 64. This demonstrates that DAPE_{1×3} not only maintains its performance advantage across different training lengths but also requires a lower computational cost.

4.7 THE PERFORMANCE WITH DIFFERENT KERNEL SIZES

Different experiment settings may have different optimal kernel sizes. Appendix F shows the performance of DAPE with various kernel sizes, including DAPE (equivalent to a 1×1 kernel size), DAPE_{1×3}, DAPE_{1×5}, and DAPE_{1×7}. For the Arxiv dataset, larger kernel sizes consistently achieve better performance, evaluating with training lengths of 128 or 512. However, for the Books3 dataset, DAPE_{1×3} performs best when the training length is 128 and evaluated at 8192, whereas DAPE_{1×5} performs best at the same evaluation level when the training length is 512. These results suggest that the optimal kernel size may vary depending on the experimental setting, ranging from 1×1 to larger kernel sizes. Although larger kernel sizes contribute to stronger expressiveness from intuition, we conjecture that the performance degradation for overly large kernel sizes results from optimization challenges.

4.8 THE PERFORMANCE ON CHE BENCHMARK WITH ACCURACY EVALUATION METRICS

Different tasks have different optimal kernel sizes, as shown in Appendix G and Appendix F. For example, on MISSING DUPLICATE task, the DAPE_{1×3}-Kerple improves the 87.57 of DAPE-Kerple to 99.65. However, on the STACK MANIPULATION task, the DAPE_{1×3}-Kerple decreases the 72.04 of DAPE-Kerple to 68.18. Also, as shown in Appendix F, the larger kernel size does not always lead to better performance. Overall, larger kernel size provides a potential way to improve the Transformer length extrapolation performance, and we usually could find a suitable kernel size (ranging from 1×1 to larger kernel sizes) to achieve better performance than without further processing attention score.

The large kernel size performance improvement is related to the baseline bias matrix. As shown in Appendix G, the best performance is usually achieved by further processing attention scores via kernel size 1 or 3. Moreover, on 11 permutation-variant tasks, the DAPE_{1×3}-Kerple achieves better performance on 8 of 11 tasks compared to Kerple. And the DAPE_{1×3}-FIRE achieves better performance on 6 of 11 tasks compared to FIRE. This suggests that the large kernel size performance improvement is related to the baseline bias matrix.

4.9 THE TIME COST

As the model size increases, the additional computational cost ratio gradually decreases. As shown in Appendix H, when the model size is 350M, the time cost for Kerple is 189.91 ms, while DAPE-Kerple takes 224.22 ms, and DAPE_{1×3}-Kerple requires 252.84 ms. Compared to DAPE_{1×3}-Kerple, the time cost ratios for Kerple and DAPE-Kerple are 0.7511 and 0.8868, respectively. As the model size increases from 350M to 2.7B and 6.7B, the time cost ratio for Kerple rises from 0.7511 to 0.8205 and 0.8918, respectively. Similarly, the time cost ratio for DAPE-Kerple increases from 0.8868 to 0.9361 and 0.9677. Therefore, as the model size increases, the time cost ratio also increases, indicating that the additional computational cost decreases progressively.

5 CONCLUSION

In this paper, we point out that the key of Transformer length extrapolation is the better and more accurate attention score. Therefore, we develop and analyze DAPE_{1×3} by processing the attention score as feature maps via convolution operation. Theoretically, we show that the associative recall tasks, which account for the most perplexity scores, can be realized by the proposed Transformer with convolution, in contrast to the vanilla Transformer. We conducted comprehensive experiments on Arxiv, Books3, and CHE to validate the effectiveness of the proposed method, where the proposed method exhibits significant superiority.

REFERENCES

- 540
541
542 Muhammad Adnan, Akhil Arunkumar, Gaurav Jain, Prashant J Nair, Ilya Soloveychik, and Pu-
543 rushotham Kamath. Keyformer: KV cache reduction through key tokens selection for efficient
544 generative inference. *arXiv preprint arXiv:2403.09054*, 2024.
- 545 Devanshu Agrawal, Shang Gao, and Martin Gajek. Can’t remember details in long documents? you
546 need some r&r. *arXiv preprint arXiv:2403.05004*, 2024.
- 547
548 Joshua Ainslie, Tao Lei, Michiel de Jong, Santiago Ontanon, Siddhartha Brahma, Yury Zemlyan-
549 ski, David Uthus, Mandy Guo, James Lee-Thorp, Yi Tay, et al. CoLT5: Faster long-range
550 transformers with conditional computation. In *The 2023 Conference on Empirical Methods in*
551 *Natural Language Processing*, 2023.
- 552 Dosovitskiy Alexey. An image is worth 16x16 words: Transformers for image recognition at scale.
553 *arXiv preprint arXiv: 2010.11929*, 2020.
- 554
555 Chenxin An, Fei Huang, Jun Zhang, Shansan Gong, Xipeng Qiu, Chang Zhou, and Lingpeng Kong.
556 Training-free long-context scaling of large language models. *arXiv preprint arXiv:2402.17463*,
557 2024.
- 558 Simran Arora, Sabri Eyuboglu, Aman Timalisina, Isys Johnson, Michael Poli, James Zou, Atri
559 Rudra, and Christopher Re. Zoology: Measuring and improving recall in efficient language
560 models. In *The Twelfth International Conference on Learning Representations*, 2024. URL
561 <https://openreview.net/forum?id=LY3ukUANko>.
- 562
563 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
564 learning to align and translate. *International Conference on Learning Representations*, 2015.
- 565 Maximilian Beck, Korbinian Pöppel, Markus Spanring, Andreas Auer, Oleksandra Prudnikova,
566 Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. xLSTM: Ex-
567 tended long short-term memory. *arXiv preprint arXiv:2405.04517*, 2024.
- 568
569 Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.
570 *arXiv preprint arXiv:2004.05150*, 2020.
- 571 Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a
572 transformer: A memory viewpoint. *Advances in Neural Information Processing Systems*, 36,
573 2024.
- 574
575 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
576 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
577 few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- 578 Guanzheng Chen, Xin Li, Zaiqiao Meng, Shangsong Liang, and Lidong Bing. CLEX: Continu-
579 ous length extrapolation for large language models. In *International Conference on Learning*
580 *Representations*, 2023a.
- 581
582 Junsong Chen, Jincheng YU, Chongjian GE, Lewei Yao, Enze Xie, Zhongdao Wang, James Kwok,
583 Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- α : Fast training of diffusion transformer
584 for photorealistic text-to-image synthesis. In *The Twelfth International Conference on Learning*
585 *Representations*, 2024. URL <https://openreview.net/forum?id=eAKmQP3ml>.
- 586
587 Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending context window
588 of large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*, 2023b.
- 589
590 Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. Lon-
591 gLoRA: Efficient fine-tuning of long-context large language models. *International Conference*
592 *on Learning Representations*, 2023c.
- 593
594 Ta-Chung Chi, Ting-Han Fan, Peter J Ramadge, and Alexander Rudnicky. KERPLE: Kernelized rel-
595 ative positional embedding for length extrapolation. *Advances in Neural Information Processing*
596 *Systems*, 35:8386–8399, 2022.

- 594 Ta-Chung Chi, Ting-Han Fan, Alexander Rudnicky, and Peter Ramadge. Dissecting transformer
595 length extrapolation via the lens of receptive field analysis. In *Proceedings of the 61st Annual*
596 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13522–
597 13537, 2023a.
- 598 Ta-Chung Chi, Ting-Han Fan, and Alexander I Rudnicky. Attention alignment and flexible positional
599 embeddings improve transformer length extrapolation. *arXiv preprint arXiv:2311.00684*, 2023b.
- 600
601 Hansul Cho, Jaeyoung Cha, Pranjal Awasthi, Srinadh Bhojanapalli, Anupam Gupta, and Chulhee
602 Yun. Position coupling: Leveraging task structure for improved length generalization of trans-
603 formers. *arXiv preprint arXiv:2405.20671*, 2024.
- 604
605 Noam Chomsky. Three models for the description of language. *IRE Transactions on Information*
606 *Theory*, 2(3):113–124, 1956.
- 607
608 Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea
609 Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser,
610 David Benjamin Belanger, Lucy J Colwell, and Adrian Weller. Rethinking attention with per-
611 formers. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=Ua6zuk0WRH>.
- 612
613 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
614 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
615 Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):
616 1–113, 2023.
- 617
618 Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov.
619 Transformer-XL: Attentive language models beyond a fixed-length context. In *Proceedings of the*
620 *57th Annual Meeting of the Association for Computational Linguistics*, pp. 2978–2988, 2019.
- 621
622 Soham De, Samuel L Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Al-
623 bert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, et al. Griffin: Mix-
624 ing gated linear recurrences with local attention for efficient language models. *arXiv preprint*
arXiv:2402.19427, 2024.
- 625
626 Gregoire Deletang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt,
627 Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, et al. Neural networks and the chomsky
628 hierarchy. In *International Conference on Learning Representations*, 2022.
- 629
630 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep
631 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*
632 *the North American Chapter of the Association for Computational Linguistics: Human Language*
Technologies, Volume 1 (Long and Short Papers), pp. 4171–4186, 2019.
- 633
634 Hantian Ding, Zijian Wang, Giovanni Paolini, Varun Kumar, Anoop Deoras, Dan Roth, and Stefano
635 Soatto. Fewer truncations improve language modeling. *arXiv preprint arXiv:2404.10830*, 2024a.
- 636
637 Yiran Ding, Li Lyna Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan
638 Yang, and Mao Yang. LongRoPE: Extending llm context window beyond 2 million tokens. *arXiv*
preprint arXiv:2402.13753, 2024b.
- 639
640 Zafeirios Fountas, Martin A Benfeghoul, Adnan Omerjee, Fenia Christopoulou, Gerasimos Lam-
641 pouras, Haitham Bou-Ammar, and Jun Wang. Human-like episodic memory for infinite context
642 llms. *arXiv preprint arXiv:2407.09450*, 2024.
- 643
644 Daniel Y Fu, Tri Dao, Khaled K Saab, Armin W Thomas, Atri Rudra, and Christopher Ré.
645 Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint*
arXiv:2212.14052, 2022.
- 646
647 Yao Fu, Rameswar Panda, Xinyao Niu, Xiang Yue, Hannaneh Hajishirzi, Yoon Kim, and Hao Peng.
Data engineering for scaling language models to 128k context. *arXiv preprint arXiv:2402.10171*,
2024.

- 648 Chaochen Gao, Xing Wu, Qi Fu, and Songlin Hu. Quest: Query-centric data synthesis approach for
649 long-context scaling of large language model. *arXiv preprint arXiv:2405.19846*, 2024.
650
- 651 Nicholas Geneva and Nicholas Zabarar. Transformers for modeling physical systems. *Neural Net-*
652 *works*, 146:272–289, 2022.
- 653 Olga Golovneva, Tianlu Wang, Jason Weston, and Sainbayar Sukhbaatar. Contextual position en-
654 coding: Learning to count what’s important. *arXiv preprint arXiv:2405.18719*, 2024.
655
- 656 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
657 *preprint arXiv:2312.00752*, 2023.
658
- 659 Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung, and
660 Yinfei Yang. LongT5: Efficient text-to-text transformer for long sequences. *Findings of the*
661 *Association for Computational Linguistics: NAACL*, 2022.
- 662 Adi Haviv, Ori Ram, Ofir Press, Peter Izsak, and Omer Levy. Transformer language models with-
663 out positional encodings still learn positional information. In *Findings of the Association for*
664 *Computational Linguistics: EMNLP 2022*, pp. 1382–1390, 2022.
665
- 666 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
667 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
668 770–778, 2016.
- 669 Zhenyu He, Guhao Feng, Shengjie Luo, Kai Yang, Di He, Jingjing Xu, Zhi Zhang, Hongxia Yang,
670 and Liwei Wang. Two stones hit one bird: Bilevel positional encoding for better length extrapo-
671 lation. *arXiv preprint arXiv:2401.16421*, 2024.
672
- 673 Peyman Hosseini, Ignacio Castro, Iacopo Ghinassi, and Matthew Purver. Efficient solutions for an
674 intriguing failure of llms: Long context window does not mean llms can analyze long sequences
675 flawlessly. *arXiv preprint arXiv:2408.01866*, 2024.
676
- 677 Zhiyuan Hu, Yuliang Liu, Jinman Zhao, Suyuchen Wang, Yan Wang, Wei Shen, Qing Gu, Anh Tuan
678 Luu, See-Kiong Ng, Zhiwei Jiang, et al. Longrecipe: Recipe for efficient long context general-
679 ization in large language models. *arXiv preprint arXiv:2409.00509*, 2024.
- 680 Hongye Jin, Xiaotian Han, Jingfeng Yang, Zhimeng Jiang, Zirui Liu, Chia-Yuan Chang, Huiyuan
681 Chen, and Xia Hu. LLM maybe LongLM: Self-extend LLM context window without tuning.
682 *arXiv preprint arXiv:2401.01325*, 2024.
683
- 684 Osman Semih Kayhan and Jan C van Gemert. On translation invariance in cnns: Convolutional
685 layers can exploit absolute spatial location. In *Proceedings of the IEEE/CVF Conference on*
686 *Computer Vision and Pattern Recognition*, pp. 14274–14285, 2020.
- 687 Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva
688 Reddy. The impact of positional encoding on length generalization in transformers. *Advances*
689 *in Neural Information Processing Systems*, 36, 2024.
690
- 691 Guolin Ke, Di He, and Tie-Yan Liu. Rethinking positional encoding in language pre-training. In
692 *International Conference on Learning Representations*, 2020.
693
- 694 Shun Kiyono, Sosuke Kobayashi, Jun Suzuki, and Kentaro Inui. SHAPE: Shifted absolute position
695 embedding for transformers. In *Proceedings of the 2021 Conference on Empirical Methods in*
696 *Natural Language Processing*, pp. 3309–3321, 2021.
- 697 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
698 lutional neural networks. *Advances in neural information processing systems*, 25, 2012.
699
- 700 Moshe Leshno, Vladimir Ya Lin, Allan Pinkus, and Shimon Schocken. Multilayer feedforward net-
701 works with a nonpolynomial activation function can approximate any function. *Neural Networks*,
6(6):861–867, 1993.

- 702 Jingyao Li, Pengguang Chen, Zexin He, Shaozuo Yu, Shu Liu, and Jiaya Jia. Rethinking out-of-
703 distribution (OOD) detection: Masked image modeling is all you need. In *Proceedings of the*
704 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11578–11589,
705 June 2023a.
- 706
707 Jingyao Li, Pengguang Chen, Shengju Qian, and Jiaya Jia. Tagclip: Improving discrimination ability
708 of open-vocabulary semantic segmentation, 2023b.
- 709
710 Jingyao Li, Pengguang Chen, and Jiaya Jia. Motcoder: Elevating large language models with mod-
711 ular of thought for challenging programming tasks, 2024a.
- 712
713 Jingyao Li, Pengguang Chen, Shaozuo Yu, Shu Liu, and Jiaya Jia. Bal: Balancing diversity and
714 novelty for active learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46
(5):3653–3664, 2024b. doi: 10.1109/TPAMI.2023.3345844.
- 715
716 Shanda Li, Chong You, Guru Guruganesh, Joshua Ainslie, Santiago Ontanon, Manzil Zaheer, Sumit
717 Sanghai, Yiming Yang, Sanjiv Kumar, and Srinadh Bhojanapalli. Functional interpolation for
718 relative positions improves long context transformers. In *International Conference on Learning*
719 *Representations*, 2023c.
- 720
721 Zhenyu Li, Yike Zhang, Tengyu Pan, Yutao Sun, Zhichao Duan, Junjie Fang, Rong Han, Zixuan
722 Wang, and Jianyong Wang. Focuslm: Scaling llm’s context by parallel decoding. *arXiv preprint*
arXiv:2408.11745, 2024c.
- 723
724 Zihan Liao, Jun Wang, Hang Yu, Lingxiao Wei, Jianguo Li, and Wei Zhang. E2llm: Encoder
725 elongated large language models for long-context understanding and reasoning. *arXiv preprint*
726 *arXiv:2409.06679*, 2024.
- 727
728 Opher Lieber, Barak Lenz, Hofit Bata, Gal Cohen, Jhonathan Osin, Itay Dalmedigos, Erez Safahi,
729 Shaked Meir, Yonatan Belinkov, Shai Shalev-Shwartz, et al. Jamba: A hybrid transformer-
mamba language model. *arXiv preprint arXiv:2403.19887*, 2024.
- 730
731 Tatiana Likhomanenko, Qiantong Xu, Gabriel Synnaeve, Ronan Collobert, and Alex Rogozhnikov.
732 CAPE: Encoding relative positions with continuous augmented positional embeddings. *Advances*
733 *in Neural Information Processing Systems*, 34:16079–16092, 2021.
- 734
735 Bin Lin, Tao Peng, Chen Zhang, Minmin Sun, Lanbo Li, Hanyu Zhao, Wencong Xiao, Qi Xu, Xiafei
736 Qiu, Shen Li, et al. Infinite-LLM: Efficient LLM service for long context with distattention and
distributed kvcache. *arXiv preprint arXiv:2401.02669*, 2024a.
- 737
738 Hongzhan Lin, Ang Lv, Yuhan Chen, Chen Zhu, Yang Song, Hengshu Zhu, and Rui Yan. Mixture
739 of in-context experts enhance llms’ long context awareness. *arXiv preprint arXiv:2406.19598*,
740 2024b.
- 741
742 Jiacheng Liu, Sewon Min, Luke Zettlemoyer, Yejin Choi, and Hannaneh Hajishirzi. Infini-
743 gram: Scaling unbounded n-gram language models to a trillion tokens. *arXiv preprint*
arXiv:2401.17377, 2024a.
- 744
745 Jiaheng Liu, Zhiqi Bai, Yuanxing Zhang, Chenchen Zhang, Yu Zhang, Ge Zhang, Jiakai Wang,
746 Haoran Que, Yukang Chen, Wenbo Su, et al. E²-LLM: Efficient and extreme length extension
747 of large language models. *arXiv preprint arXiv:2401.06951*, 2024b.
- 748
749 Xiaoran Liu, Hang Yan, Chenxin An, Xipeng Qiu, and Dahua Lin. Scaling laws of RoPE-based
extrapolation. In *International Conference on Learning Representations*, 2023.
- 750
751 Xuanqing Liu, Hsiang-Fu Yu, Inderjit Dhillon, and Cho-Jui Hsieh. Learning to encode position for
752 transformer with continuous dynamical model. In *International Conference on Machine Learn-*
753 *ing*, pp. 6327–6335. PMLR, 2020.
- 754
755 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the*
IEEE/CVF international conference on computer vision, pp. 10012–10022, 2021a.

- 756 Zhuang Liu, Degen Huang, Kaiyu Huang, Zhuang Li, and Jun Zhao. Finbert: A pre-trained financial
757 language representation model for financial text mining. In *Proceedings of the Twenty-ninth*
758 *International Conference on International Joint Conferences on Artificial Intelligence*, pp. 4513–
759 4519, 2021b.
- 760 Shengjie Luo, Shanda Li, Tianle Cai, Di He, Dinglan Peng, Shuxin Zheng, Guolin Ke, Liwei Wang,
761 and Tie-Yan Liu. Stable, fast and accurate: Kernelized attention with relative positional encoding.
762 *Advances in Neural Information Processing Systems*, 34:22795–22807, 2021.
- 763
764 Shengjie Luo, Shanda Li, Shuxin Zheng, Tie-Yan Liu, Liwei Wang, and Di He. Your transformer
765 may not be as powerful as you expect. *Advances in Neural Information Processing Systems*, 35:
766 4301–4315, 2022.
- 767 Xindian Ma, Wenyuan Liu, Peng Zhang, and Nan Xu. 3d-rpe: Enhancing long-context modeling
768 through 3d rotary position encoding. *arXiv preprint arXiv:2406.09897*, 2024.
- 769
770 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
771 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style,
772 high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32,
773 2019.
- 774 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
775 *the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- 776
777 Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Huanqi Cao, Xin
778 Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, et al. RWKV: Reinventing RNNs
779 for the transformer era. *Findings of the Association for Computational Linguistics: EMNLP*,
780 2023a.
- 781 Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. YaRN: Efficient context window
782 extension of large language models. In *International Conference on Learning Representations*,
783 2023b.
- 784 Ofir Press, Noah Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables
785 input length extrapolation. In *International Conference on Learning Representations*, 2021.
- 786
787 Zhen Qin, Weigao Sun, Dong Li, Xuyang Shen, Weixuan Sun, and Yiran Zhong. Lightning
788 attention-2: A free lunch for handling unlimited sequence lengths in large language models. *arXiv*
789 *preprint arXiv:2401.04658*, 2024a.
- 790 Zhen Qin, Yiran Zhong, and Hui Deng. Exploring transformer extrapolation. In *Proceedings of the*
791 *AAAI Conference on Artificial Intelligence*, volume 38, pp. 18897–18905, 2024b.
- 792
793 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
794 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 795
796 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
797 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
798 transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- 799 Anian Ruoss, Grégoire Delétang, Tim Genewein, Jordi Grau-Moya, Róbert Csordás, Mehdi Ben-
800 nani, Shane Legg, and Joel Veness. Randomized positional encodings boost length generalization
801 of transformers. In *Proceedings of the 61st Annual Meeting of the Association for Computational*
802 *Linguistics (Volume 2: Short Papers)*, pp. 1889–1903, 2023.
- 803
804 Mahdi Sabbaghi, George Pappas, Hamed Hassani, and Surbhi Goel. Explicitly encoding structural
805 symmetry is key to length generalization in arithmetic tasks. *arXiv preprint arXiv:2406.01895*,
806 2024.
- 807 Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representa-
808 tions. In *Proceedings of the 2018 Conference of the North American Chapter of the Association*
809 *for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp.
464–468, 2018.

- 810 Noam Shazeer, Zhenzhong Lan, Youlong Cheng, Nan Ding, and Le Hou. Talking-heads attention.
811 *arXiv preprint arXiv:2003.02436*, 2020.
812
- 813 Weijia Shi, Sewon Min, Maria Lomeli, Chunting Zhou, Margaret Li, Victoria Lin, Noah A Smith,
814 Luke Zettlemoyer, Scott Yih, and Mike Lewis. In-context pretraining: Language modeling be-
815 yond document boundaries. *International Conference on Learning Representations*, 2023.
816
- 817 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image
818 recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 819 Matt Stallone, Vaibhav Saxena, Leonid Karlinsky, Bridget McGinn, Tim Bula, Mayank Mishra,
820 Adriana Meza Soria, Gaoyuan Zhang, Aditya Prasad, Yikang Shen, et al. Scaling granite code
821 models to 128k context. *arXiv preprint arXiv:2407.13739*, 2024.
822
- 823 Konrad Staniszewski, Szymon Tworkowski, Sebastian Jaszczur, Henryk Michalewski, Łukasz
824 Kuciński, and Piotr Miłoś. Structured packing in LLM training improves long context utiliza-
825 tion. *arXiv preprint arXiv:2312.17296*, 2023.
- 826 Jianlin Su, Murtadha Ahmed, Luo Ao, Mingren Zhu, Yunfeng Liu, et al. Naive bayes-based context
827 extension for large language models. *arXiv preprint arXiv:2403.17552*, 2024a.
828
- 829 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-
830 hanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024b.
831
- 832 Jiankai Sun, Chuanyang Zheng, Enze Xie, Zhengying Liu, Ruihang Chu, Jianing Qiu, Jiaqi Xu,
833 Mingyu Ding, Hongyang Li, Mengzhe Geng, et al. A survey of reasoning with foundation models.
834 *arXiv preprint arXiv:2312.11562*, 2023a.
- 835 Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav Chaud-
836 hary, Xia Song, and Furu Wei. A length-extrapolatable transformer. *Proceedings of the*
837 *61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Pa-*
838 *pers)*, pp. 14590–14604, July 2023b. doi: 10.18653/v1/2023.acl-long.816. URL [https://](https://aclanthology.org/2023.acl-long.816)
839 aclanthology.org/2023.acl-long.816.
- 840 Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question
841 answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference*
842 *of the North American Chapter of the Association for Computational Linguistics: Human Lan-*
843 *guage Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, 2019.
844
- 845 Yi Tay, Dara Bahri, Liu Yang, Donald Metzler, and Da-Cheng Juan. Sparse sinkhorn attention. In
846 *International Conference on Machine Learning*, pp. 9438–9447. PMLR, 2020.
847
- 848 Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia,
849 Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for
850 science. *arXiv preprint arXiv:2211.09085*, 2022.
- 851 Junfeng Tian, Da Zheng, Yang Cheng, Rui Wang, Colin Zhang, and Debing Zhang. Untie the knots:
852 An efficient data augmentation strategy for long-context pre-training in language models. *arXiv*
853 *preprint arXiv:2409.04774*, 2024.
854
- 855 Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and
856 Hervé Jégou. Training data-efficient image transformers & distillation through attention. In
857 *International conference on machine learning*, pp. 10347–10357. PMLR, 2021.
- 858 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
859 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and
860 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
861
- 862 Szymon Tworkowski, Konrad Staniszewski, Mikołaj Patek, Yuhuai Wu, Henryk Michalewski, and
863 Piotr Miłoś. Focused transformer: Contrastive training for context scaling. *Advances in Neural*
Information Processing Systems, 36, 2024.

- 864 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
865 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Informa-*
866 *tion Processing Systems*, 30, 2017.
- 867
- 868 Benyou Wang, Lifeng Shang, Christina Lioma, Xin Jiang, Hao Yang, Qun Liu, and Jakob Grue
869 Simonsen. On position embeddings in BERT. In *International Conference on Learning Repre-*
870 *sentations*, 2020.
- 871
- 872 Huadong Wang, Xin Shen, Mei Tu, Yimeng Zhuang, and Zhiyuan Liu. Improved transformer with
873 multi-head dense collaboration. *IEEE/ACM Transactions on Audio, Speech, and Language Pro-*
874 *cessing*, 30:2754–2767, 2022.
- 875
- 876 Jie Wang, Tao Ji, Yuanbin Wu, Hang Yan, Tao Gui, Qi Zhang, Xuanjing Huang, and Xiaoling
877 Wang. Length generalization of causal transformers without position encoding. *arXiv preprint*
arXiv:2404.12224, 2024a.
- 878
- 879 Suyuchen Wang, Ivan Kobyzev, Peng Lu, Mehdi Rezagholizadeh, and Bang Liu. Resonance
880 RoPE: Improving context length generalization of large language models. *arXiv preprint*
arXiv:2403.00071, 2024b.
- 881
- 882 Y Wang, D Ma, and D Cai. With greater text comes greater necessity: Inference-time training helps
883 long text generation. *arXiv preprint arXiv:2401.11504*, 2024c.
- 884
- 885 BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić,
886 Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. Bloom:
887 A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*,
888 2022.
- 889
- 890 Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prab-
891 hanjan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model
for finance. *arXiv preprint arXiv:2303.17564*, 2023.
- 892
- 893 Tong Wu, Yanpeng Zhao, and Zilong Zheng. Never miss a beat: An efficient recipe for context win-
894 dow extension of large language models with consistent” middle” enhancement. *arXiv preprint*
arXiv:2406.07138, 2024a.
- 895
- 896 Wenhao Wu, Yizhong Wang, Yao Fu, Xiang Yue, Dawei Zhu, and Sujian Li. Long context alignment
897 with short instructions and synthesized positions. *arXiv preprint arXiv:2405.03939*, 2024b.
- 898
- 899 Chaojun Xiao, Pengle Zhang, Xu Han, Guangxuan Xiao, Yankai Lin, Zhengyan Zhang, Zhiyuan
900 Liu, Song Han, and Maosong Sun. InfLLM: Unveiling the intrinsic capacity of LLMs for under-
901 standing extremely long sequences with training-free memory. *arXiv preprint arXiv:2402.04617*,
902 2024a.
- 903
- 904 Da Xiao, Qingye Meng, Shengping Li, and Xingyuan Yuan. Improving transformers with dynami-
905 cally composable multi-head attention. *arXiv preprint arXiv:2405.08553*, 2024b.
- 906
- 907 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
908 language models with attention sinks. In *International Conference on Learning Representations*,
2024c. URL <https://openreview.net/forum?id=NG7sS51zVF>.
- 909
- 910 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
911 language models with attention sinks. In *The Twelfth International Conference on Learning Rep-*
resentations, 2024d. URL <https://openreview.net/forum?id=NG7sS51zVF>.
- 912
- 913 Wenhan Xiong, Jingyu Liu, Igor Molybog, Hejia Zhang, Prajjwal Bhargava, Rui Hou, Louis Martin,
914 Rashi Rungta, Karthik Abinav Sankararaman, Barlas Oguuz, et al. Effective long-context scaling
915 of foundation models. *arXiv preprint arXiv:2309.16039*, 2023.
- 916
- 917 Peng Xu, Wei Ping, Xianchao Wu, Zihan Liu, Mohammad Shoeybi, and Bryan Catanzaro. Chatqa
2: Bridging the gap to proprietary llms in long context and rag capabilities. *arXiv preprint*
arXiv:2407.14482, 2024.

- 918 Kai Yang, Jan Ackermann, Zhenyu He, Guhao Feng, Bohang Zhang, Yunzhen Feng, Qiwei Ye,
919 Di He, and Liwei Wang. Do efficient transformers really save computation? *International Con-*
920 *ference on Machine Learning*, 2024.
- 921
922 Howard Yen, Tianyu Gao, and Danqi Chen. Long-context language modeling with parallel context
923 encoding. *arXiv preprint arXiv:2402.16617*, 2024.
- 924 Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. PEGASUS: Pre-training with ex-
925 tracted gap-sentences for abstractive summarization. In *International Conference on Machine*
926 *Learning*, pp. 11328–11339. PMLR, 2020.
- 927
928 Zhenyu Zhang, Runjin Chen, Shiwei Liu, Zhewei Yao, Olatunji Ruwase, Beidi Chen, Xiaoxia Wu,
929 and Zhangyang Wang. Found in the middle: How language models use long contexts better via
930 plug-and-play positional encoding. *arXiv preprint arXiv:2403.04797*, 2024.
- 931 Liang Zhao, Xiaocheng Feng, Xiachong Feng, Bin Qin, and Ting Liu. Length extrapolation of trans-
932 formers: A survey from the perspective of position encoding. *arXiv preprint arXiv:2312.17044*,
933 2023.
- 934 Liang Zhao, Tianwen Wei, Liang Zeng, Cheng Cheng, Liu Yang, Peng Cheng, Lijie Wang, Chenxia
935 Li, Xuejie Wu, Bo Zhu, et al. Longskywork: A training recipe for efficiently extending context
936 length in large language models. *arXiv preprint arXiv:2406.00605*, 2024.
- 937
938 Chuanyang Zheng, Zhengying Liu, Enze Xie, Zhenguo Li, and Yu Li. Progressive-hint prompting
939 improves reasoning in large language models. *arXiv preprint arXiv:2304.09797*, 2023.
- 940
941 Chuanyang Zheng, Yihang Gao, Han Shi, Minbin Huang, Jingyao Li, Jing Xiong, Xiaozhe Ren,
942 Michael Ng, Xin Jiang, Zhenguo Li, et al. Dape: Data-adaptive positional encoding for length
943 extrapolation. *Advances in Neural Information Processing Systems*, 2024.
- 944
945 Hattie Zhou, Arwen Bradley, Etai Littwin, Noam Razin, Omid Saremi, Joshua M. Susskind, Samy
946 Bengio, and Preetum Nakkiran. What algorithms can transformers learn? a study in length
947 generalization. In *International Conference on Learning Representations*, 2024a. URL [https://
openreview.net/forum?id=AssIuHnmHX](https://openreview.net/forum?id=AssIuHnmHX).
- 948
949 Jin Peng Zhou, Charles E Staats, Wenda Li, Christian Szegedy, Kilian Q Weinberger, and Yuhuai
950 Wu. Don’t trust: Verify – grounding LLM quantitative reasoning with autoformalization. In
951 *The Twelfth International Conference on Learning Representations*, 2024b. URL [https://
openreview.net/forum?id=V5tdi14ple](https://openreview.net/forum?id=V5tdi14ple).
- 952
953 Yongchao Zhou, Uri Alon, Xinyun Chen, Xuezhi Wang, Rishabh Agarwal, and Denny Zhou. Trans-
954 formers can achieve length generalization but not robustly. *arXiv preprint arXiv:2402.09371*,
955 2024c.
- 956
957 Dawei Zhu, Nan Yang, Liang Wang, Yifan Song, Wenhao Wu, Furu Wei, and Sujian Li. PoSE:
958 Efficient context window extension of llms via positional skip-wise training. In *International
Conference on Learning Representations*, 2023a.
- 959
960 Qianchao Zhu, Jiangfei Duan, Chang Chen, Siran Liu, Xiuhong Li, Guanyu Feng, Xin Lv, Huanqi
961 Cao, Xiao Chuanfu, Xingcheng Zhang, et al. Near-lossless acceleration of long context llm
962 inference with adaptive structured sparse attention. *arXiv preprint arXiv:2406.15486*, 2024.
- 963
964 Shiyi Zhu, Jing Ye, Wei Jiang, Qi Zhang, Yifan Wu, and Jianguo Li. CoCA: Fusing position em-
965 bedding with collinear constrained attention for fine-tuning free context window extending. *arXiv
e-prints*, pp. arXiv–2309, 2023b.
- 966
967
968
969
970
971

A Δ PERPLEXITY FOR LENGTH EXTRAPOLATION EVALUATIONTable 1: The ΔP on Book dataset with training length 512, compared to baselines.

Method	RoPE	ALiBi	Kerple	DAPE-Kerple	DAPE _{1×3} -Kerple
$P(M(x_{512}), T_{train})$	19.74	20.04	19.83	19.25	18.95
$P(M(x_{1024}), T_{train})$	261.39	19.74	19.19	18.28	17.92
$P(M(x_{1024}[-T_{train}] :), T_{train})$	19.51	19.79	19.58	19.03	18.74
ΔP_{1024}	-241.88	0.05	0.39	0.75	0.82
$P(M(x_{2048}), T_{train})$	411.23	20.17	20.48	17.20	16.79
$P(M(x_{2048}[-T_{train}] :), T_{train})$	18.74	19.03	19.84	18.28	18.01
ΔP_{2048}	-392.49	-1.14	-0.64	1.08	1.22
$P(M(x_{4096}), T_{train})$	635.80	20.50	28.33	17.58	17.05
$P(M(x_{4096}[-T_{train}] :), T_{train})$	19.11	19.35	19.07	18.59	18.19
ΔP_{4096}	-616.69	-1.15	-9.26	1.01	1.14
$P(M(x_{8192}), T_{train})$	762.86	21.30	40.94	17.85	17.20
$P(M(x_{8192}[-T_{train}] :), T_{train})$	19.78	20.02	19.85	19.38	18.98
ΔP_{8192}	-743.08	-1.28	-21.09	1.53	1.78

In this discussion, we explore how to effectively use perplexity as a metric, incorporating concepts of information gain and entropy. Let $L(\cdot)$ represent the process for calculating loss, and $M(x)$ denote the logit output generated by the model after processing an input sequence x . For evaluating model performance, we define $P(M(x), K)$ as follows:

1. Process the entire sequence x using $M(x)$.
2. Compute the perplexity on the last K tokens of the sequence.

To interpret information gain, we consider the training sequence length T_{train} . Given an input x , we calculate the change in loss/perplexity, ΔP , as:

$$\Delta P = P(M(x[-T_{train} :]), T_{train}) - P(M(x), T_{train}) \quad (6)$$

The term ΔP provides insights into the model’s information gain relative to local and global context, allowing us to quantify entropy in terms of model uncertainty reduction. We interpret ΔL as follows:

- When $\Delta P = 0$: The model’s information gain from the full sequence is negligible, indicating an entropy level comparable to local attention (e.g., models like ALiBi when the evaluation length is 1024). This suggests the model does not leverage context beyond a limited range.
- When $\Delta P < 0$: Processing the entire sequence increases entropy, resulting in worse performance than focusing only on the last T_{train} tokens. This implies negative information gain and limited extrapolation capability (e.g. such as RoPE), as the model may overfit to recent tokens without capturing broader context effectively.
- When $\Delta P > 0$: The model benefits from the information within $x[: T_{train}]$, achieving a reduction in entropy that reflects positive information gain. This suggests the model leverages contextual information beyond the training sequence, indicating extrapolation capability.
- Our suggestion of bias matrix. The Kerple is a good choice for almost all settings, the FIRE may need longer training length/tokens to present its ability, and do not use ALiBi unless necessary . It is easy to train Kerple, as Kerple usually has few trainable parameters compared to FIRE. If you do not know which one to use, directly use Kerple. FIRE may have better performance, but may need longer training length (diverges at 128 but works well at 512, with DAPE). FIRE $b(i, j) = f_{\theta} \left(\frac{\psi(i-j)}{\psi(\max\{L, i\})} \right)$ so that we may need longer training length or more training tokens to well-train the neural network f_{θ} . Do not use ALiBi unless necessary. The ALiBi will quickly become local attention as the sequence length increases.

1026 By examining ΔP , we can evaluate the model's ability to reduce entropy and gain information from
1027 extended sequences, providing a measure of its extrapolative power.
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

B COMPARE WITH BASELINE ON ARXIV DATASET WITH TRAINING LENGTH 1024

Table 2: The performance (ppl) on Arxiv dataset with training length 1024, compared to baselines.

Method	1024	2048	4096	8192
NoPE (Kazemnejad et al., 2024)	4.16	42.27	1854.73	17167.32
RoPE (Su et al., 2024b)	4.07	86.20	237.67	256.12
T5’s bias (Raffel et al., 2020)	4.03	4.28	13.07	79.55
ALiBi (Press et al., 2021)	4.09	4.53	4.45	4.22
Kerple (Chi et al., 2022)	4.06	4.09	4.68	6.951
FIRE (Li et al., 2023c)	4.06	9.21	236.18	440.60
DAPE-Kerple (Zheng et al., 2024)	3.98	3.91	3.68	3.41
DAPE _{1×3} -Kerple	3.93	3.86	3.61	3.37

C LARGE MODEL SIZE

Table 3: The performance (ppl) under large model size 2.7B on Books3 dataset.

Method	512	1024	2048	4096
RoPE	21.01	25.00	48.13	160.59
T5’s bias	21.10	21.88	23.59	33.23
Kerple	21.14	22.08	23.38	27.21
DAPE-Kerple	20.52	21.01	20.23	19.67
DAPE _{1×3} -Kerple (kernel size 1x3)	20.16	20.54	19.80	19.02

D THE PERFORMANCE OF $\text{DAPE}_{1 \times 3}$ WITH INFORMATION LEAKAGE

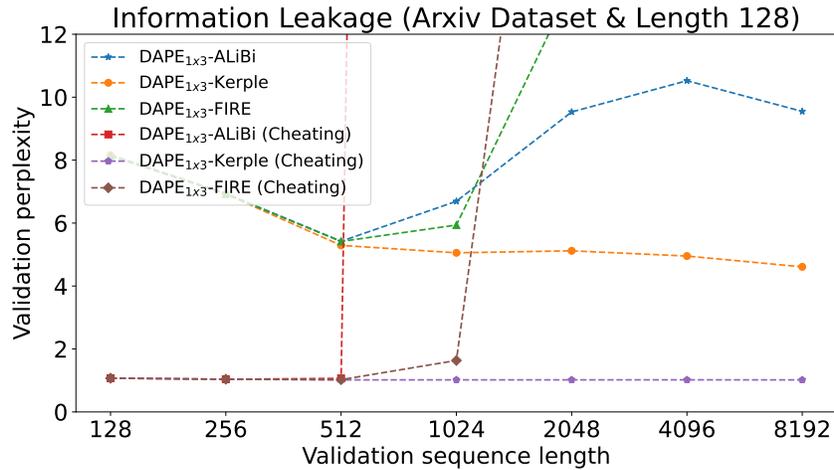


Figure 6: **Result with information leakage.**

E COMPARE $\text{DAPE}_{1 \times 3}$ AND DAPE WITH APPROXIMATE COMPUTATIONAL COST

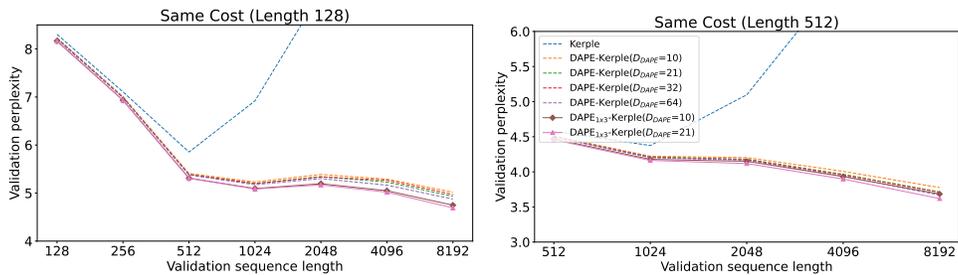


Figure 7: **Compare $\text{DAPE}_{1 \times 3}$ and DAPE with the approximately same cost on Arxiv Dataset.** We compare the $\text{DAPE}_{1 \times 3}$ and DAPE with approximate cost and different D_{DAPE} . As the kernel size of is $\text{DAPE}_{1 \times 3}$ 1×3 , the proposed $\text{DAPE}_{1 \times 3}$ is the triple computation cost of DAPE , with the same D_{DAPE} .

F THE PERFORMANCE WITH DIFFERENT KERNEL SIZE

Table 4: The performance with different kernel sizes, with training length 128 and evaluation from length 128 to 8192. For different datasets and training length, the optimal kernel size may not always be the largest one, especially when the evaluation length is larger.

Dataset	Method	128	256	512	1024	2048	4096	8192
Arxiv	Kerple	8.30	7.10	5.85	6.91	9.17	11.48	12.59
	DAPE-Kerple (Kernel Size 1x1)	8.21	6.98	5.38	5.20	5.33	5.26	4.97
	DAPE _{1x3} -Kerple (Kernel Size 1x3)	8.15	6.92	5.29	5.05	5.11	4.95	4.60
	DAPE _{1x5} -Kerple (Kernel Size 1x5)	8.13	6.91	5.27	5.04	5.10	4.91	4.57
	DAPE _{1x7} -Kerple (Kernel Size 1x7)	8.12	6.89	5.26	5.02	5.09	4.91	4.57
Books3	Kerple	32.10	29.09	28.10	35.75	44.68	56.39	66.23
	DAPE-Kerple (Kernel Size 1x1)	31.49	28.27	24.93	24.31	23.34	24.38	25.01
	DAPE _{1x3} -Kerple (Kernel Size 1x3)	31.07	27.81	24.38	23.57	22.40	23.19	23.52
	DAPE _{1x5} -Kerple (Kernel Size 1x5)	31.02	27.79	24.36	23.57	22.41	23.32	23.71
	DAPE _{1x7} -Kerple (Kernel Size 1x7)	30.98	27.76	24.31	23.47	22.30	23.00	23.57

Table 5: The performance with different kernel size, with training length 512 and evaluation from length 512 to 8192. For different datasets and training length, the optimal kernel size may not always be the largest one, especially when the evaluation length is larger.

Dataset	Method	512	1024	2048	4096	8192
Arxiv	Kerple	4.57	4.37	5.09	6.80	9.08
	DAPE-Kerple (Kernel Size 1x1)	4.49	4.20	4.17	3.95	3.70
	DAPE _{1x3} -Kerple (Kernel Size 1x3)	4.44	4.14	4.09	3.87	3.58
	DAPE _{1x5} -Kerple (Kernel Size 1x5)	4.44	4.14	4.10	3.85	3.59
	DAPE _{1x7} -Kerple (Kernel Size 1x7)	4.43	4.13	4.08	3.85	3.57
Books3	Kerple	19.83	19.19	20.48	28.33	40.94
	DAPE-Kerple (Kernel Size 1x1)	19.25	18.28	17.20	17.58	17.85
	DAPE _{1x3} -Kerple (Kernel Size 1x3)	18.95	17.92	16.79	17.05	17.20
	DAPE _{1x5} -Kerple (Kernel Size 1x5)	18.89	17.87	16.76	17.09	17.10
	DAPE _{1x7} -Kerple (Kernel Size 1x7)	18.86	17.82	16.70	17.01	17.16

G THE PERFORMANCE OF DAPE_{1×3} ON CHE BENCHMARK

Table 6: Train on length 40 with 200k steps, and test from lengths 41 to 500. The random accuracy is 50%, except for MODULAR ARITHMETIC (SIMPLE), CYCLE NAVIGATION, BUCKET SORT, SOLVE EQUATION and MODULAR ARITHMETIC, where it is 20%. ††† denotes permutation-invariant tasks, which are expected to be solved without positional information. The dataset comes from Choromanski et al. (2021), with experiment setting from Randomized PE(Ruoss et al., 2023).

Level	Task	Baseline					DAPE (Kernel Size 1)			DAPE (Kernel Size 3)		
		RoPE	Relative	ALiBi	Kerple	FIRE	ALiBi	Kerple	FIRE	ALiBi	Kerple	FIRE
R	EVEN PAIRS	99.98	96.60	73.52	57.50	73.86	99.99	99.58	100	99.99	100	100
	MODULAR ARITHMETIC (SIMPLE)	21.35	20.84	20.02	21.79	21.09	23.58	24.47	24.46	21.48	23.90	23.43
	PARITY CHECK†††	50.05	50.09	50.09	50.07	50.97	50.30	50.07	50.04	50.13	52.51	50.11
	CYCLE NAVIGATION†††	27.63	26.95	24.64	29.47	28.41	22.99	34.53	27.54	24.43	24.32	24.34
DCF	STACK MANIPULATION	61.49	64.73	66.42	66.93	69.33	68.18	72.04	70.90	58.90	68.18	60.90
	REVERSE STRING	65.23	65.59	71.09	71.54	65.89	73.37	70.74	76.40	56.61	81.84	70.11
	MODULAR ARITHMETIC	31.25	31.74	30.56	24.79	30.92	31.34	32.37	31.50	29.46	26.13	27.00
	SOLVE EQUATION	21.85	22.93	19.92	21.15	22.06	20.03	22.49	22.42	20.26	23.95	23.62
CS	DUPLICATE STRING	64.97	67.66	65.13	66.72	69.03	70.84	72.95	72.71	52.96	57.03	66.01
	MISSING DUPLICATE	63.37	72.34	74.21	79.06	79.27	83.41	87.57	89.17	59.33	99.65	74.83
	ODDS FIRST	61.00	61.57	59.88	62.59	63.28	63.78	67.08	66.34	57.35	56.87	56.57
	BINARY ADDITION	55.59	56.96	54.72	56.35	55.70	59.71	60.88	56.62	57.49	55.32	57.86
	COMPUTE SQRT	51.88	51.63	50.63	51.11	50.80	51.64	51.33	52.46	52.08	51.76	51.93
	BUCKET SORT†††	98.12	99.31	98.45	99.38	99.57	99.38	98.81	99.37	96.61	99.06	98.56

H DAPE_{1×3} TIME COST

Table 7: The time cost (millisecond) under different testing lengths, with D_{DAPE} as 32 and default batch size 1, with training length 512.

Method	350M Total	Ratio	2.7B Total	Ratio	6.7B Total	Ratio
RoPE (Su et al., 2024b)	210.01	0.8306	472.63	1.0472	635.57	0.8564
T5’s bias (Raffel et al., 2020)	355.16	1.4046	537.62	1.1912	808.85	1.0899
ALiBi (Press et al., 2021)	172.60	0.6826	325.95	0.7222	596.77	0.8041
Kerple (Chi et al., 2022)	189.91	0.7511	370.32	0.8205	661.82	0.8918
FIRE (Li et al., 2023c)	248.13	0.9813	432.63	0.9586	797.68	1.0748
DAPE-Kerple (Zheng et al., 2024)	224.22	0.8868	422.48	0.9361	717.46	0.9667
DAPE _{1×3} -Kerple	252.84	1.0000	451.29	1.0000	742.10	1.0000

1296 I MODEL CONFIGURATION

1297

1298 All experiments are conducted on 8 GPUs. The 125M and 350M model configuration is the follow-
1299 ing.

1300

1301

Table 8: **Model Configurations.**

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318 J DATA-ADAPTIVE RELATED POSITION ENCODING PERFORMANCE
1319 COMPARISON

1320

1321 Table 9: The performance comparison between data-related position encoding, with dataset Books3
1322 and training length 128.

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

Method	128	256	512	1024	2048	4096	8192
Transformer-XL	31.57	28.49	26.07	26.98	27.90	32.76	41.12
CoPE	31.61	28.41	25.79	27.96	33.80	54.08	90.66
DAPE-Kerple (Kernel Size 1x1)	31.49	28.27	24.93	24.31	23.34	24.38	25.01
DAPE _{1x3} -Kerple (Kernel Size 1x3)	31.07	27.81	24.38	23.57	22.40	23.19	23.52

K DAPE_{1×3} VISUALIZATION

The model is trained with DAPE_{1×3}-Kerple on length 512. Compared to DAPE (Zheng et al., 2024), it seems that the DAPE_{1×3} presents a more obvious attention sink (Xiao et al., 2024d).

K.1 VISUALIZATION ON LENGTH 512

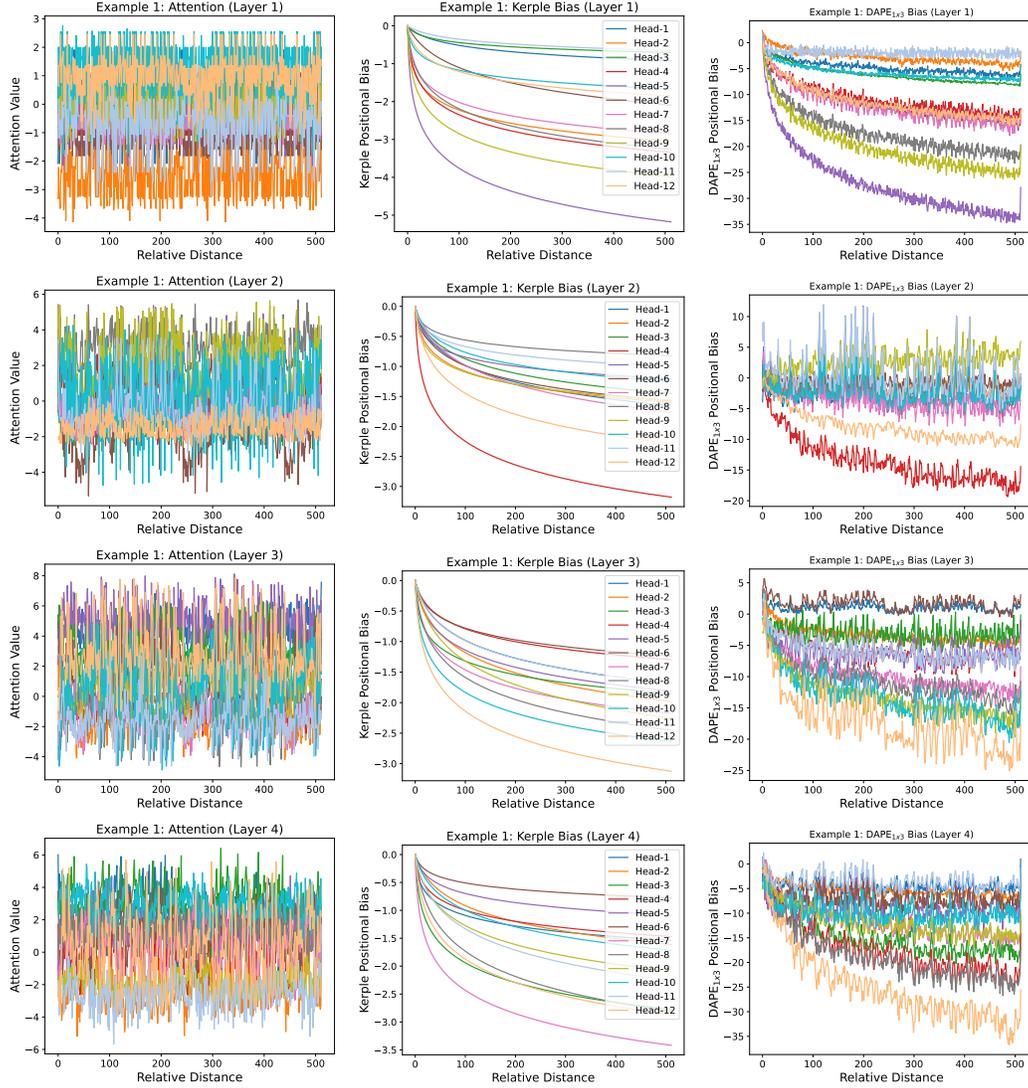


Figure 8: Evaluation Length 512 Example 1: Part 1. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1×3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

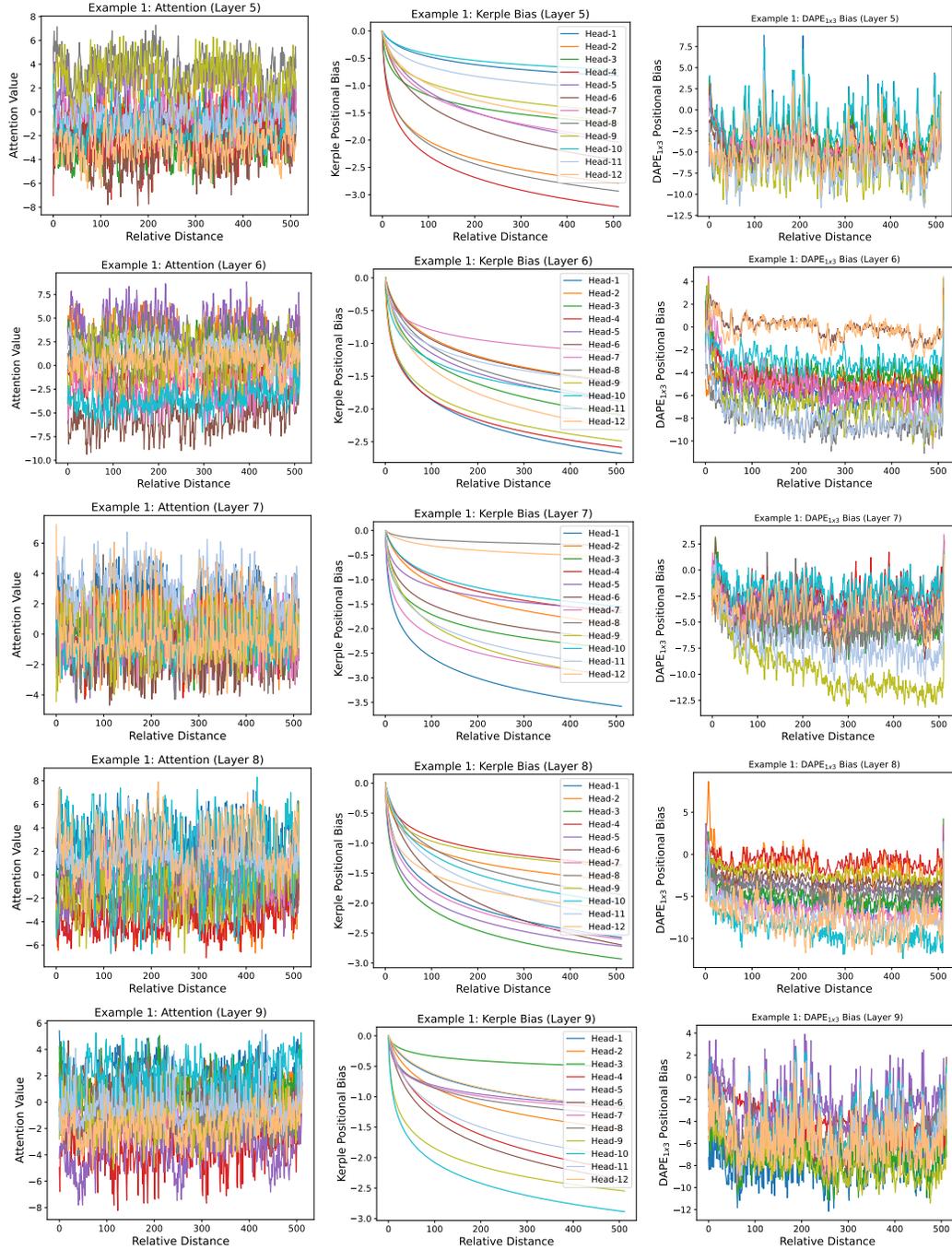


Figure 9: Evaluation Length 512 Example 1: Part 2. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

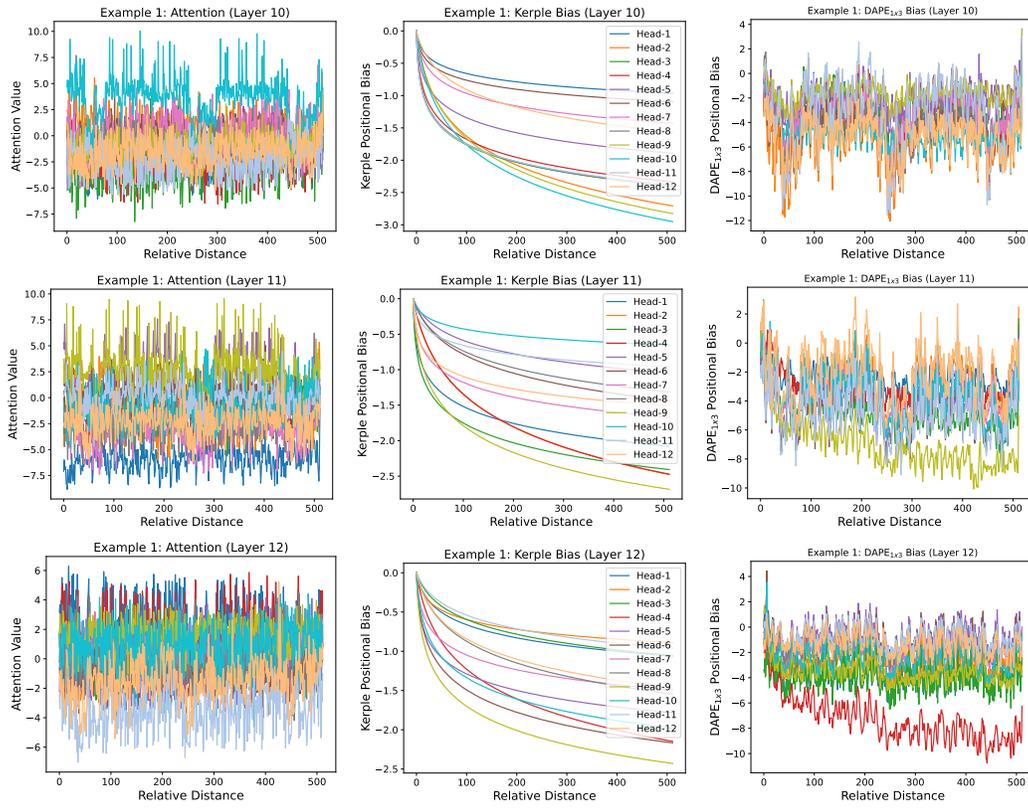


Figure 10: Evaluation Length 512 Example 1: Part 3. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

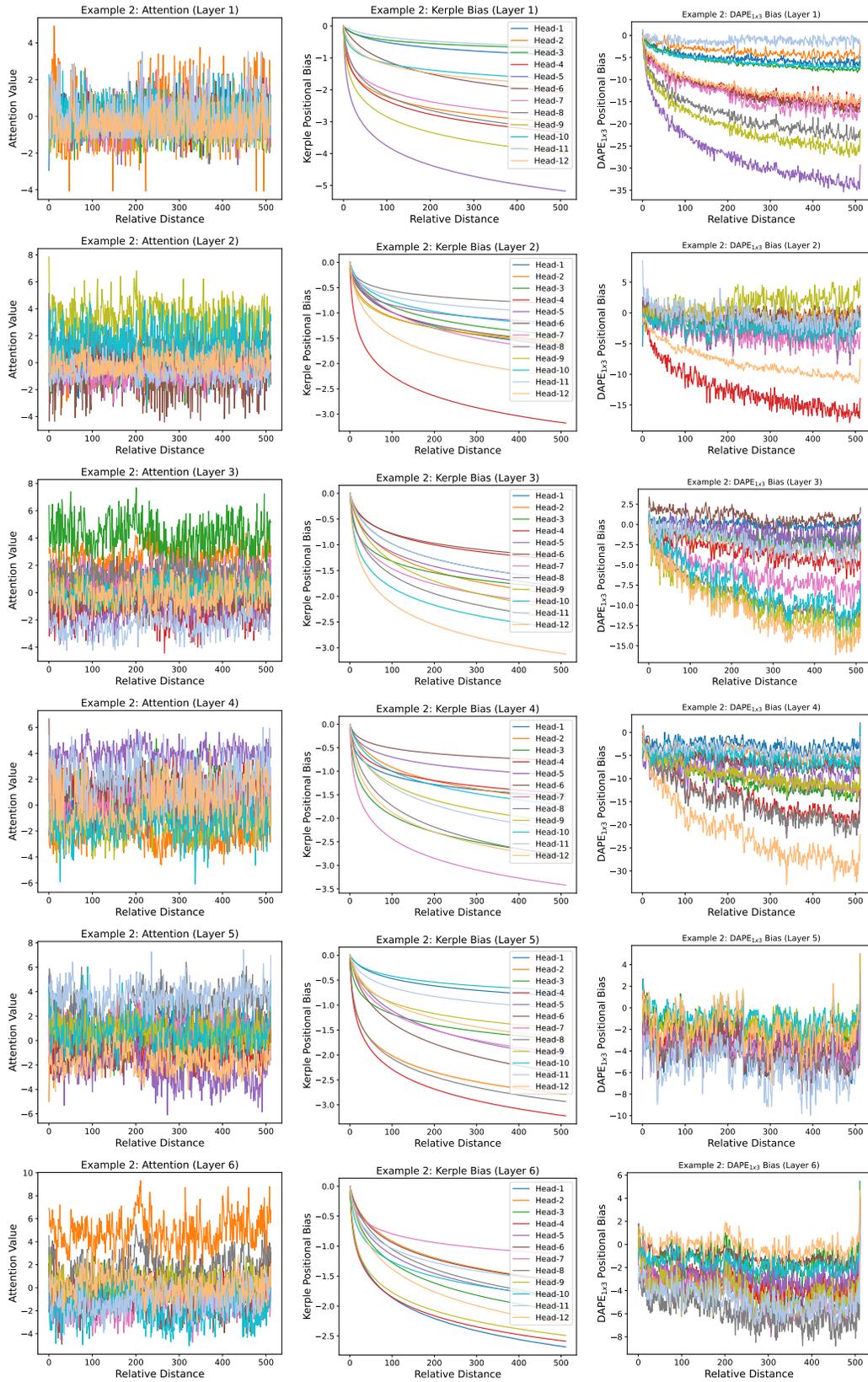
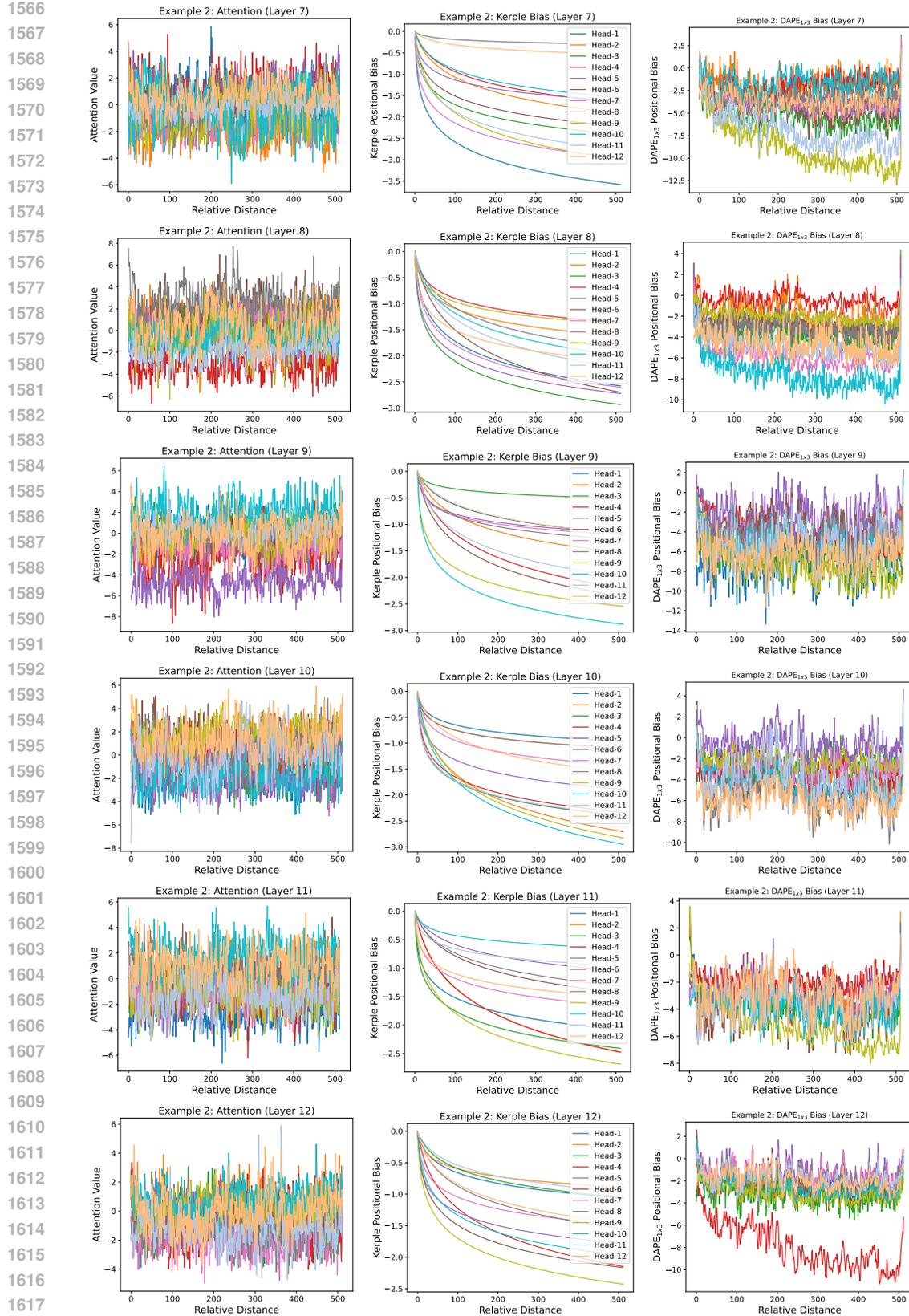


Figure 11: Evaluation Length 512 Example 2: Part 1. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.



1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619

Figure 12: Evaluation Length 512 Example 2: Part 2. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

K.2 VISUALIZATION ON LENGTH 2048

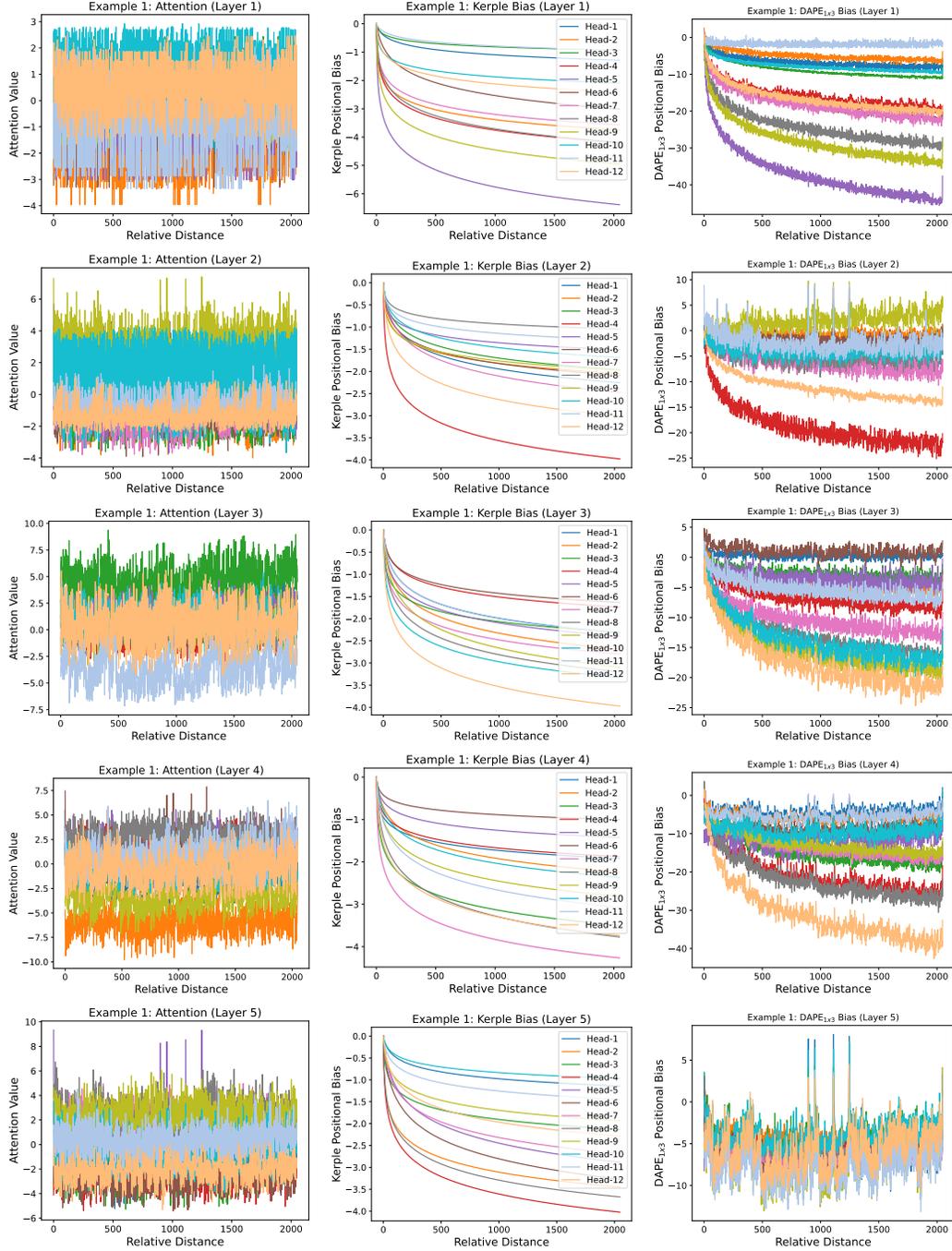


Figure 13: Evaluation Length 2048 Example 1: Part 1. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

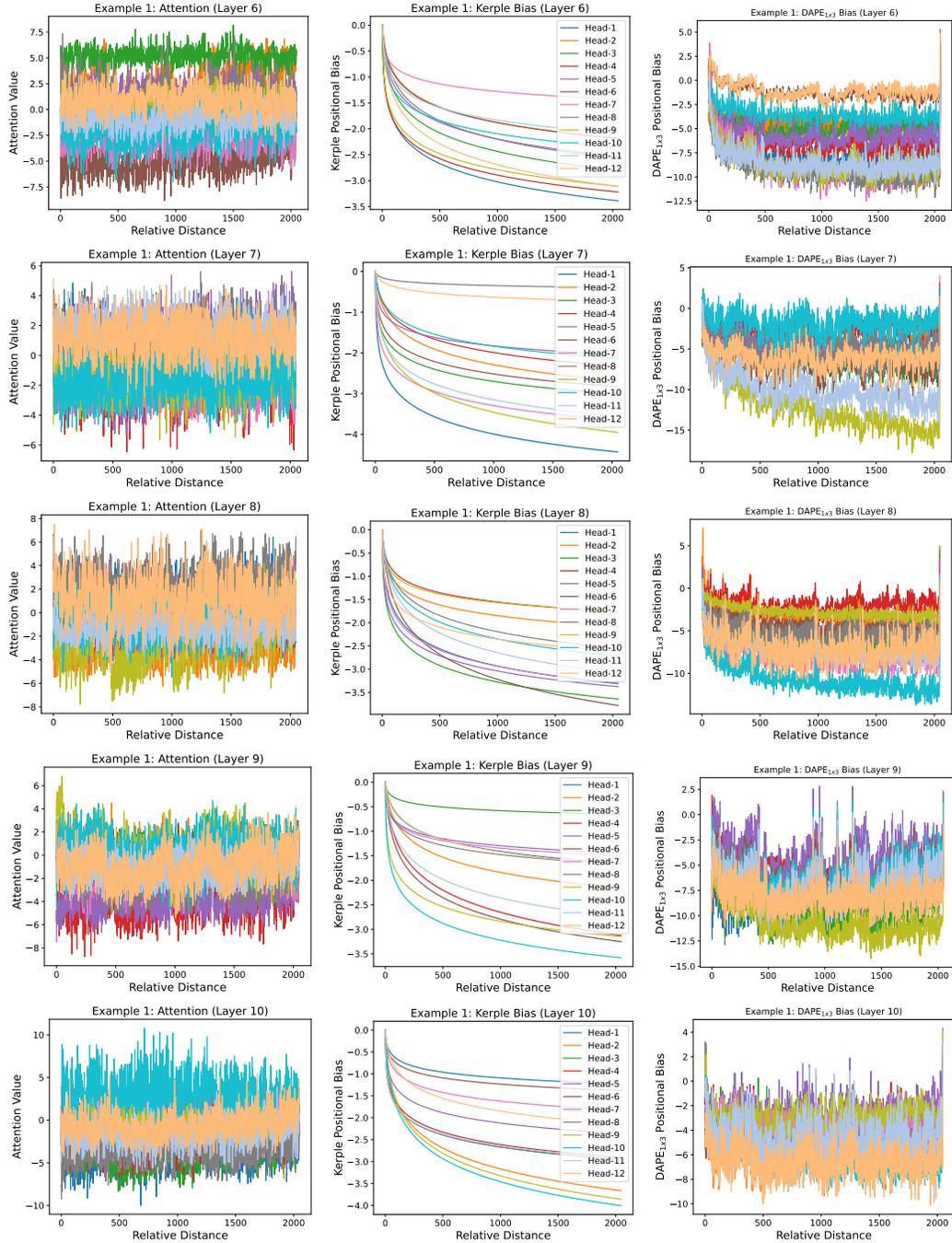


Figure 14: Evaluation Length 2048 Example 1: Part 2. From Left to Right: (1) The Attention is $XW_Q(XW_K)^\top$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^\top, B)$.

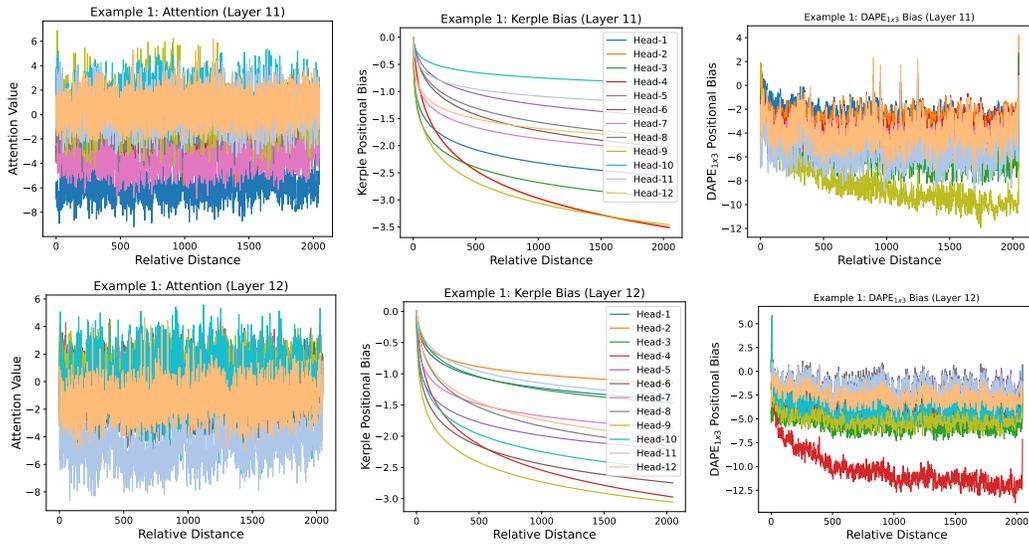


Figure 15: Evaluation Length 2048 Example 1: Part 3. From Left to Right: (1) The Attention is $XW_Q(XW_K)^\top$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^\top, B)$.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

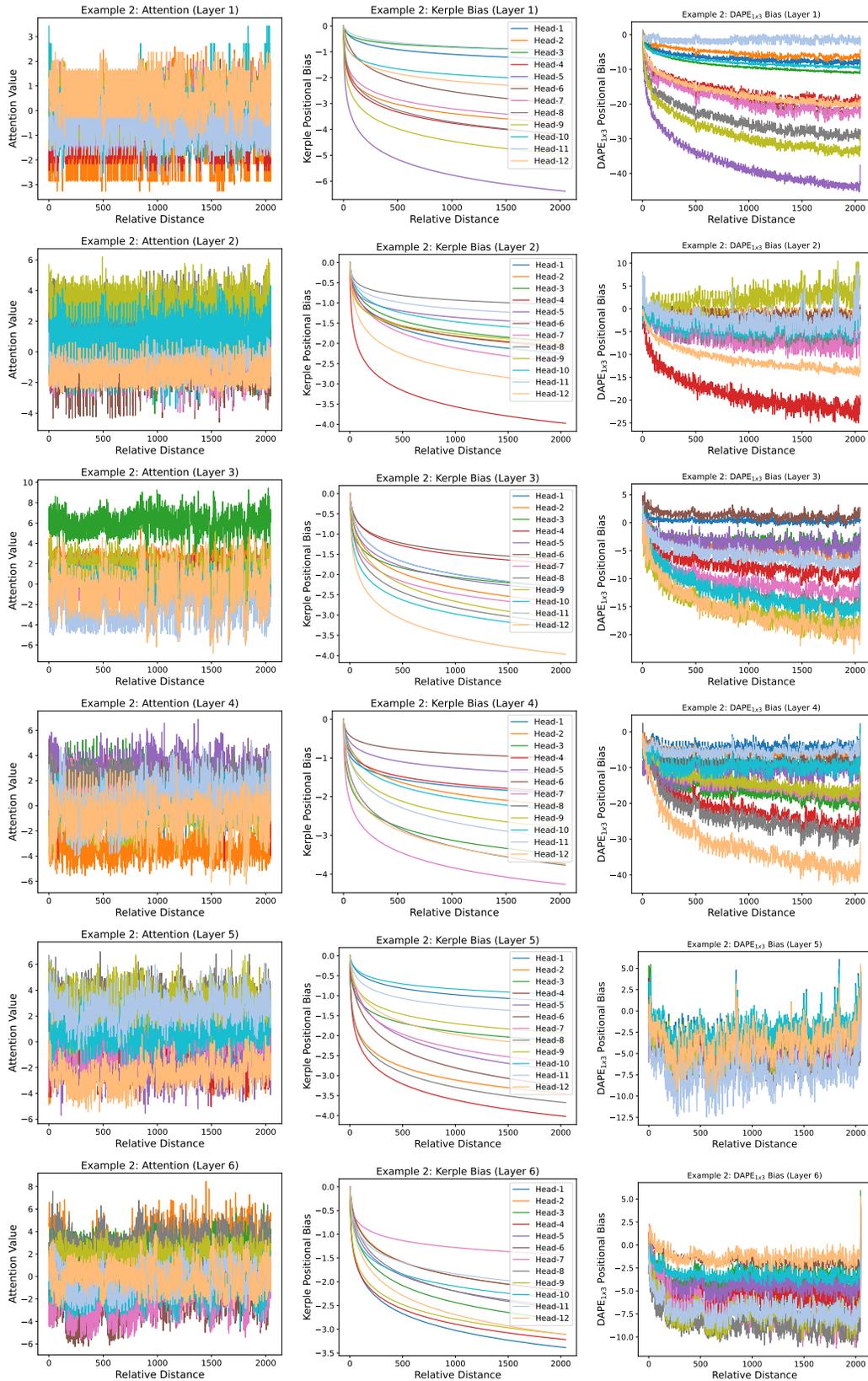


Figure 16: Evaluation Length 2048 Example 2: Part 1. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

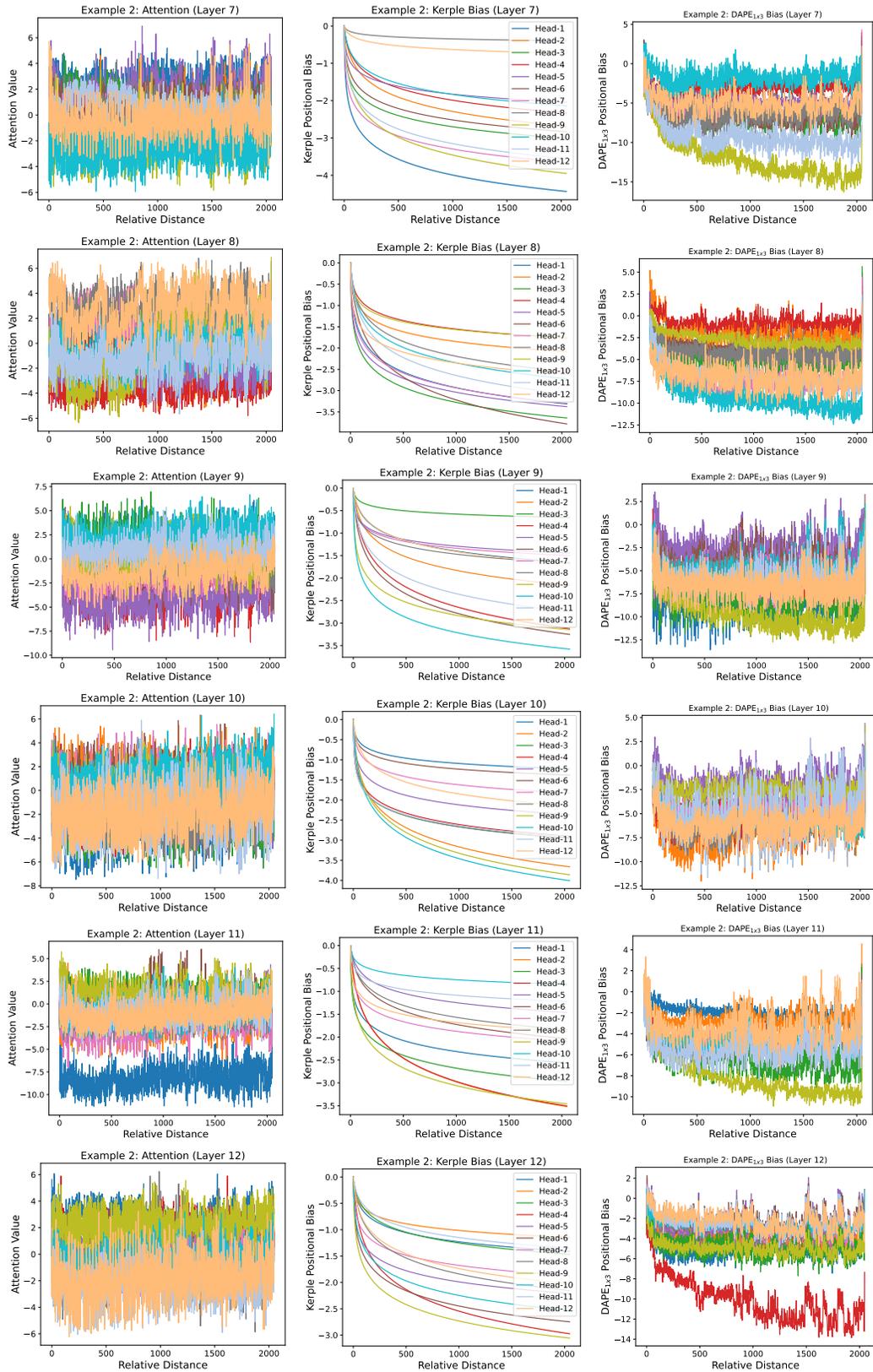


Figure 17: Evaluation Length 2048 Example 2: Part 2. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

K.3 VISUALIZATION ON LENGTH 8192

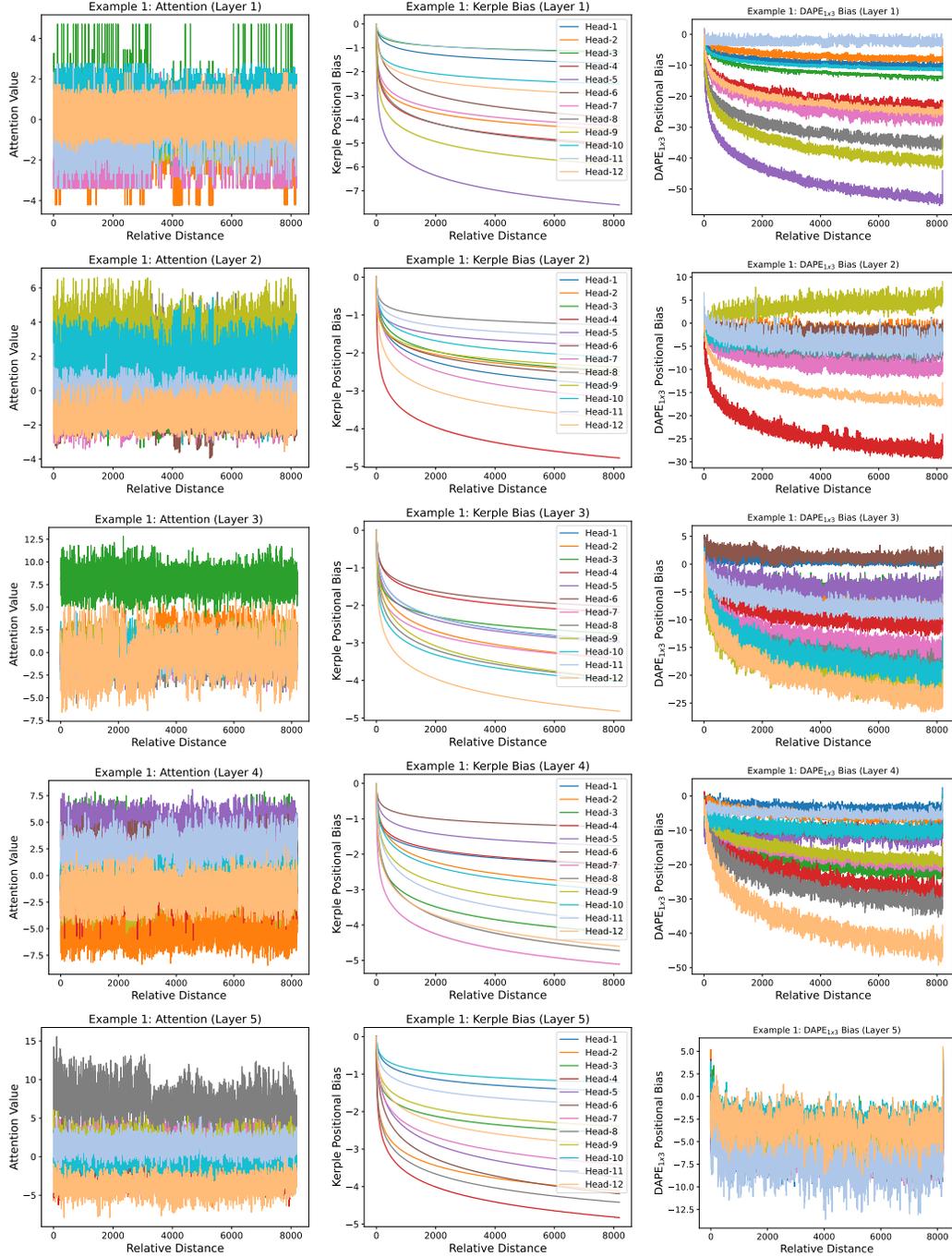


Figure 18: Evaluation Length 8192 Example 1: Part 1. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

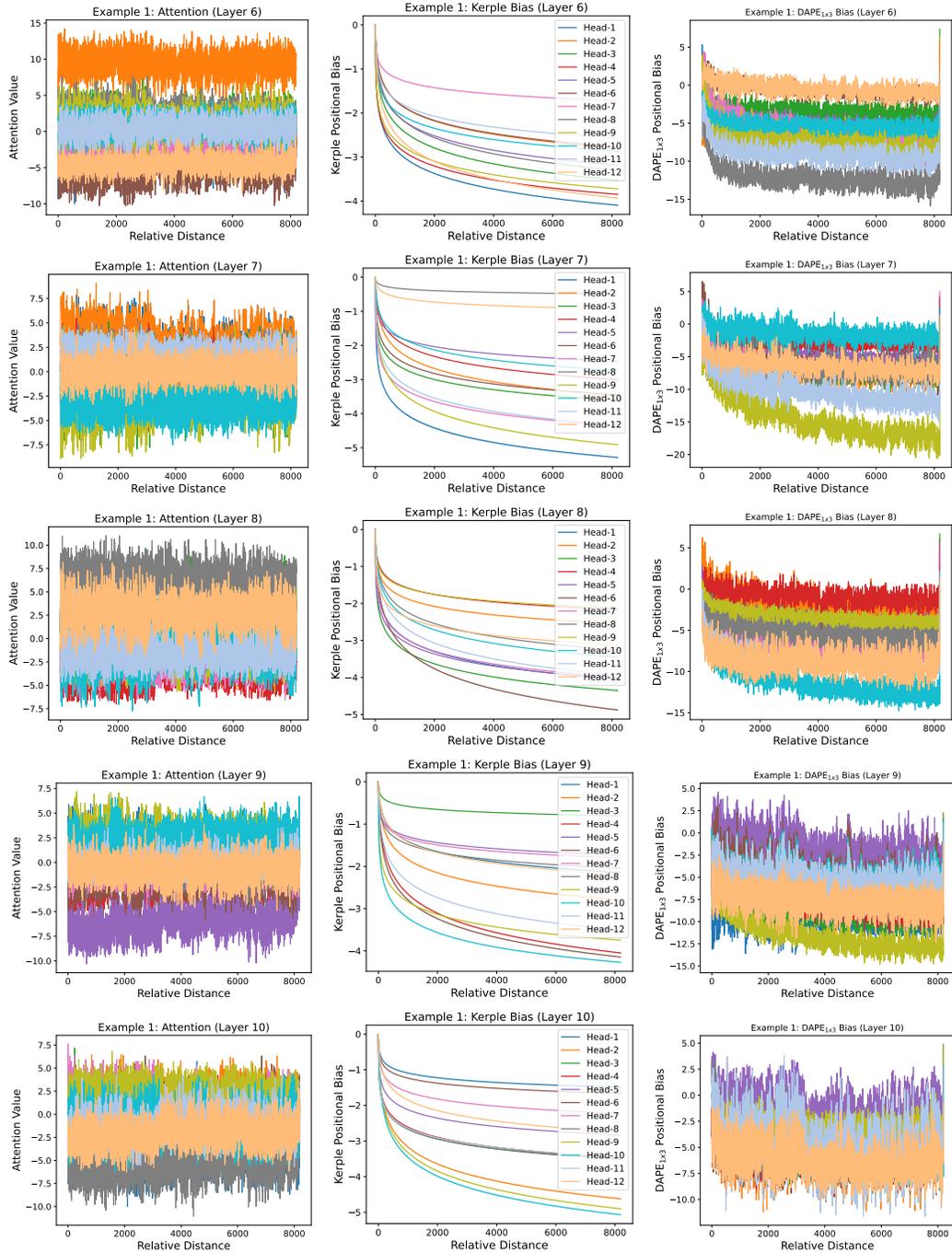


Figure 19: Evaluation Length 8192 Example 1: Part 2. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

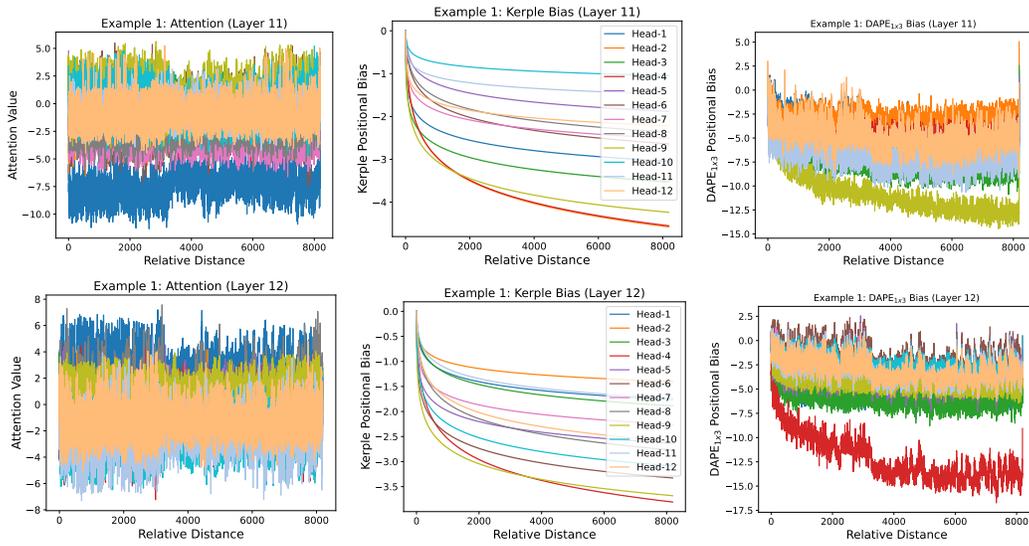


Figure 20: Evaluation Length 8192 Example 1: Part 3. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

2052
 2053
 2054
 2055
 2056
 2057
 2058
 2059
 2060
 2061
 2062
 2063
 2064
 2065
 2066
 2067
 2068
 2069
 2070
 2071
 2072
 2073
 2074
 2075
 2076
 2077
 2078
 2079
 2080
 2081
 2082
 2083
 2084
 2085
 2086
 2087
 2088
 2089
 2090
 2091
 2092
 2093
 2094
 2095
 2096
 2097
 2098
 2099
 2100
 2101
 2102
 2103
 2104
 2105

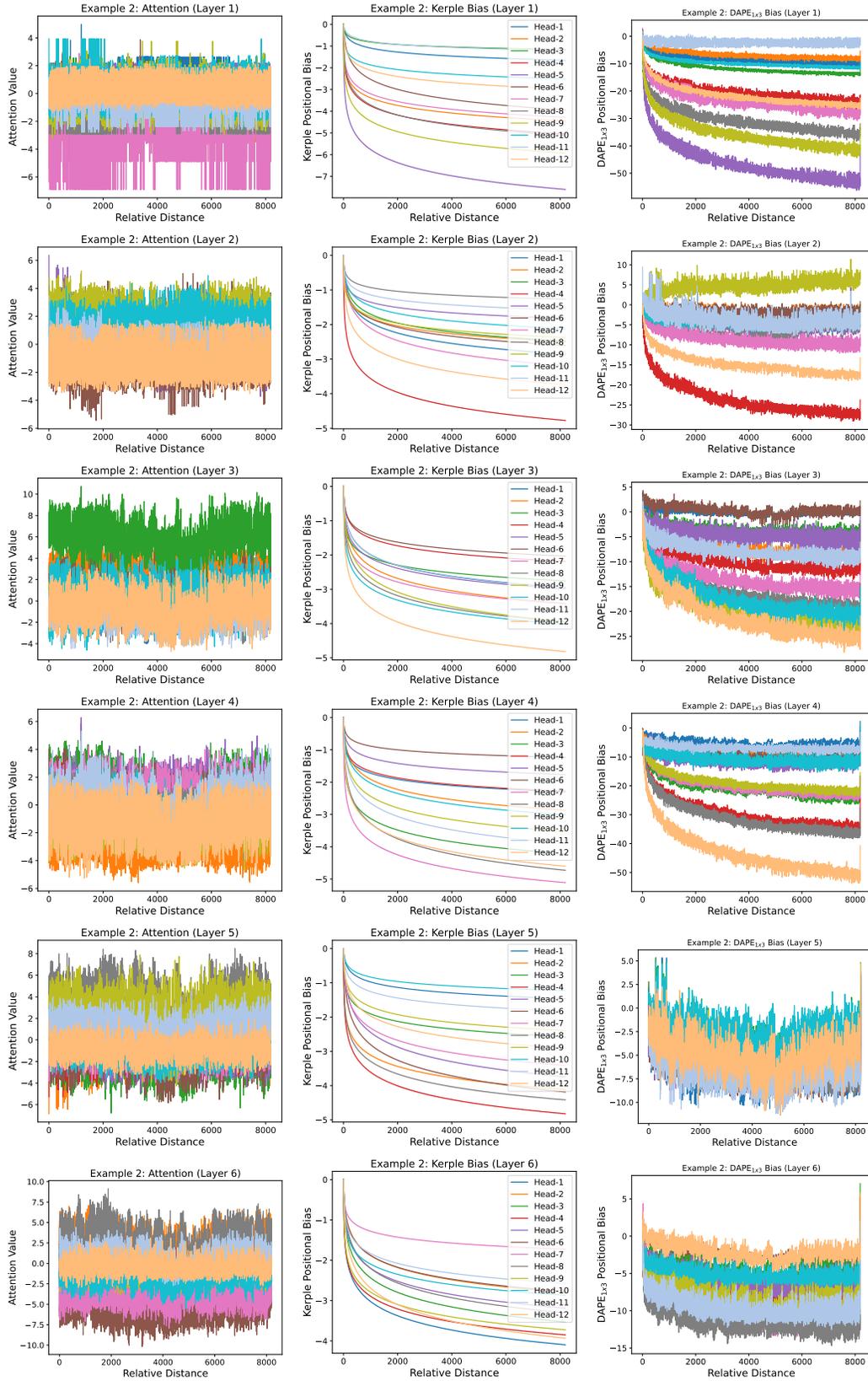


Figure 21: Evaluation Length 8192 Example 2: Part 1. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

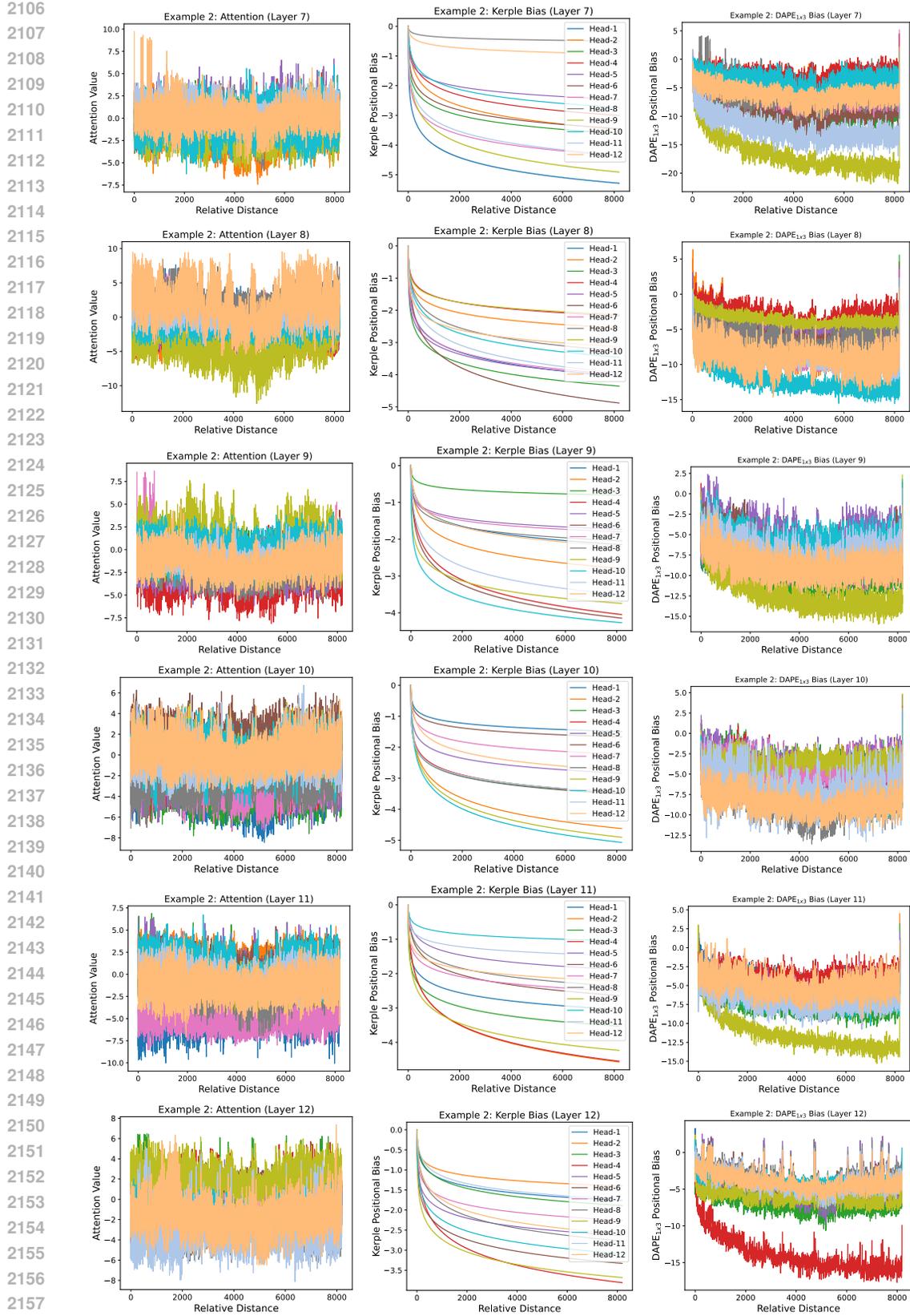


Figure 22: Evaluation Length 8192 Example 2: Part 2. From Left to Right: (1) The Attention is $XW_Q(XW_K)^T$; (2) The Kerple bias is B ; (3) The DAPE_{1x3} (with Kerple) bias is $f(XW_Q(XW_K)^T, B)$.

L IMPLEMENTATION

In this section, we present the implementation of the proposed DAPE_{1×3} module in PyTorch (Paszke et al., 2019).

```

2165 import torch
2166 import torch.nn as nn
2167
2168 class DAPEV2(nn.Module):
2169     def __init__(self, head_number=12, mlp_width=32, kernel_size=3):
2170         """
2171         DAPEV2 attention bias module.
2172
2173         Args:
2174         num_heads: number of attention heads.
2175         mlp_width: Width of MLP.
2176         kernel_size: convolution kernel size.
2177         """
2178         super(DAPEV2, self).__init__()
2179
2180         self.mlp = nn.Sequential(
2181             nn.Conv2d(in_channels=head_number*2, out_channels=
2182                       mlp_width, kernel_size=(1, kernel_size), stride=(1, 1),
2183                       padding=(0, kernel_size // 2), dilation=(1, 1)),
2184             nn.LeakyReLU(),
2185             nn.Conv2d(in_channels=mlp_width, out_channels=
2186                       head_number, kernel_size=(1, kernel_size), stride=(1, 1),
2187                       padding=(0, kernel_size // 2), dilation=(1, 1)))
2188
2189     def forward(self, attention: torch.Tensor, bias: torch.Tensor):
2190         """
2191         Args:
2192         attention: input sequence, which is q^T * k,
2193                   shape [bsz, num_heads, seq_len, seq_len]
2194         bias: bias matrix, which can be generated by ALiBi, Kerple
2195              FIRE or other additive position encodings
2196              shape [1, num_heads, seq_len, seq_len]
2197
2198         Returns:
2199         attention with DAPEV2,
2200         shape [bsz, num_heads, seq_len, seq_len]
2201         """
2202         bias_tile=torch.tile(fire_bias, (x.shape[0],1,1,1) )
2203         attention_bias_concat=torch.cat( (attention, bias_tile), dim=1)
2204         attention_bias_concat=torch.tril(attention_bias_concat)
2205         attention_bias_concat=self.mlp(attention_bias_concat)
2206
2207         return attention+bias+attention_bias_concat

```