

Ultron: A High-Performance Sequence-Processing Model Based on the Composite Mapping Layer

Lingkai Hu^{a,*}, Feng Zhan^{b,*}, Wenkai Huang^{a,**}, Weiming Gan^a, Haoxiang Hu^c

^a*School of Mechanical and Electrical Engineering, Guangzhou University, Guangzhou, 510006, China*

^b*The Hong Kong University of Science and Technology (Guangzhou), Smart Manufacturing Thrust, Nansha, Guangzhou, 511400, China*

^c*College of Information Science and Technology, Beijing University of Chemical Technology, Beijing, 102202, China*

ARTICLE INFO

Keywords:

Language modeling

Composite mapping

Sequence processing

Long-term dependency

ABSTRACT

A crucial task in deep-learning research, sequence processing involves extracting semantic or structural relationships between distant elements in a sequence. However, existing sequence-processing methods have several limitations, such as high levels of complexity, difficulty in modeling long-term dependencies, and insufficient generalization ability. To address these issues, Ultron, a high-performance sequence-processing model based on the composite mapping layer (CM layer), is proposed. The CM layer is a novel neural layer that uses a composite function of multiple information-mapping functions to capture long-term information dependencies. Compared with the self-attention mechanism, the CM layer has a lower complexity of $O(L \log L)$. Ultron adopts the CM layer as the core component of the encoder and decoder structures, achieving the same parallel-computing capability as the Transformer model and improving its fitting ability and robustness. A high-quality benchmark dataset for Chinese language modeling, zhwiki520, is also constructed based on the complete text of Chinese Wikipedia as of May 20, 2023. Ultron is evaluated on the long-range arena, wikitext-103, zhwiki520, and enwik8. The experimental results show that Ultron surpasses Transformer and its variants in all the tasks and meets the current criteria for state-of-the-art performance on the enwik8 dataset. This paper contributes a novel and effective method to the sequence-processing field and provides a valuable evaluation metric for Chinese language modeling tasks.

1. Introduction

Sequence processing is a fundamental problem in various fields, such as natural language processing (NLP) [1–3], speech recognition [4–6], machine translation [7–11], and bioinformatics [12–15]. It involves tasks that require the analysis, understanding, generation, or transformation of sequences of different modalities (e.g., text, audio, or image) [4]. The main challenge of sequence processing is to effectively capture the long-term dependencies in sequences, which are the semantic or structural relations between distant elements [16–18].

Traditional sequence-processing models rely mainly on either recurrent neural networks (RNNs) [19, 20] or convolutional neural networks (CNNs) [21, 22]. RNNs can handle arbitrary-length sequences by recursively updating their hidden states, but they suffer from gradient vanishing or exploding, parallelization difficulty, and ineffective long-term dependency modeling [23, 24]. CNNs can extract local features through multiple convolutional operations and expand their receptive field through pooling or skip connections, thus achieving parallelization and long-term dependency modeling. However, they also have limitations, including a fixed receptive field size and sequential-information loss [25, 26].

In recent years, Transformers [27] have gained widespread attention and application. They are sequence-processing models that use self-attention mechanisms, which can model global dependencies by computing the relevance between any two positions in a sequence and have efficient parallel training and inference capabilities. Transformer-based models have improved performance significantly in various NLP tasks and have inspired many variants and extensions, such as BERT [28], GPT [29], and XLNet [30].

However, Transformer-based models face several challenges and issues [31]. First, the self-attention mechanism has a quadratic time complexity and space complexity with respect to the input-sequence length, which requires a large number of parameters and computational resources. This hinders their performance in long-sequence tasks [16].

*Co-first author.

**Corresponding author.

E-mail address: smallkat@gzhu.edu.cn (Wenkai Huang)

Second, they tend to disregard the relative relationships and dependencies among different positions in the sequence, which may affect the model’s sensitivity to long-term dependencies and structural information [32]. Third, they may have limited generalizability or transferability when dealing with texts from diverse languages or domains [33, 34]. Therefore, designing a more efficient, robust, and scalable sequence-processing model is a valuable research endeavor.

This study proposes a novel neural layer for sequence processing called the composite mapping layer (CM layer). By applying a composite function of multiple information-mapping functions, the CM layer can capture long-term information dependencies. Based on the CM layer, a high-performance neural network called Ultron is designed. Ultron has lower complexity than Transformer and outperforms all the Transformer variants in the experiments on fitting ability (Fig .1). Most existing research on language modeling has focused on English datasets, such as Penn Treebank [35], wikitext-103 [36], and enwik8 [37]. However, Chinese, a widely used language, lacks a general benchmark dataset for this task. To address this gap, zhwiki520¹, a high-quality benchmark dataset for Chinese language modeling, is constructed based on the full text of Chinese Wikipedia as of May 20, 2023 [38]. It is preprocessed using Su’s script [39] and released for public use.

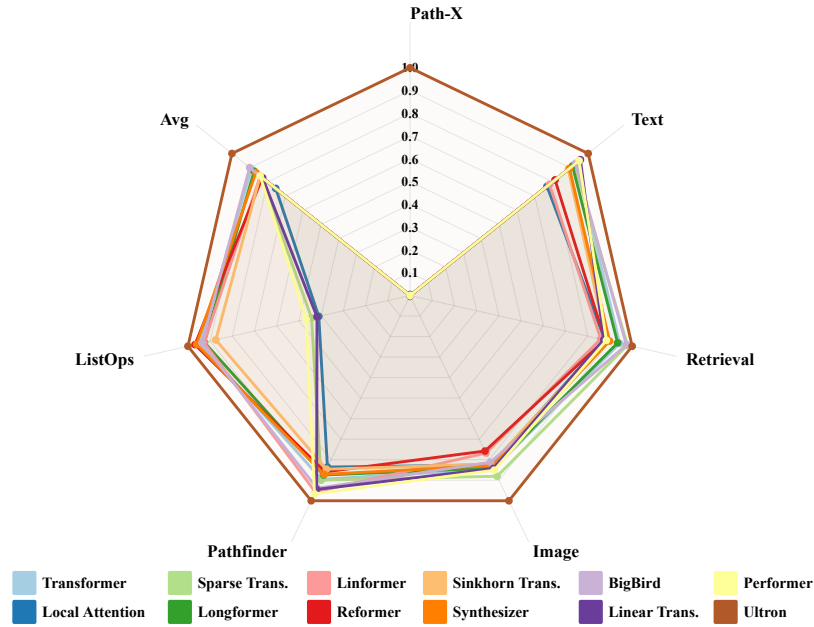


Fig. 1. Accuracy comparison of Ultron and Transformer variants in the LRA [40] dataset.

This paper makes the following contributions:

1. A novel neural layer for sequence processing called the CM layer is proposed, which models long-term information dependencies through the composition of multiple information-mapping functions. It has a lower complexity of $O(L \log L)$ than the self-attention mechanism.
2. A novel model called Ultron is proposed, which incorporates CM layers into both the encoder and decoder structures. Like the Transformer model, the Ultron model can perform parallel computation, and it has better fitting ability and robustness than the Transformer model.
3. The Ultron model is evaluated on four datasets: long-range arena (LRA), wikitext-103, zhwiki520, and enwik8. The experimental results demonstrate that the Ultron model outperforms Transformer and its variants in all the tasks and meets the current criteria for state-of-the-art performance on the enwik8 dataset.

The paper is structured as follows. Section 2 introduces the research background of the sequence-processing model. Section 3 examines the principle and implementation of the CM layer as well as the structure and details of the Ultron model. Section 4 outlines the experimental settings and results. Section 5 concludes the paper by summarizing its contributions and suggesting directions for future work.

¹Our dataset zhwiki520 is published on: <https://kaggle.com/datasets/lingkaihu/zhwiki520>

2. Preliminaries

Sequence processing is a crucial machine learning problem with various applications in NLP, speech recognition, machine translation, bioinformatics, and other domains. However, sequence processing also poses some challenges, the most important of which concerns how to model the long-term dependencies in the sequence effectively.

The current mainstream sequence-processing models fall into three categories: RNN-based models, CNN-based models, and self-attention-based models [41, 42]. Each of these types of models has a unique set of strengths and weaknesses. This section will briefly introduce their principles and performance and discuss some of their problems and limitations in long-sequence tasks.

2.1. Models based on RNNs

RNNs have a neural-network structure that can handle sequences of arbitrary length by recurrently updating the hidden state, which captures the sequential and contextual information in the sequence [19]. RNNs have achieved good results in many sequence-processing tasks, but they also have some drawbacks, such as those identified by [23, 24]:

1. RNNs suffer from the problem of gradient vanishing or exploding, which means that the gradient decays or grows exponentially as the time step increases in the backpropagation process, leading to difficult or unstable training.
2. RNNs are difficult to parallelize, because each time step depends on the calculation result of the previous time step, which limits the training speed and scale of RNN.
3. RNNs cannot effectively model long-term dependencies, because the hidden state gradually loses or confuses past information as the time step increases.

Some improved and extended models, such as long short-term memory [43] and gated recurrent units [44], have been proposed to overcome the shortcomings of RNNs. These models introduce gate structures to control the information flow and forgetting degree, thereby alleviating the gradient-vanishing and long-term dependency problems. However, when the sequence length is extremely long, these improvements are still weak. Further, some studies have introduced cross-time skip connections to alleviate memory decay [45, 46], but this approach breaks the continuity of the sequence, reduces the recognition accuracy of the neural network, and increases the complexity of the neural network. Other studies have incorporated attention into RNN computation to establish long-term dependencies [47, 48], but the resulting models have the complexity and robustness problems that characterize other attention-based techniques. In general, these improved models are still limited by the RNN structure itself, and they still have significant drawbacks when the input sequence is excessively long. Moreover, they require more parameters and computational resources [49].

2.2. Models based on CNNs

CNNs are neural networks that can extract local features and model spatial- or temporal-translation invariance by performing multiple convolutional operations. They also expand the receptive field by using pooling or skip connections [21]. CNNs have been widely applied in image processing, computer vision, and other fields. They can also process the sequential data used for tasks such as text classification and semantic analysis. In contrast to RNNs, CNNs can achieve parallelization, because each convolution kernel depends only on local input rather than on the entire sequence. This improves the neural network's training speed. Moreover, CNNs can capture features of different scales by using convolution kernels of different sizes, thus achieving multiscale modeling [22]. However, CNNs have some drawbacks, the most important of which are as follows [25, 50]:

1. The receptive-field size of CNNs is fixed, meaning that each output depends only on a fixed length of input. This may cause CNNs to fail to capture long-term dependencies that lie beyond the range of the receptive field.
2. CNNs cannot capture the sequential information in a sequence, meaning that the correspondence between each output and input is fixed. This may cause CNNs to fail to handle position changes or structure changes in the sequence.
3. CNNs require a large number of convolution kernels to cover features of different positions and different scales. This increases the parameter and computational complexity of CNNs.

To overcome the shortcomings of CNNs, some improved and extended models have been proposed, such as dynamic convolutional neural networks [51], dilated convolution [50], and deformable convolutional networks [52].

Yang et al. proposed the Shuffle network, which uses a depthwise separable CNN to form a new group convolutional fusion unit. It also uses shuffling and an intraskip-connection mechanism to expand the receptive field of CNNs [53]. Xue et al. proposed a novel fault-diagnosis method based on a local binary temporal CNN, which uses a new temporal module with dilated causal convolution to establish long-term dependencies. It also adopts a local binary convolution layer to reduce computational parameters and improve robustness [54]. These models introduce new types of convolution kernels to adapt to features of different positions and different scales, thereby enhancing the flexibility and adaptability of CNNs. However, these models cannot solve the fundamental problem of CNNs, that is, their inability to obtain a global receptive field [55].

2.3. Models based on self-attention

The self-attention mechanism is a technique that computes and aggregates the relevance between any two positions in a sequence using matrix-multiplication and softmax -normalization operations [27]. It has been widely applied in NLP, machine translation, and other fields and has been used to construct sequence-processing models based purely on self-attention mechanisms, such as Transformer. The self-attention mechanism has the following advantages over RNNs and CNNs [56, 57]:

1. The self-attention mechanism can model global dependency relationships, meaning that each output can depend on the entire sequence of inputs rather than on local or fixed-length inputs. This improves the self-attention mechanism's sensitivity and expressiveness for long-term dependency relationships.
2. The self-attention mechanism can achieve parallelization, because each output depends only on the input matrix and weight matrix rather than on the previous output or previous time step. This improves the self-attention mechanism's training speed and scale.
3. The self-attention mechanism can capture position information and structure information by using position encoding or relative-position encoding to assign different weights or biases to different positions. This enhances the understanding of the self-attention mechanism for sequential-order information and structure information.

However, the self-attention mechanism has some drawbacks, the most important of which are as follows [16, 34, 40, 58]:

1. Both the time complexity and space complexity of the self-attention mechanism are quadratic functions of the input-sequence length. Consequently, a large number of parameters and extensive computational resources are required for the self-attention mechanism to achieve efficient sequence modeling.
2. The self-attention mechanism tends to ignore the relative relationship and dependency between different positions in the sequence, meaning that the relevance between each output and input depends only on their content rather than on their position. As a result, the self-attention mechanism may be insensitive to long-term dependency and structure information in the sequence.
3. The self-attention mechanism may have insufficient generalizability or transferability when dealing with different languages or domains of text, meaning that the relevance between each output and input depends only on their semantics rather than on their syntax or pragmatics. As a result, the self-attention mechanism may not fully capture the language features and domain knowledge in the sequence.

To overcome the self-attention mechanism's drawbacks, some studies have proposed improved and extended models, such as Linformer, Performer, and Longformer [59–69]. Essentially, these models compress the information in the sequence into lower-dimensional data and either perform information mapping based on the compressed data [59, 70, 71] or sparsify or group queries and keys [72–74], thereby reducing the complexity. However, this compression process may lead to information loss or limit the attention mechanism to local information, which prevents the neural network from obtaining complete, global information. Therefore, in practical performance comparisons, these models often have a weaker fitting ability than the original Transformer [40].

3. Methodology

3.1. CM layer

The fitting process of artificial neural networks involves adjusting the network parameters to approximate or simulate an unknown-function relationship. For example, in a unidirectional-sequence-prediction model, the hidden

layer function relationship is defined as $Y = F(X)$, where X is an input sequence of length L . The output of Y at the n th time step can be expressed as

$$Y_n = F_n (X_1, \dots, X_n). \quad (1)$$

This equation can be rewritten as

$$H_n = \sum_1^n F_{(i,n)} (X_i) \quad (2)$$

$$Y_n = G_n (H_n), \quad (3)$$

where (2) represents the fusion of the context features of the entire sequence, $F_{(i,n)}$ is the information-mapping function from the i th time step to the n th time step, and (3) is the feature-mapping function of the fused information. The complexity of this calculation is

$$O(F(X)) = O\left(\sum_1^L i\right) = O(L^2/2). \quad (4)$$

To model the information relationship between any two-time points, both the number and complexity of parameters in this calculation are quadratic functions of the sequence length, which consumes a large amount of computation. One method of simplifying the computation is to assume that the information-mapping relationship between any two-time points with the same distance is the same, such as:

$$F_{(i,n)} = F_{(i+T,n+T)}, \quad (5)$$

and assume that the long-term-mapping relationship can be expressed as the iteration of the unit-mapping function of the distance number, such as:

$$F_{(i,n)} = F_1^{n-i}. \quad (6)$$

In (6), we use F_1 to represent any unit-mapping function, such as $F_{(i,i+1)}$. By iterating multiple unit-information-mapping functions F_1 , we simulate the long-term information-mapping function, such as:

$$Y_{n+1} = F_1 (Y_n, X_{n+1}). \quad (7)$$

Eq. (5) is a typical Markov process, in which the output of each step is a function of the output and input of the previous step. RNNs can be considered a special case of it. In this computation method, due to parameter sharing, the number of parameters is greatly reduced; due to the presence of only L iterations, the complexity is reduced to $O(L)$. This method improves the model's computational efficiency to some extent, but the iterative structure prevents the model operation from being parallelized. When the sequence length is long, the information tends to decay in the iteration, making it difficult to establish long-term dependency relationships.

To address the limitations of existing sequence-processing models, this study proposes a novel algorithm improvement called the CM layer. The CM layer consists of several information-mapping functions with distances of 2 to the power of integers, such as: $F_1, F_2, F_4, F_8, \dots$.

These functions can be composed to form any information-mapping function with an arbitrary distance, such as

$$F_{(3,10)} = F_7 = F_1 \circ F_2 \circ F_4 \quad (8)$$

$$F_{(1,15)} = F_{14} = F_2 \circ F_4 \circ F_8 \quad (9)$$

$$F_{(7,18)} = F_{11} = F_1 \circ F_2 \circ F_8 \quad (10)$$

$$F_{(6,21)} = F_{15} = F_1 \circ F_2 \circ F_4 \circ F_8. \quad (11)$$

This technique enables long-term information mapping through shorter computation paths. It can be efficiently implemented with parallel computing, and its specific implementation technique is expressed as Algorithm 1.

Algorithm 1 Information Mapping in Decoder

Input: $X \in \mathbb{R}^{L \times d}$
Output: $Y \in \mathbb{R}^{L \times d}$
for $i = 0$ to $\lceil \log_2 L \rceil$ **do**
 $X[2^i :] += X[: -2^i] W_{2^i}$
end for
 $Y = X W^Y$

In this paper, the mapping function is defined as a linear layer with parameters $W_{2^i} \in \mathbb{R}^{d \times d}$ for each information-mapping function and $W^Y \in \mathbb{R}^{d \times d}$ for the feature-mapping function in (3). Thus, all the information-mapping functions can be simplified as a single linear transformation by multiplying several weight matrices. Therefore, the length of the information-propagation path in this structure is 1.

To enable neurons to access complete contextual information in the classification model, a bidirectional CM layer is designed, as shown in Algorithm 2.

Algorithm 2 Information Mapping in Encoder

Input: $X \in \mathbb{R}^{L \times d}$
Output: $Y \in \mathbb{R}^{L \times d}$
 $\vec{X} = X[:, : d/2]$
 $\leftarrow{X} = X[:, d/2 :]$
for $i = 0$ to $\lceil \log_2 L \rceil$ **do**
 $\vec{X}[2^i :] += \vec{X}[: -2^i] \vec{W}_{2^i}$
 $\leftarrow{X}[: -2^i] += \leftarrow{X}[2^i :] \leftarrow{W}_{2^i}$
end for
 $X = \text{Concat}(\vec{X}, \leftarrow{X})$
 $Y = X W^Y$

This algorithm allows neurons to obtain comprehensive contextual information.

3.2. Feedforward networks

To introduce nonlinearity, a nonlinear feedforward network is defined between each CM layer. This network comprises two linear layers with a Leaky ReLU [75] activation function. The Leaky ReLU is defined as

$$\text{Leaky ReLU}(X) = \text{Max}(0.01X, X). \quad (12)$$

The feedforward network is defined as

$$\text{FFN}(X) = \text{Leaky ReLU}(XW_1 + b_1)W_2 + b_2. \quad (13)$$

Leaky ReLU is selected as the activation function due to its nonzero gradient in the negative region, which helps prevent neurons from becoming inactive. The input and output embedding dimensions of the feedforward network are denoted as d , while the intermediate layer embedding dimension is d_{ff} .

3.3. Model architecture

Conventional sequence-processing models typically adopt an encoder–decoder architecture for various tasks. Ultron follows this design and adds learnable position encoding as an input layer, enabling the model to capture relative-position relationships. Fig. 2 shows its encoder and decoder structures, where Fig. 2A is the encoder and Fig. 2B is the decoder:

Encoder: The encoder consists of N layers, each containing two sublayers. The first sublayer is a bidirectional CM layer that models both forward and backward contexts. The second sublayer is a feedforward network that applies a nonlinear transformation to the output of the CM layer. Each sublayer has layer normalization and residual connection for gradient stability.

Decoder: The decoder also consists of N layers, each containing two sublayers. In the odd-numbered layers, the first sublayer is a unidirectional CM layer that models only the forward context. The second sublayer is a feedforward network that applies a nonlinear transformation to the output of the CM layer. In the even-numbered layers, the first sublayer is a unidirectional self-attention layer that computes the relevance between each output token and all the previous output tokens. The second sublayer is a feedforward network that applies a nonlinear transformation to the output of the self-attention layer.

The decoder interleaves the CM layer and the self-attention layer because better accuracy is achieved by this strategy than using only the CM layer or only the self-attention layer in our experiments. It is conjectured that the CM layer excels at modeling relative-position dependency whereas the self-attention layer excels at computing long-term relevance. This mixed strategy leverages the strengths of both structures.

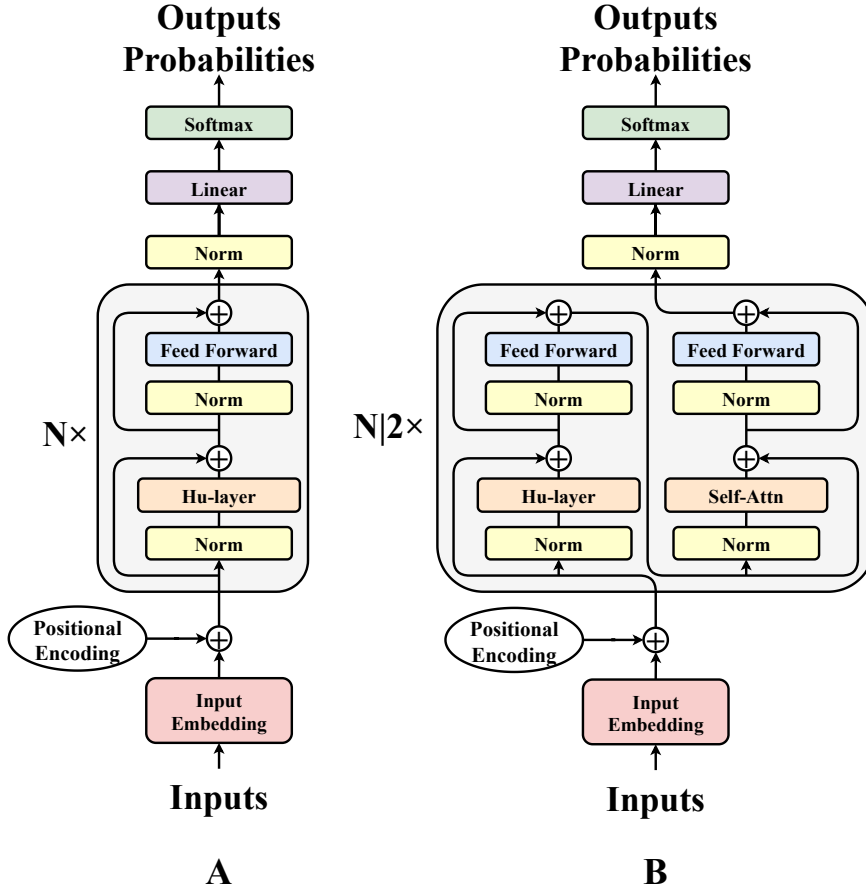


Fig. 2. The overall architecture of the Ultron.

3.4. Complexity analysis

The complexity of the CM layer arises from its iterative computation, which requires approximately $\log_2 L$ iterations. The complexity can be expressed as

$$O\left(\int_0^{\lceil \log_2 L \rceil} d^2 (L - 2^x) dx\right) = O\left(\frac{d^2}{\ln 2} (L \ln L - L + 1)\right) \quad (14)$$

$$\approx O(L d^2 \log_2 L). \quad (15)$$

Hence, the CM layer exhibits a complexity of $L \log_2 L$, which is lower than that of the self-attention mechanism. The CM layer demonstrates a more pronounced advantage when L is large. Table 1 compares the performance metrics of common neural network layers and highlights the CM layer's relatively high scores in these metrics.

Table 1

Comparison of complexity, maximum path length, and perceptual field of various neural network layers [16].

Layer Type	Complexity per Layer	Maximum Path Length	Receptive Field
Self-Attention	$O(L^2 D)$	$O(1)$	$O(L)$
CM Layer	$O(LD^2 \log_2 L)$	$O(1)$	$O(L)$
Fully Connected	$O(L^2 D^2)$	$O(1)$	$O(L)$
Convolutional	$O(KLD^2)$	$O(\log_K L)$	$O(K)$
Recurrent	$O(LD^2)$	$O(1)$	$O(L)$

“L” is the sequence length, “D” is the number of hidden state dimensions, “K” is the convolution kernel size, and “T” is the sampling period of the cuneate layer.

4. Experiments

In the previous section, a detailed account of Ultron’s structure, which is based on the CM layer, was provided, and the feasibility and effectiveness of the CM layer were theoretically analyzed. In this section, the performance of our model is evaluated on several benchmark experiments, including LRA, wikitext-103, enwik8, and zhwiki520, a benchmark for Chinese language modeling that was constructed. The neural-network models were implemented using PyTorch, and Adam [76] was employed as the optimizer. The training and testing of the models were conducted on a single NVIDIA GeForce RTX 4090 GPU.

4.1. Long-range modeling on LRA

This study evaluated the encoder of our model on the LRA dataset, a systematic method for assessing sequence-processing models. The LRA dataset consists of six tasks and datasets that test various aspects of sequence-processing models, such as generalization ability, computational efficiency, memory usage, and reasoning skills, on long sequences. The LRA dataset covers diverse data types and modalities, including text, natural images, synthetic images, and mathematical expressions. The sequence lengths range from 1 to 16 K, and the data require similarity, structural, or visual-spatial reasoning.

Table 2

Hyperparameters of neural networks in the LRA experiment.

Task	N	d_{model}	d_{ff}	Batch size	Iterations	Epochs
ListOps	4	512	1024	32	5000	\
Text	4	256	1024	32	20000	\
Retrieval	4	128	512	32	5000	\
Image	1	32	64	256	\	200
Pathfinder	1	32	64	512	\	200
Path-X	1	32	64	64	\	200

In this experiment, the encoder of Ultron was used as the neural network for classification and the temporal average of its output as its output layer. Ultron’s accuracy, generalization ability, computational efficiency, and memory usage were compared with those of some Transformer-based models. To ensure a fair comparison of the various models, the same model configuration as [40] was followed, including structure and number of iterations. Table 2 summarizes the neural-network hyperparameters used in each task. Table 3 shows the accuracy of each model in each subtask of LRA. Even though Ultron was limited by the low number of iterations, which prevented it from converging well in these tasks, its accuracy still surpassed that of all the Transformer-based models in all the tasks.

Table 3

Experimental results for the LRA benchmark. The best Model is in Boldface and the Second Best is Underlined.

Model	ListOps	Text	Retrieval	Image	Pathfinder	Path-X	Avg
Transformer	36.37	64.27	57.46	42.44	71.40	FAIL	54.39
Local Attention [60]	15.82	52.98	53.39	41.46	66.63	FAIL	46.06
Sparse Trans [61]	17.07	63.58	59.59	44.24	71.71	FAIL	51.24
Longformer [62]	35.63	62.85	56.89	42.22	69.71	FAIL	53.46
Linformer [63]	35.70	53.94	52.27	38.56	76.34	FAIL	51.36
Reformer [64]	37.27	56.10	53.40	38.07	68.50	FAIL	50.67
Sinkhorn Trans [65]	33.67	61.20	53.83	41.23	67.45	FAIL	51.39
Synthesizer [66]	<u>36.99</u>	61.68	54.67	41.61	69.45	FAIL	52.88
BigBird [67]	<u>36.05</u>	64.02	<u>59.29</u>	40.83	74.87	FAIL	<u>55.01</u>
Linear Trans [68]	16.13	<u>65.90</u>	53.09	42.34	75.30	FAIL	50.55
Performer [69]	18.01	<u>65.40</u>	53.82	<u>42.77</u>	<u>77.05</u>	FAIL	51.41
Our Work	38.49	68.91	60.88	50.24	79.71	68.27	61.08

Following [40]’s experimental setup, all models were evaluated on byte-level classification tasks using the same batch size. Figs. 3 and 4 present a comparison of speed and memory usage among various models with different input lengths (1 K, 2 K, 3 K, and 4 K). Ultron outperformed some of the Transformer-based sequence-processing models in terms of computational speed and memory usage.

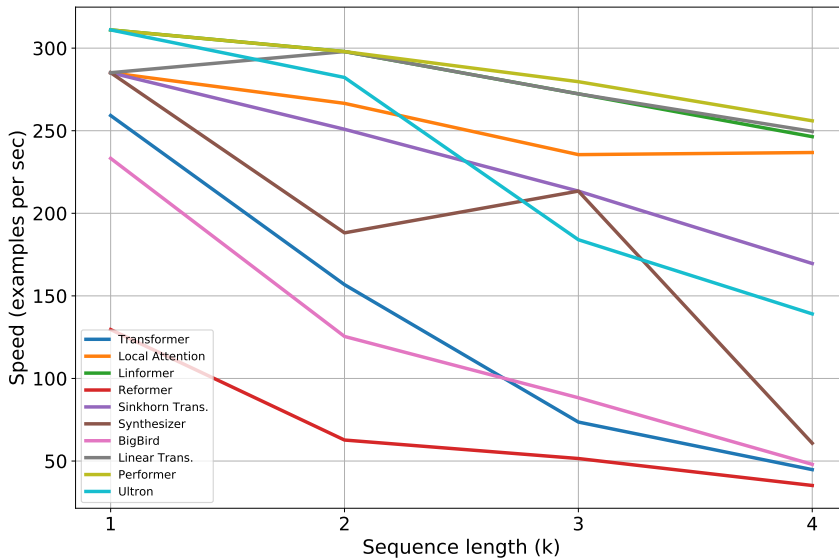


Fig. 3. Comparison of speed among various models with different input lengths.

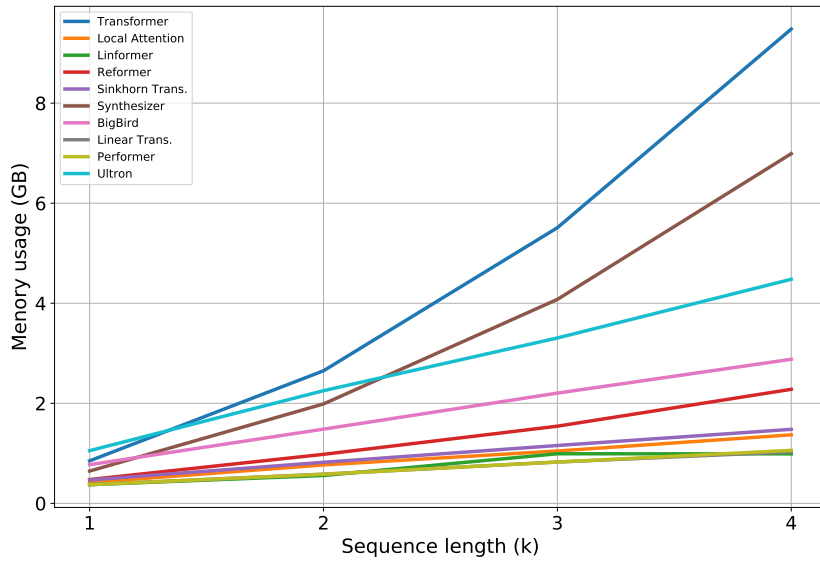


Fig. 4. Comparison of memory usage among various models with different input lengths.

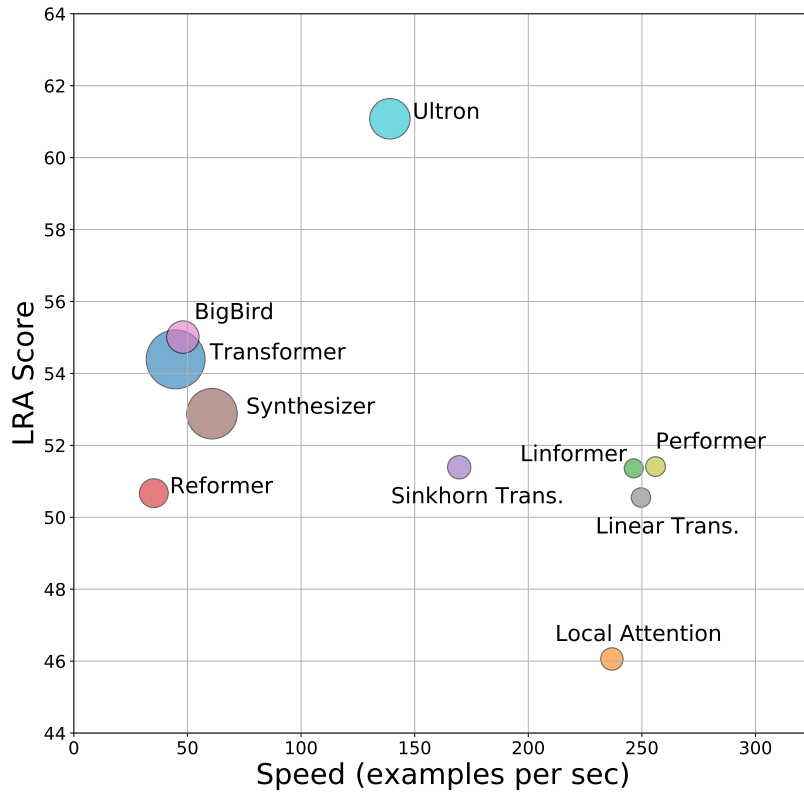


Fig. 5. Performance (y-axis), speed (x-axis), and memory footprint (size of the circles) of various models.

Fig. 5 shows the trade-off between memory usage, computational speed, and accuracy among various models. Most models that reduce the complexity of Transformer-based models sacrifice accuracy. However, our model achieves both higher accuracy and lower time and space complexity than Transformer. A noticeable drawback of Transformer and its variants is their low level of robustness. Unlike attention-based models, the Ultron model uses position-related features on sequences, which have a higher inductive bias for information that is close in spatial distance. This makes it robust in long-sequence tasks. Table 4 shows the training and testing accuracy of various models on the image-classification benchmark. The neural networks in the comparison use strong regularization algorithms, but they suffer from severe overfitting, with a large gap between their training and testing accuracy. In contrast, Ultron does not use any regularization algorithms in this experiment, but it has the highest testing accuracy, and its training and testing accuracy are very close. This shows that Ultron has better robustness than the other models.

Table 4

Testing and training accuracy of various models on the image-classification task.

Model	Test Accuracy	Train Accuracy
Transformer	42.44	69.45
Local Attention	41.46	63.19
Linformer	38.56	97.23
Reformer	38.07	68.45
Sinkhorn Trans	41.23	69.21
Synthesizer	41.61	97.31
BigBird	40.83	71.49
Linear Trans	42.34	65.61
Performer	42.77	73.90
Our Work	50.24	52.60

4.2. Natural language modeling

The language modeling performance of the proposed model’s decoder was compared with that of the GPT model on various natural language datasets, including wikitext-103, enwik8, and zhwiki520, which is a Chinese language dataset created for this study. Three metrics were used to evaluate the models: loss, computational speed, and memory usage. The model sizes were varied and the metrics were measured on each dataset. The results show that the proposed model outperforms the GPT model in all metrics with the same structural parameters.

4.2.1. Token-level language modeling on wikitext-103

In this section, three sets of comparative experiments are presented, each based on one of the three structures of GPT-2. The GPT-2 model is a pretrained Transformer-based model that, like other GPT models, consists of self-attention mechanisms and feedforward neural networks. The performance of these models is evaluated using the wikitext-103 dataset, which is a large English corpus dataset containing 1.03 million tokens from English Wikipedia articles. 1024 consecutive tokens are randomly selected from the documents as the input sequence of the neural network.

Table 5

Structure parameters and experimental results for Ultron and GPT in wikitext-103.

Model	N	d_{model}	d_{ff}	PPL
GPT-small	12	768	3072	37.50
GPT-medium	24	1024	4096	26.37
GPT-large	32	1280	5120	22.05
Ultron-small	12	768	3072	23.69
Ultron-medium	24	1024	4096	21.22
Ultron-large	32	1280	5120	19.54

Table 5 shows the three model-size groups that are compared in this experiment. No additional training datasets are used for Ultron. The results indicate that Ultron outperforms GPT-2 in all the experiments, achieving lower perplexity (PPL) and demonstrating superior performance in English token-level text-generation tasks.

4.3. Byte-level language modeling

This study investigated the challenges of language modeling for Chinese, a language that lacks explicit word boundaries. Approaches that use words as the smallest input units depend on the performance and efficiency of segmentation tools, which may affect the model’s performance. Chinese vocabulary also exhibits many variations, such as synonyms, near-synonyms, homophones, and variant characters, which may lead to semantic loss or ambiguity due to tokenization [77]. Therefore, the gb18030 encoding method was adopted to encode Chinese text into byte-level-sequence data. Byte-level units are character-based, so they do not require segmentation tools. They can theoretically handle any language, and they avoid out-of-vocabulary tokens. They also have a high degree of flexibility. Moreover, because the vocabulary size is 256, they can significantly reduce the parameter size and complexity of the word-embedding module.

The text from the main body of the entries in Chinese Wikipedia on May 20, 2023 was collected, and text information was extracted based on Su’s script [39]. A total of 1 355 711 entries were obtained, of which 13 557 entries were split as a test set, and the rest were used as a train set. The character `\x00` was used as a separator between each entry in the dataset. After preprocessing, the training set contained 1 868 574 744 characters and the test set contained 18 817 080 characters. For performance testing, three scales of neural networks with the same parameter settings presented in Table 5 were used. The neural network input was 1024 consecutive characters randomly cropped from the document. Fig. 6 shows the experimental results.

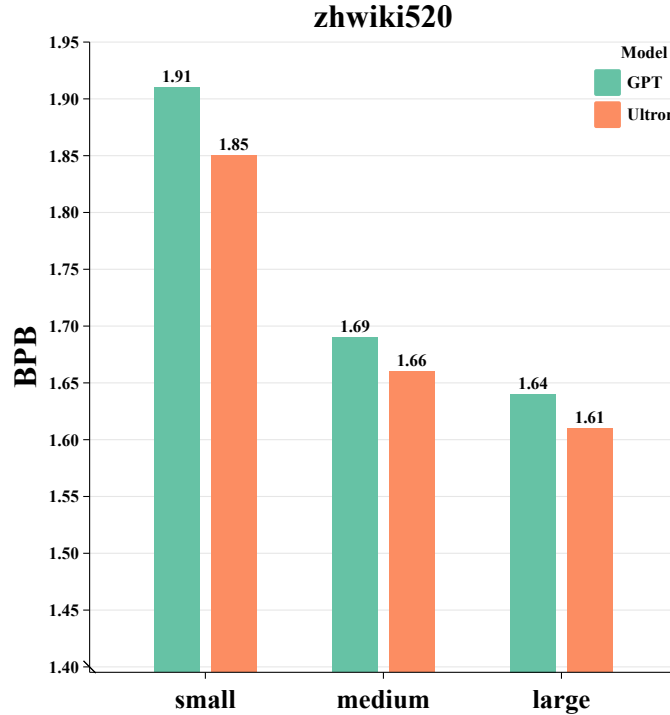


Fig. 6. Experimental results for Ultron and GPT in zhwiki520.

We also evaluated this model’s byte-level English language modeling performance. The enwik9 dataset [78] was used to train the model, and the enwik8 dataset’s last 10 M characters were employed for testing. The content of the test set was intentionally excluded from the training set. Fig. 7 shows the experimental results. As Figs. 6 and 7 show, Ultron outperformed the GPT model in both Chinese and English byte-level language modeling tasks. In particular, its bit per character (BPB) achieved state-of-the-art performance (0.93) [79] in the enwik8 test. This demonstrates Ultron’s superior performance in byte-level language modeling.

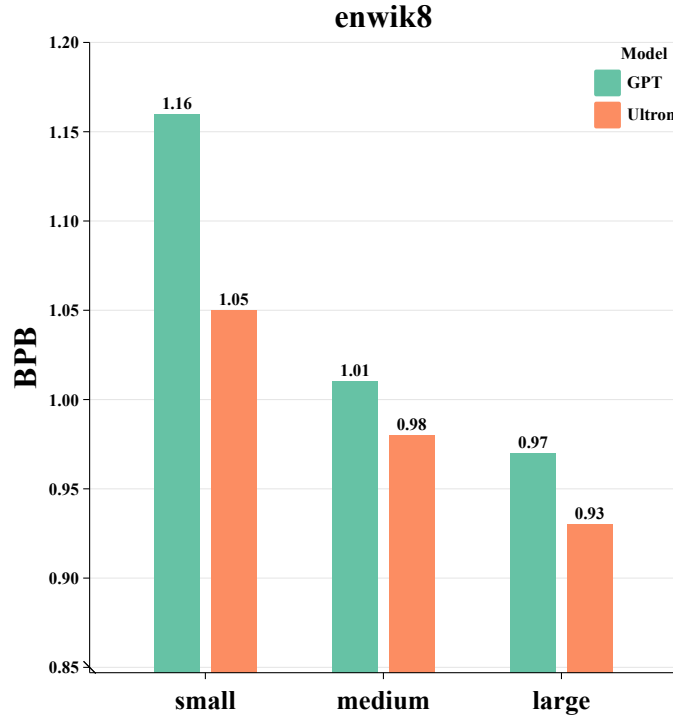


Fig. 7. Experimental results for Ultron and GPT in enwik8.

5. Conclusion and future work

This paper presents Ultron, a novel sequence-processing model that surpasses Transformer in terms of complexity, robustness, and fitting ability. Ultron’s CM layer captures long-term information dependencies by combining mapping functions, which reduces the complexity to $O(L \log L)$ and enables parallel computing. This advantage becomes more evident as the sequence length increases. In the encoder experiments, Ultron outperforms Transformer and its variants in all the tasks. In the decoder experiments, Ultron has a lower degree of loss than the GPT model of the same scale on English and Chinese language modeling tasks. On the enwik8 dataset, Ultron meets the current criteria for state-of-the-art performance. Our study demonstrates the effectiveness and potential of Ultron, which have important implications for NLP.

In future research, the aim is to apply Ultron to a broader array of domains, including image classification, signal detection, and translation. Given its exceptional fitting capability, Ultron demonstrates versatility across a wide spectrum of applications. At the algorithmic level, efforts will be directed toward further optimization of Ultron’s memory utilization and computing speed, as well as the development of a more tailored regularization algorithm. It is acknowledged that the experiments conducted are constrained by the available computing resources, and plans are in place to assess Ultron’s performance on larger parameter sizes and longer sequences in forthcoming research. The aspiration is that this paper will serve as a source of inspiration for subsequent studies in this field.

CRedit authorship contribution statement

Lingkai Hu: Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Writing – original draft. **Feng Zhan:** Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Writing – original draft. **Wenkai Huang:** Conceptualization, Methodology, Supervision, Writing – review & editing. **Weiming Gan:** Methodology, Validation, Formal analysis, Investigation, Data curation, Writing – original draft. **Haoliang Hu:** Validation, Formal analysis, Investigation, Data curation.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

The data are available based on a request from the corresponding author.

Acknowledgments

This work was supported by Guangzhou Youth Science and Technology Education Project under Grant KP2024403.

References

- [1] Puneet Kumar and Balasubramanian Raman. A bert based dual-channel explainable text emotion recognition system. *Neural Networks*, 150:392–407, 2022.
- [2] Wei Fang, Zhaofei Yu, Zhaokun Zhou, Ding Chen, Yanqi Chen, Zhengyu Ma, Timothée Masquelier, and Yonghong Tian. Parallel spiking neurons with high efficiency and ability to learn long-term dependencies. *Advances in Neural Information Processing Systems*, 36, 2024.
- [3] Jia Gong, Lin Geng Foo, Yixuan He, Hossein Rahmani, and Jun Liu. Llms are good sign language translators. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18362–18372, 2024.
- [4] Chao Li, Ning Bian, Ziping Zhao, Haishuai Wang, and Björn W Schuller. Multi-view domain-adaptive representation learning for eeg-based emotion recognition. *Information Fusion*, 104:102156, 2024.
- [5] Shih-Lun Wu, Chris Donahue, Shinji Watanabe, and Nicholas J Bryan. Music controlnet: Multiple time-varying controls for music generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32:2692–2703, 2024.
- [6] Abdinabi Mukhamadiyev, Mukhridin Mukhiddinov, Ilyos Khujayarov, Mannon Ochilov, and Jinsoo Cho. Development of language models for continuous uzbek speech recognition system. *Sensors*, 23(3):1145, 2023.
- [7] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27, 2014.
- [8] Biao Zhang, Deyi Xiong, Jun Xie, and Jinsong Su. Neural machine translation with gru-gated attention model. *IEEE transactions on neural networks and learning systems*, 31(11):4688–4698, 2020.
- [9] Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. Beyond english-centric multilingual machine translation. *The Journal of Machine Learning Research*, 22(1):4839–4886, 2021.
- [10] Ben Saunders, Necati Cihan Camgoz, and Richard Bowden. Continuous 3d multi-channel sign language production via progressive transformers and mixture density networks. *International journal of computer vision*, 129(7):2113–2135, 2021.
- [11] Shuang Li, Jiangjie Chen, Siyu Yuan, Xinyi Wu, Hao Yang, Shimin Tao, and Yanghua Xiao. Translate meanings, not just words: Idiomkb’s role in optimizing idiomatic translation with language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18554–18563, 2024.
- [12] Daria Grechishnikova. Transformer neural network for protein-specific de novo drug generation as a machine translation problem. *Scientific reports*, 11(1):321, 2021.
- [13] Shaun Mahony, Panayiotis V Benos, Terry J Smith, and Aaron Golden. Self-organizing neural networks to support the discovery of dna-binding motifs. *Neural Networks*, 19(6-7):950–962, 2006.
- [14] Jing Li, Botong Wu, Xinwei Sun, and Yizhou Wang. Causal hidden markov model for time series disease forecasting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12105–12114, 2021.
- [15] Jie Ren, Kai Song, Chao Deng, Nathan A Ahlgren, Jed A Fuhrman, Yi Li, Xiaohui Xie, Ryan Poplin, and Fengzhu Sun. Identifying viruses from metagenomic data using deep learning. *Quantitative Biology*, 8:64–77, 2020.
- [16] Tianyang Lin, Yuxin Wang, Xiangyang Liu, and Xipeng Qiu. A survey of transformers. *AI Open*, 2022.
- [17] Abhijit Mahalunkar and John D Kelleher. Understanding recurrent neural architectures by analyzing and synthesizing long distance dependencies in benchmark sequential datasets. *arXiv preprint arXiv:1810.02966*, 2018.
- [18] Mostafa M Amin, Rui Mao, Erik Cambria, and Björn W Schuller. A wide evaluation of chatgpt on affective computing tasks. *IEEE Transactions on Affective Computing*, 2024.
- [19] Robin M Schmidt. Recurrent neural networks (rnns): A gentle introduction and overview. *arXiv preprint arXiv:1912.05911*, 2019.
- [20] Yong Yu, Xiaosheng Si, Changhua Hu, and Jianxun Zhang. A review of recurrent neural networks: Lstm cells and network architectures. *Neural computation*, 31(7):1235–1270, 2019.
- [21] Colin Lea, Michael D Flynn, Rene Vidal, Austin Reiter, and Gregory D Hager. Temporal convolutional networks for action segmentation and detection. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 156–165, 2017.
- [22] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- [23] Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. Gated feedback recurrent neural networks. In *International conference on machine learning*, pages 2067–2075. PMLR, 2015.

- [24] Mehmet Ozgur Turkoglu, Stefano D'Aronco, Jan Dirk Wegner, and Konrad Schindler. Gating revisited: Deep multi-layer rnns that can be trained. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(8):4081–4092, 2021.
- [25] Ossama Abdel-Hamid, Abdel-rahman Mohamed, Hui Jiang, Li Deng, Gerald Penn, and Dong Yu. Convolutional neural networks for speech recognition. *IEEE/ACM Transactions on audio, speech, and language processing*, 22(10):1533–1545, 2014.
- [26] Yihao Xue, Rui Yang, Xiaohan Chen, Weibo Liu, Zidong Wang, and Xiaohui Liu. A review on transferability estimation in deep transfer learning. *IEEE Transactions on Artificial Intelligence*, 2024.
- [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [28] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [29] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [30] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32, 2019.
- [31] Sneha Chaudhari, Varun Mithal, Gungor Polatkan, and Rohan Ramanath. An attentive survey of attention models. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 12(5):1–32, 2021.
- [32] Dmitrii Kochkov, Janni Yuval, Ian Langmore, Peter Norgaard, Jamie Smith, Griffin Mooers, Milan Klöwer, James Lottes, Stephan Rasp, Peter Düben, et al. Neural general circulation models for weather and climate. *Nature*, pages 1–7, 2024.
- [33] Ibomoiye Domor Mienye, Theo G Swart, and George Obaido. Recurrent neural networks: A comprehensive review of architectures, variants, and applications. *Information*, 15(9):517, 2024.
- [34] Dušan Variš and Ondřej Bojar. Sequence length is a domain: Length-based overfitting in transformer models. *arXiv preprint arXiv:2109.07276*, 2021.
- [35] Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernocký, and Sanjeev Khudanpur. Recurrent neural network based language model. In *Interspeech*, volume 2, pages 1045–1048. Makuhari, 2010.
- [36] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*, 2016.
- [37] Matt Mahoney. Large text compression benchmark, 2011.
- [38] Wikimedia Foundation. Index of zhwiki. <https://dumps.wikimedia.org/zhwiki/>, 2023. Accessed on Sep 20, 2023.
- [39] Jianlin Su. Obtain and process chinese wikipedia corpus, Jan 2017.
- [40] Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena: A benchmark for efficient transformers. *arXiv preprint arXiv:2011.04006*, 2020.
- [41] Zhici Wang, Qiancheng Yu, Jinyun Wang, Zhiyong Hu, and Aoqiang Wang. Grammar correction for multiple errors in chinese based on prompt templates. *Applied Sciences*, 13(15):8858, 2023.
- [42] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [43] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [44] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [45] Jaeyoung Kim, Mostafa El-Khamy, and Jungwon Lee. Residual lstm: Design of a deep recurrent architecture for distant speech recognition. *arXiv preprint arXiv:1701.03360*, 2017.
- [46] Boxuan Yue, Junwei Fu, and Jun Liang. Residual recurrent neural networks for learning sequential representations. *Information*, 9(3):56, 2018.
- [47] Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, et al. Rkv: Reinventing rnns for the transformer era. *arXiv preprint arXiv:2305.13048*, 2023.
- [48] Zhicai Zhao, Na Lv, Runquan Xiao, and Shanben Chen. A novel penetration state recognition method based on lstm with auditory attention during pulsed gtaw. *IEEE Transactions on Industrial Informatics*, 2022.
- [49] Wei Fang, Yupeng Chen, and Qiongying Xue. Survey on research of rnn-based spatio-temporal sequence prediction algorithms. *Journal on Big Data*, 3(3):97, 2021.
- [50] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
- [51] Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. A convolutional neural network for modelling sentences. *arXiv preprint arXiv:1404.2188*, 2014.
- [52] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017.
- [53] Yang Xian, Yang Sun, Wenwu Wang, and Syed Mohsen Naqvi. Convolutional fusion network for monaural speech enhancement. *Neural Networks*, 143:97–107, 2021.
- [54] Yihao Xue, Rui Yang, Xiaohan Chen, Zhongbei Tian, and Zidong Wang. A novel local binary temporal convolutional neural network for bearing fault diagnosis. *IEEE Transactions on Instrumentation and Measurement*, 2023.
- [55] Ismail Khalfauui-Hassani, Thomas Pellegrini, and Timothée Masquelier. Dilated convolution with learnable spacings. *arXiv preprint arXiv:2112.03740*, 2021.
- [56] Tao Shen, Tianyi Zhou, Guodong Long, Jing Jiang, Shirui Pan, and Chengqi Zhang. Disan: Directional self-attention network for rnn/cnn-free language understanding. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [57] Zhaoyang Niu, Guoqiang Zhong, and Hui Yu. A review on the attention mechanism of deep learning. *Neurocomputing*, 452:48–62, 2021.
- [58] Michael Hahn. Theoretical limitations of self-attention in neural sequence models. *Transactions of the Association for Computational Linguistics*, 8:156–171, 2020.

- [59] Chuhan Wu, Fangzhao Wu, Tao Qi, Yongfeng Huang, and Xing Xie. Fastformer: Additive attention can be all you need. *arXiv preprint arXiv:2108.09084*, 2021.
- [60] Moab Arar, Ariel Shamir, and Amit H Bermano. Learned queries for efficient local attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10841–10852, 2022.
- [61] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- [62] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- [63] Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- [64] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*, 2020.
- [65] Yi Tay, Dara Bahri, Liu Yang, Donald Metzler, and Da-Cheng Juan. Sparse sinkhorn attention. In *International Conference on Machine Learning*, pages 9438–9447. PMLR, 2020.
- [66] Y Tay, D Bahri, D Metzler, D Juan, Z Zhao, and C Zheng. Synthesizer: Rethinking self-attention in transformer models. *arxiv* 2020. *arXiv preprint arXiv:2005.00743*, 2, 2020.
- [67] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. *Advances in neural information processing systems*, 33:17283–17297, 2020.
- [68] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, pages 5156–5165. PMLR, 2020.
- [69] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, David Belanger, Lucy Colwell, et al. Masked language modeling for proteins via linearly scalable long-context transformers. *arXiv preprint arXiv:2006.03555*, 2020.
- [70] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International conference on machine learning*, pages 3744–3753. PMLR, 2019.
- [71] Xuezhe Ma, Xiang Kong, Sinong Wang, Chunting Zhou, Jonathan May, Hao Ma, and Luke Zettlemoyer. Luna: Linear unified nested attention. *Advances in Neural Information Processing Systems*, 34:2441–2453, 2021.
- [72] Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang, Shaohan Huang, Wenhui Wang, and Furu Wei. Longnet: Scaling transformers to 1,000,000,000 tokens. *arXiv preprint arXiv:2307.02486*, 2023.
- [73] Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. Retentive network: A successor to transformer for large language models. *arXiv preprint arXiv:2307.08621*, 2023.
- [74] Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse attention with routing transformers. *Transactions of the Association for Computational Linguistics*, 9:53–68, 2021.
- [75] Andrew L Maas, Awni Y Hannun, Andrew Y Ng, et al. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, page 3. Atlanta, GA, 2013.
- [76] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [77] Chu-Ren Huang, Petr Šimon, Shu-Kai Hsieh, and Laurent Prévot. Rethinking chinese word segmentation: tokenization, character classification, or wordbreak identification. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 69–72, 2007.
- [78] Marcus Hutter. Hutter prize for lossless compression of human knowledge, 2006. <http://prize.hutter1.net/>, Last accessed on 2023-05-26.
- [79] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.