

Attentive Multi-Layer Perceptron for Non-autoregressive Generation

Shuyang Jiang¹, Jun Zhang², Jiangtao Feng³, Lin Zheng⁴, and Lingpeng Kong⁴(✉)

¹ Shanghai Jiao Tong University, Shanghai, China
jiangshuyang@sjtu.edu.cn

² Shanghai Artificial Intelligence Laboratory, Shanghai, China
zhangjun@pjlab.org.cn

³ jiangtaofeng0906@gmail.com

⁴ The University of Hong Kong, Hong Kong, China
{linzheng@connect, lpk@cs}.hku.hk

Abstract. Autoregressive (AR) generation almost dominates sequence generation for its efficacy. Recently, non-autoregressive (NAR) generation gains increasing popularity for its efficiency and growing efficacy. However, its efficiency is still bottlenecked by quadratic complexity in sequence lengths, which is prohibitive for scaling to long sequence generation and few works have been done to mitigate this problem. In this paper, we propose a novel MLP variant, **Attentive Multi-Layer Perceptron (AMLP)**, to produce a generation model with linear time and space complexity. Different from classic MLP with static and learnable projection matrices, AMLP leverages adaptive projections computed from inputs in an attentive mode. The sample-aware adaptive projections enable communications among tokens in a sequence, and model the measurement between the query and key space. Furthermore, we marry AMLP with popular NAR models, deriving a highly efficient NAR-AMLP architecture with linear time and space complexity. Empirical results show that such marriage architecture surpasses competitive efficient NAR models, by a significant margin on text-to-speech synthesis and machine translation. We also test AMLP’s self- and cross-attention ability separately with extensive ablation experiments, and find them comparable or even superior to the other efficient models. The efficiency analysis further shows that AMLP extremely reduces the memory cost against vanilla non-autoregressive models for long sequences.

Keywords: AMLP · Multi-Layer Perceptron · Attention Mechanism · Non-Autoregressive Model.

1 Introduction

Attention-based sequence generation methods have achieved great success and gained increasing popularity in machine learning [53,30,35,11]. A large body of research in neural architectures has been devoted to the autoregressive (AR)

method [40,41], where tokens are generated one after another in an iterative manner. The computational overhead in decoding can thus be prohibitive, especially for long sequences. Recently, non-autoregressive (NAR) generation attracts more attention for its efficiency and growing efficacy [17,18,42,43,46,7]. In a non-autoregressive model, the decoder generates the target sequence all at once, significantly reducing its computational overhead at the inference stage. Nevertheless, relatively little research has been done on the attention architecture in non-autoregressive models. In particular, the conventionally adopted softmax attention comes with a quadratic time and memory cost. It is therefore still difficult to scale up non-autoregressive models to long sequence generation tasks.

In this paper, we propose Attentive Multi-Layer Perceptron (§2.2; AMLP) to integrate the attention mechanism with the multi-layer perceptron (MLP) in non-autoregressive architecture, resulting in a fully parallelizable sequence generation model with linear complexity. Unlike the widely-used MLP whose weights are invariant across different sequences, we compute the weights in AMLP through adaptive projections from (multiple) input tokens and model their interactions in an attentive manner. Specifically, we put forward two methods (§2.3) to compute the adaptive projections in AMLP, which implicitly model the association between the query and key space. We utilize the simplicity and efficiency of MLP while obtaining the strong modeling capability of AMLP for input tokens' communication. Finally, we present a hybrid NAR-AMLP model (§2.4) to achieve both linear complexity and high parallelism.

We evaluate the AMLP architecture on text-to-speech synthesis for a relatively long sequence scenario and machine translation for a relatively short sequence scenario. Experiments show that AMLP achieves more superior scores with objective measurements compared with the strong softmax attention counterpart (§3.1) on text-to-speech synthesis, with less computational cost (§3.3). On machine translation, AMLP performs competitive with vanilla attention but achieves the best result among efficient NAR and AR models with linear complexity (§3.1). Further, we test the self- and cross-attention ability of AMLP on super resolution and long sequence time-series forecasting tasks, respectively. Empirical results show that AMLP is on par with other efficient attention in self-attention and achieves the best performance in cross-attention scenarios (§3.2). Additionally, when scaling to long sequence, AMLP reduces the memory footprint substantially and further improves the inference speed in NAR models (§3.3). The code is available in <https://github.com/Shark-NLP/AttentiveMLP>.

2 Non-Autoregressive Generation with Attentive MLP

In this section, we first give a brief introduction to autoregressive (AR) and non-autoregressive (NAR) generation, and then delve into the nuances that differentiate the attention mechanisms utilized in autoregressive (AR) and non-autoregressive (NAR) models. After that, we present the AMLP architecture to model the communication among sequence tokens. Finally, we build up an NAR-AMLP architecture with linear time and space complexity.

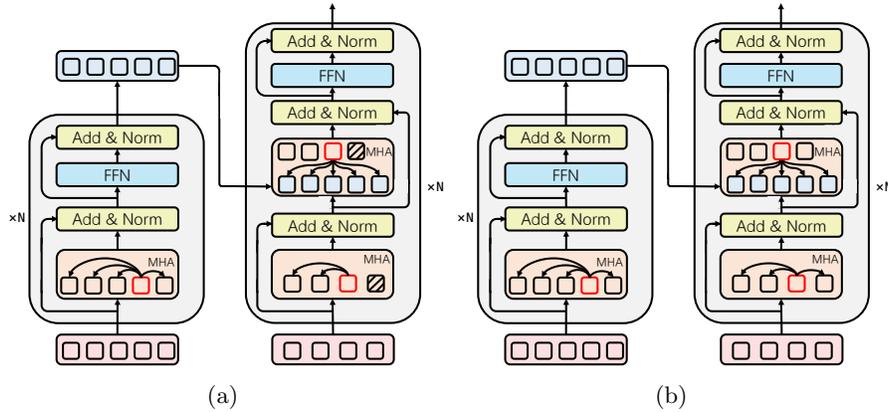


Fig. 1: AR (a) and NAR (b) encoder-decoder architectures. “MHA” stands for multi-head attention. Blocks with red rims represent the current state token. Shaded blocks represent future tokens that are invisible to the current state.

2.1 Background: Autoregressive and Non-Autoregressive Generation

Given a source sequence $X_{1:m}$, conditional sequence generation targets to predict a target sequence $Y_{1:n}$ by modeling the conditional probability $p(Y|X)$. Autoregressive generation decomposes the probability $p(Y|X)$ as:

$$p(Y|X) = \prod_{i=1..n} p(Y_i|Y_{<i}, X), Y_{<1} = \emptyset. \quad (1)$$

which is implemented as a typical encoder-decoder architecture shown in Fig. 1a. Although such decomposition is proved effective, it suffers from two main drawbacks: efficiency and exposure bias. On the one hand, the autoregressive decoding process, where each token depends on the previous predicted ones, prevents the model from fast inference in usage. On the other hand, teacher-forcing exposes ground truth tokens in network inputs during the training process, where the exposed tokens are unable to observe in inference. Such exposure creates an inconsistency between the training and inference, and harms the prediction quality.

Recently, non-autoregressive generation, depicted as Fig. 1b, shows its capability of sequence modeling in terms of both efficiency and efficacy, which decomposes the conditional probability $p(Y|X)$ via a Naïve Bayes assumption:

$$p(Y|X) = \prod_{i=1..n} p(Y_i|X) \quad (2)$$

The NAR decomposition enables parallel decoding for each token, and speeds up the inference process substantially. Although NAR generation is much faster than AR generation, its speed is still limited by the $O(n^2 + nm + m^2)$ time complexity of the multi-head softmax attention module. This is especially problematic in modeling long sequences.

Attention Types in AR & NAR Models Although autoregressive and non-autoregressive models differ from each other in sequence generation paradigms, their underlying attention mechanisms in their architectures are also different. The token-by-token generation of AR models requires a causal decoder that forces tokens to attend to only previous features. A typical causal decoder utilizes causal softmax attentions both in self-attention and cross-attention. The attention causality entails that during the computation, it is important to ensure that the query token does not attend to the context on its right side, just as the shaded blocks in Fig. 1a. In contrast, the NAR model, which allows for parallel generation of the output sequence and global contextualization using attention, employs a noncausal decoder in Fig. 1b. The self-attention in the NAR model can attend to both side contexts of a given token, which makes it suitable for tasks that require a broader contextual understanding. NAR architectures also reduce the design restrictions on cross-attention, making query tokens attend to key tokens in a holistic view. This modeling feature of attention emphasizes both global and local contextualization modeling for attention modules. In practice, causality in vanilla softmax self-attention is ensured by leveraging a lower triangular mask in AR models, while linearized attention requires more sophisticated implementation. Since no causality is required in NAR models, designing an efficient attention mechanism is much more flexible.

2.2 Attentive Multi-Layer Perceptron

Modeling interactions between tokens is crucial and challenging in sequence generation. Transformer [53] stacks the MLP, which aims to learn features of individual tokens, on top of the attention block, which is responsible for modeling the communication within the sequence. In AR generation, the attention needs to be recomputed for each time step through the recurrent process, as the key and value set is changing. However, this procedure is non-causal in NAR generation. We therefore are able to integrate the modeling of token interactions into the MLP architecture and make the whole architecture fully parallelizable and more efficient.

Given a sequence representation $\mathbf{X} \in \mathbb{R}^{n \times d}$, where n is the sequence length and d is dimensionality of the feature space, the conventional MLP models the feature of individual token $\mathbf{X}_i \in \mathbb{R}^d$ as:

$$\text{MLP}(\mathbf{X}_i) = \sigma(\mathbf{X}_i W_1) W_2 \quad (3)$$

where $W_1 \in \mathbb{R}^{d \times d_h}$, $W_2 \in \mathbb{R}^{d_h \times d}$ are learnable parameters d_h is the dimensionality of hidden space. $\sigma(\cdot)$ is a non-linear activation function such as $\text{ReLU}(\cdot)$. However, it disables the communication between tokens in the sequence, and prevents the model from learning contextualized token representations.

A widely-used approach to enable communication between each token in a sequence is the attention mechanism [53]. Vanilla attention learns to incorporate source sequence features $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{m \times d}$ into target $\mathbf{Q} \in \mathbb{R}^{n \times d}$ with an attention matrix

$$\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(\mathbf{QK}^\top) \mathbf{V} \quad (4)$$

where m, n are the source and target length respectively. Here we omit the input projections for $\mathbf{Q}, \mathbf{K}, \mathbf{V}$, the output projection, and the scaling factor $1/\sqrt{d}$ for simplicity.

The motivation of Attentive Multi-Layer Perceptron (AMLP) starts from the fact that the vanilla softmax attention can be viewed as a projection function as $\text{SA}(\cdot|\mathbf{K}, \mathbf{V}) : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ which projects the original $\mathbf{Q} \in \mathbb{R}^{n \times d}$ with \mathbf{K} and \mathbf{V} features as its context while preserving \mathbf{Q} 's shape. In vanilla attention, $\text{softmax}(\mathbf{Q}\mathbf{K}^\top)$ is a softmax kernel which can be decomposed into a multiplication of two kernel functions: $\phi(\mathbf{Q}) \cdot \phi(\mathbf{K})^\top$, which is verified in Performer [10], cosFormer [44] and LARA [60]. Meanwhile, the low-rank factorization of the attention matrix, $\text{softmax}(\mathbf{Q}\mathbf{K}^\top)$, does not impact the performance much, which is verified by Nyströmformer [57]. Based on their findings, we propose an alternative modeling solution by fusing key $\mathbf{K} \in \mathbb{R}^{m \times d}$ and value $\mathbf{V} \in \mathbb{R}^{m \times d}$ information into query $\mathbf{Q} \in \mathbb{R}^{n \times d}$, via a symmetric and positive semi-definite distance matrix $\mathbf{\Sigma} \in \mathbb{R}^{d \times d}$ on \mathbf{Q} and \mathbf{K} space. The contextualizing process on \mathbf{Q} can be formulated as:

$$f(\mathbf{Q}; \mathbf{K}, \mathbf{V}) = \mathbf{Q}\mathbf{\Sigma}\mathbf{K}^\top\mathbf{V} \quad (5)$$

where $\mathbf{\Sigma}$ is computed from \mathbf{Q} and \mathbf{K} .

With similar functionality to [10,57], the matrix $\mathbf{Q}\mathbf{\Sigma}\mathbf{K}^\top$ can also enjoy lower computation costs from low-rank approximation while maintaining strong modeling capability. Without taking any low-rank assumptions on input \mathbf{Q}, \mathbf{K} , we decompose the distance matrix as:

$$\mathbf{\Sigma} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top = \mathbf{U}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{U}^\top \approx \mathbf{U}\hat{\mathbf{\Lambda}}^{\frac{1}{2}}\hat{\mathbf{\Lambda}}^{\frac{1}{2}}\mathbf{U}^\top = (\mathbf{U}\hat{\mathbf{\Lambda}}^{\frac{1}{2}})(\mathbf{U}\hat{\mathbf{\Lambda}}^{\frac{1}{2}})^\top = \mathbf{L}\mathbf{L}^\top \quad (6)$$

where \mathbf{U} is the orthogonal eigenvector of matrix and $\mathbf{\Lambda}$ is the diagonal eigenvalues matrix. $\hat{\mathbf{\Lambda}}$ here is an approximation to $\mathbf{\Lambda}$ by keeping largest- c eigen-values and masking the others with 0, where c is a hyper-parameter in AMLP. Thus we derive a decomposition equation $\mathbf{\Sigma} \approx \mathbf{L}\mathbf{L}^\top$ where $\mathbf{L} = \kappa(\mathbf{Q}, \mathbf{K})^\top \in \mathbb{R}^{d \times c}$ indicates a low-rank matrix. We will show two different methods for parameterization of \mathbf{L} , resulting in two different AMLP variants. We rewrite Eq. 5 by decomposing the distance matrix $\mathbf{\Sigma}$ as:

$$f(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \approx \mathbf{Q}\mathbf{L}\mathbf{L}^\top\mathbf{K}^\top\mathbf{V} \quad (7)$$

Now Eq. 5 could be approximated with Eq. 7 by linearly projecting the original \mathbf{Q} with adaptive weights twice. By reordering the computation and adding nonlinearity into Eq. 5, we derive a general form of AMLP model as:

$$\text{AMLP}(\mathbf{Q}; \mathbf{K}, \mathbf{V}) = \sigma_1(\mathbf{Q}W_{\mathbf{Q},\mathbf{K}})W_{\mathbf{Q},\mathbf{K},\mathbf{V}} \quad (8)$$

where the nonlinear function $\sigma_1(\cdot)$ can be adjusted arbitrarily. Eq. 8 address the general form of AMLP, and the adaptive weights $W_{\mathbf{Q},\mathbf{K}}$ and $W_{\mathbf{Q},\mathbf{K},\mathbf{V}}$ can be specified in various ways. Following the form of Eq. 8, we will further introduce two AMLP variants in § 2.3, by specifying $\mathbf{L} = W_{\mathbf{Q},\mathbf{K}} = \kappa(\mathbf{Q}, \mathbf{K})$, computational order and nonlinear function.

The computation of adaptive weights in AMLP fuses token-level communication, while MLP models tokens in a sequence independently. Therefore, AMLP

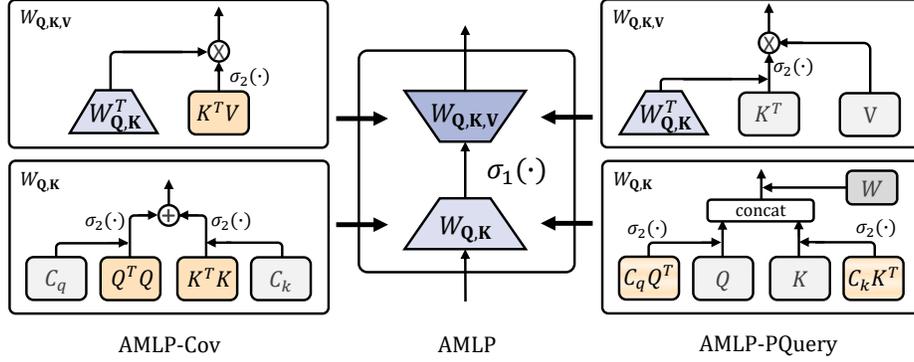


Fig. 2: Computation diagram of two AMLP variants. The middle part shows the computation of basic AMLP. The Left and right figures show the detailed computation of two adaptive weight matrices in AMLP-Cov and AMLP-PQuery.

enables the communication between tokens in a sequence. And different from vanilla softmax attention, AMLP utilizes a distance matrix Σ between \mathbf{Q} and \mathbf{K} spaces to fuse information among their contexts and outputs a contextualized \mathbf{Q} . Through this distance matrix, AMLP computes the similarity between \mathbf{Q} and \mathbf{K} like softmax attention, and leverages it to aggregate \mathbf{V} .

2.3 Parameterization

In this section, we describe two methods for the parameterization of two adaptive weight matrices $W_{\mathbf{Q},\mathbf{K}}$ and $W_{\mathbf{Q},\mathbf{K},\mathbf{V}}$. Fig. 2 illustrates the computation graph of these two methods.⁵

Cross-Covariance We present AMLP-Cov, a variant that adopts cross-covariance to parameterize $W_{\mathbf{Q},\mathbf{K}}$ and $W_{\mathbf{Q},\mathbf{K},\mathbf{V}}$. One challenge of AMLP is to fuse information of $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ of different shapes into static-shaped projection matrices $W_{\mathbf{Q},\mathbf{K}}$ and $W_{\mathbf{Q},\mathbf{K},\mathbf{V}}$. Inspired by [1], we propose to use \mathbf{Q}, \mathbf{K} 's covariance and the cross-covariance between \mathbf{K} and \mathbf{V} in AMLP. To obtain $\mathbf{L} = \kappa(\mathbf{Q}, \mathbf{K})^\top$, we separately compute \mathbf{Q} 's and \mathbf{K} 's covariance matrices and combines them with learned down-sampling projection matrices $C_q \in \mathbb{R}^{c \times d}$ and $C_k \in \mathbb{R}^{c \times d}$:

$$\kappa(\mathbf{Q}, \mathbf{K}) = C_q (\sigma_2(\mathbf{Q}^\top \mathbf{Q})) + C_k (\sigma_2(\mathbf{K}^\top \mathbf{K})) \quad (9)$$

where $\sigma_2(\cdot)$ is set to softmax function as [1] suggest. The covariance matrices of \mathbf{Q}, \mathbf{K} are of the same shape and can be directly fused. We add the softmax function as a non-linear activation to enhance the expressiveness. For $W_{\mathbf{Q},\mathbf{K},\mathbf{V}}$, we notice the shapes of \mathbf{K} and \mathbf{V} are usually identical, and we hence use their

⁵ AMLP is implemented with multiple heads [53], but for simplicity and without loss of generality, we will discuss our AMLP computation process in a single-head setting.

cross-covariance $\mathbf{K}^\top \mathbf{V}$ for computation in Eq. 8. $W_{\mathbf{Q},\mathbf{K},\mathbf{V}}$ is then formulated by transforming the cross-covariance $\mathbf{K}^\top \mathbf{V}$ to query space by \mathbf{L} as:

$$W_{\mathbf{Q},\mathbf{K},\mathbf{V}} = \mathbf{L}^\top \sigma_2(\mathbf{K}^\top \mathbf{V}) \quad (10)$$

Pseudo-Queries AMLP-PQuery first uses Exponential Moving Average (EMA) to compute the contextualized query via a hyperparameter β : $\hat{\mathbf{q}}_i = \beta \cdot \hat{\mathbf{q}}_{i-1} + (1 - \beta) \cdot \mathbf{q}_i$, which has been proved to model local context well [37]. To further improve the communication between target and source sequences in a long sequence view, AMLP-PQuery treats learnable C_q , C_k and \mathbf{L}^\top as pseudo attention queries. Specifically, it estimates $W_{\mathbf{Q},\mathbf{K}}$ by fusing features from query and key to the hidden space with an extra learnable weight $W \in \mathbb{R}^{2d \times d}$:

$$W_{\mathbf{Q},\mathbf{K}} = \mathbf{L}^\top = \left[\sigma_2(C_q \hat{\mathbf{Q}}^\top) \hat{\mathbf{Q}}; \sigma_2(C_k \mathbf{K}^\top) \mathbf{K} \right] W \quad (11)$$

where $\sigma_2(\cdot)$ is set to softmax as AMLP-Cov. For $W_{\mathbf{Q},\mathbf{K},\mathbf{V}}$, we notice that \mathbf{L}^\top has fused features from $\hat{\mathbf{Q}}$. So we again treat \mathbf{L}^\top as a pseudo query to fuse features from the source sequence:

$$W_{\mathbf{Q},\mathbf{K},\mathbf{V}} = \sigma_2(\mathbf{L}^\top \mathbf{K}^\top) \mathbf{V} \quad (12)$$

With explicit communication between $\hat{\mathbf{Q}}$ and \mathbf{K} in $W_{\mathbf{Q},\mathbf{K},\mathbf{V}}$, the alignment between different sequences is enhanced; therefore, AMLP-PQuery is more adaptive to cross-attention.

2.4 Linear NAR: A Hybrid Architecture of NAR and AMLP

We combine AMLP with NAR for lower memory costs, faster inference speed and higher parallelism because AMLP and NAR are mutually reinforcing.

AMLP boosts NAR On one hand, NAR parallelizes the inference process, but its efficiency is still hindered by vanilla attention. AMLP, as a plug-in efficient attentive module, mitigates the inefficiency effortlessly. On the other hand, the non-autoregressive pipeline provides a non-causal encoding framework, with which the computation of AMLP avoids fine-grained operations.

NAR augments AMLP We present the specific computation steps of AMLP in AR scenario and explain the drawbacks of AR-AMLP. We take AMLP-Cov as an example. Given an query token \mathbf{q}_t , the covariances $\mathbf{S}_t^{\mathbf{Q}}$ and $\mathbf{S}_t^{\mathbf{K}}$ of \mathbf{K}_t and \mathbf{Q}_t , and the cross-covariance \mathbf{z}_t of \mathbf{K}_t and \mathbf{V}_t , $W_{\mathbf{Q},\mathbf{K}}$ and $W_{\mathbf{Q},\mathbf{K},\mathbf{V}}$ are formulated as:

$$W_{\mathbf{Q}_t,\mathbf{K}_t} = \mathbf{L}_t^\top = C_q(\sigma_2(\mathbf{S}_t^{\mathbf{Q}})) + C_k(\sigma_2(\mathbf{S}_t^{\mathbf{K}})) \quad (13)$$

$$W_{\mathbf{Q}_t,\mathbf{K}_t,\mathbf{V}_t} = \mathbf{L}_t^\top \sigma_2(\mathbf{z}_t) \quad (14)$$

where $\mathbf{S}_t^{\mathbf{Q}} = \mathbf{S}_{t-1}^{\mathbf{Q}} + \mathbf{q}_t^{\top} \mathbf{q}_t$, $\mathbf{S}_t^{\mathbf{K}} = \mathbf{S}_{t-1}^{\mathbf{K}} + \mathbf{k}_t^{\top} \mathbf{k}_t$ and $\mathbf{z}_t = \mathbf{z}_{t-1} + \mathbf{k}_t^{\top} \mathbf{v}_t$. These computation steps increase heavy memory costs and large time consumption in the training phase, with an additional $O(ncd)$ costs beyond the overall computation. Recurrent computation also harms the parallelism and further slows down the training process, which is avoided naturally in NAR models. Moreover, CAB [59] points out that most existing efficient architectures suffer a great performance drop in causal-self or causal-cross pattern of AR models. Combining the two drawbacks brought by the fusion of efficient architecture and AR models, we decide to incorporate AMLP into NAR to produce a powerful and efficient model.

2.5 Complexity Analysis

Without loss of generality, we focus on the complexity in the typical encoder-decoder architecture and omit the independent factor *w.r.t.* target length n and source length m for simplicity.

AMLP-Cov & AMLP-PQuery Note that the inner dimension c is a constant to both m and n . The sequential computation of two adaptive projection matrices and the overall MLP computation in Eq. 8 are all of $O(n+m)$. The exclusive EMA submodule in AMLP-PQuery is $O(n)$ as well. Therefore, the time and memory complexity of AMLP (both AMLP-Cov and AMLP-PQuery) is $O(n+m)$.

NAR-AMLP Non-autoregressive models have one encoder self-attention, one decoder self-attention, and an encoder-decoder cross-attention. Due to the quadratic complexity of softmax attention, the complexities of the three attentions are $O(m^2)$, $O(n^2)$ and $O(nm)$, respectively. Therefore, the complexity of the entire model architecture is $O(n^2 + nm + m^2)$. To reduce the inefficiency bottlenecked by softmax attention, we replace softmax modules in non-autoregressive models with AMLP, deriving an NAR-AMLP architecture with linear time and space complexity.

3 Experiments

We conduct extensive experiments, covering the fields of speech, natural language processing, time-series and computer vision.⁶ For fair comparison between models, we select the typical hyperparameter setting for each efficient attention on each task, which is shown in Table 1 in detail. Specifically, we first apply our hybrid architecture NAR-AMLP in two tasks: Text-to-Speech Synthesis and Machine Translation. Then we assess AMLP’s self-attention and cross-attention abilities on super resolution and long sequence time-series forecasting tasks, respectively. Finally, we conduct ablation studies to show the hidden philosophy of AMLP and explore how efficient AMLP scales to long-sequence modeling.

⁶ In experiments, we take $\text{softmax}(\cdot)$ as the nonlinear function $\sigma_1(\cdot)$ unless otherwise specified.

Table 1: Hyperparameters of different tasks.

Task	TTS	MT	SR	LSTF
Backbone	FastSpeech 2/ Transformer-TTS	Transformer/CMLMC	SR	Informer
<i>Training hyperparameters</i>				
Batch Size	48	–	4	32
Number of Steps (epochs)	20K	100K/300K	1M	6 (epochs)
Warmup Steps	4K	4K	–	–
Peak Learning Rate	5e-4	5e-4	1e-4	1e-4
Scheduler	Inverse Sqrt	Inverse Sqrt	Linear	Exponential Decay
Optimizer	AdamW	AdamW	AdamW	AdamW
Adam	(0.9, 0.98)	(0.9, 0.98)	(0.9, 0.999)	(0.9, 0.999)
Clip Norm	5.0	5.0	0	0
Attention Dropout	0.1	0.3	0.2	0.05
Weight Decay	0.01	0.0001	0	0
Max Tokens	–	65536	–	–
Iteration	–	–	–	5
Evaluation Checkpoint	best	average last 10	average last 5	last
<i>Attention hyperparameters</i>				
wsize (local)	15	5	15	15
landmarks (ABC)	64	16	64	64
ffn_dim (AMLP)	64	16	64	64
approx_dim (Performer)	64	16	64	64

3.1 Main Results of NAR-AMLP

Text-to-Speech We select LJSpeech [25] dataset for this task, and use FastSpeech 2 (FS2) [46] and Transformer-TTS (Tr-TTS) [30] as the backbone models for NAR and AR, respectively. For both backbones, we replace all softmax attentions with efficient ones to achieve linear complexity. We use AMLP-Cov variant and $\text{ReLU}(\cdot)$ as $\sigma_1(\cdot)$ in Eq. 8. The alignment tool “g2pE” [54] is applied to train FastSpeech 2. For reproducibility, we use two widely-used objective evaluation metrics, Mel Cepstral Distortion (MCD) and Mel Spectral Distortion (MSD), to assess the quality of synthesized audio clips. We compare AMLP with gMLP [33], XCA [1], ABC [40] and local attention [36]. The details of training hyperparameters are shown in Table 1. We demonstrate the results in Table 2. AMLP substantially lowers the MCD and MSD values by a great margin up to 0.15 MCD with even lower complexity compared to vanilla models. Additionally, AMLP also outperforms other efficient models. Notably, we have significantly lower MCD than XCA which also leverages (cross-)covariance matrices.

Machine Translation To verify AMLP’s capability on short sequence modeling, we launch Machine Translation (MT) experiments on WMT 2014 English-German (WMT’14 En-De) and German-English (WMT’14 De-En) datasets [6]. We adopt AMLP-PQuery variant to CMLMC [23], which is a powerful fully NAR architecture without extra decoding algorithms. For completeness, we include widely-used AR architecture Transformer (Tr) [53] with competitive linear atten-

Table 2: Automatic evaluation metric on LJSpeech dataset. All models are trained by ourselves. n, m are the target and source sequence lengths. Colored rows represent NAR models.

Arch Model	#Params	LJSpeech	
		MCD↓	MSD↓
<i>Complexity: $O(n^2)$ or $O(n^2 + nm + m^2)$</i>			
AR	Tr-TTS	54.40M	4.095 2.199
NAR	FS2	41.23M	3.475 1.974
<i>Complexity: $O(n)$ or $O(n + m)$</i>			
AR	Tr-TTS (ABC)	54.60M	5.130 2.596
	FS2 (local)	41.23M	3.419 1.970
	FS2 (ABC)	41.36M	3.392 1.966
NAR	FS2 (XCA)	41.23M	3.500 2.024
	FS2 (gMLP)	44.90M	3.402 1.964
	FS2 (AMLP)	41.49M	3.327 1.940

Table 3: BLEU4 scores on WMT14 EN-DE and WMT14 DE-EN dataset. All models for comparison are implemented by ourselves. n, m are the target and source sequence lengths. Colored rows represent NAR models.

Arch Model	#Params	WMT' 14	
		En-De	De-En
<i>Complexity: $O(n^2 + nm + m^2)$</i>			
AR	Tr	86.74M	27.38 31.26
NAR	CMLMC	73.14M	27.91 31.43
<i>Complexity: $O(n + m)$</i>			
	Tr (local)	86.74M	24.77 28.21
AR	Tr (ABC)	86.77M	25.86 29.09
	CMLMC (ABC)	73.16M	27.37 31.30
NAR	CMLMC (local)	73.16M	27.05 30.33
	CMLMC (AMLP)	73.44M	27.60 31.50

tions. We exclude the AR-reranking process to make a fully linear-complexity generation process. Similar to TTS, we replace self/cross-attention modules in the decoder of Transformer and CMLMC to obtain their efficient variants. We use hyperparameters as CMLMC and Transformer suggest, which is present in Table 1. We report BLEU-4 [39] scores as the performance metric. Because XCA and gMLP do not support cross-attention, we here only compare AMLP with the strong ABC and local baselines. As translation has implicit token alignment between sequences, local attention can do cross-attention in this task.

Results in Table 3 indicate that the NAR-AMLP architecture achieves the best result among efficient NAR and AR models with linear complexity. Among the NAR models, the AMLP model outperforms a strong linear attention model, ABC, on both datasets, with a lead of 0.23 and 0.20 BLEU, respectively. It indicates that AMLP effectively captures short-term dependencies and produces more accurate translations than ABC. We also compare AMLP with vanilla attention, and the results indicate that AMLP outperforms vanilla attention on the de-en dataset, with only a 0.31BLEU lag compared to vanilla attention on the en-de one. This suggests that AMLP can achieve comparable performance to vanilla NAR models in certain scenarios. In comparison to AR models on both datasets, AMLP demonstrates superior performance (with at least 0.22 and 0.24 BLEU improvement), providing further evidence of the efficacy of NAR-AMLP as an architecture.

3.2 Self- and Cross-Attention Ablation

Self-attention We evaluate the self-encoding ability of AMLP on Super Resolution (SR) task. SR aims to convert low-resolution (16×16) images into high-resolution (128×128) ones. We base on a powerful backbone — SR3 [49]

Table 5: Cross-attention ablation on ETT-h1, ETT-h1, and ETT-m1 datasets. n, m are the target and source lengths. Avg. m is computed over three subdatasets.

Complex.	Methods	#Params	ETT <h>1</h>		ETT <h>2</h>		ETT <h>m</h>		Avg.	
			MSE↓	MAE↓	MSE↓	MAE↓	MSE↓	MAE↓	MSE↓	MAE↓
$O(n^2 + nm)$	vanilla	11.33M	0.754	0.573	1.907	1.036	0.754	0.716	1.138	0.775
	ABC	11.33M	0.845	0.728	1.862	1.013	0.734	0.685	1.147	0.809
	Performer	11.33M	0.861	0.703	2.137	1.091	0.764	0.663	1.254	0.819
$O(n + m)$	cosFormer	11.33M	0.848	0.723	2.094	1.067	0.715	0.680	1.219	0.823
	AMLP	11.33M	0.797	0.702	1.504	0.864	0.718	0.684	1.006	0.750

and add attention layers after each residual block to follow CAB [59] settings. We replace the softmax self-attention with five efficient architectures, i.e., local, gMLP, XCA, ABC and AMLP to compare. Following [49], we use the Flickr-Faces-HQ (FFHQ) dataset [27] for the training set and CelebA-HQ dataset [26] for the evaluation set. We use Peak Signal-to-Noise Ratio (PSNR) and Structural SIMilarity (SSIM) [55] to measure efficient models. Experiment results are shown in Table 4. AMLP improves the performance of SR3 to 23.28 (+0.10) on PSNR and 0.684 (+0.09) on SSIM against the vanilla baseline, indicating that AMLP has a strong self-encoding ability. When compared to gMLP, AMLP also has a slight performance gain. AMLP outperforms covariance-based architecture XCA by 0.20 and 0.14 on PSNR and SSIM, respectively.

Cross-Attention We test the cross-attention ability on the long sequence time-series forecasting (LSTF) task. We take Informer [61] as the backbone neural networks and evaluate efficient models on Electricity Transformer Temperature (ETT) dataset, which contains three sub-datasets ETT-h1, ETT-h2, and ETT-m1. We follow [61] to conduct univariate and multivariate evaluations on three sub-datasets and average their Mean Square Error (MSE) and Mean Absolute Error (MAE) to obtain final scores. Except for vanilla attention, we also compare AMLP with other three efficient models with strong cross-alignment abilities: ABC [40], Performer [10] and cosFormer [44]. We exclude local attention as it does not work for cross attention without explicit token alignment in the time-series forecasting task. The results performed on three sub-datasets are shown in Table 5. AMLP, in contrast to the vanilla counterpart, achieves lower MSE and MAE as well as more efficient complexity. Moreover, we notice that all other efficient models perform poorly compared to vanilla attention. It suggests that AMLP has a solid ability to model non-homologous information.

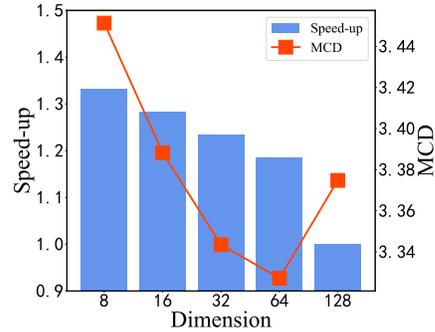
3.3 Analysis

In this section, we conduct substantial analysis experiments to dig out the efficiency and superiority of our AMLP mechanism. We first present our analysis in comparison with other efficient attention modules on the TTS task. Then we

Table 4: PSNR and SSIM on CelebA-HQ dataset. n is the pixel number of the images.

Celeb-HQ			
Model	#Params	PSNR \uparrow	SSMI \uparrow
<i>Complexity: $O(n^2)$</i>			
vanilla	99.55M	23.18	0.675
<i>Complexity: $O(n)$</i>			
local	99.55M	23.33	0.682
gMLP	101.66M	23.24	0.679
XCiT	99.55M	23.08	0.67
ABC	99.72M	22.54	0.635
AMLP	99.73M	23.28	0.684

Fig. 3: Trade-off of MCD value and speed-up of different intermediate dimension c values in text-to-speech task.



show that our approximation $c < d$ in Eq. 6 does not deteriorate the performance of speech generation. Finally, we elucidate the outstanding generation speed and GPU peak usage of our AMLP in the NAR scenario.

Intermediate Dimension Analysis The approximation of eigenvalues in Eq. 6 prompts us to know whether such approximation is feasible and whether the exorbitant approximation will deteriorate the generation performance. To this end, we test several values of c in AMLP and report each corresponding performance on TTS and the decoding speed when adopted to FastSpeech 2, in Fig. 3. Except for c value, we adopt the same setting in §3.1.

From Fig. 3, we can see that AMLP with approximation rank c can achieve as well as no approximation setting ($c = d = 128$) and does not impact the performance greatly. But with a lower c value, AMLP can achieve better decoding speed. Specifically, in contrast to $c = 64$, a higher MCD when setting c to d also indicates that maintaining the whole eigenvalues in Eq. 6 may even lead to overparameterization and impair the overall decoding efficacy. It verifies the feasibility to approximate Σ with fewer eigenspectrums in AMLP.

Efficiency Analysis To further understand the performance of NAR-AMLP architecture in inference, we set up a simulation experiment to test its efficiency. The simulation experiment evaluates NAR-AMLP efficiency from running time and memory usage with respect to sequence length from 256 to 8,192, compared with AR model and vanilla NAR model. We simulate the generation process with a single efficient module. For AR, we test its causal attention, which is its bottleneck in generation. For AMLP, we use 64 as the inner dimension with ReLU activation function for σ_1 in Eq. 8. AMLP-Cov and AMLP-PQuery shares the same complexity, so we use “AMLP” to denote the two variants. The experiments

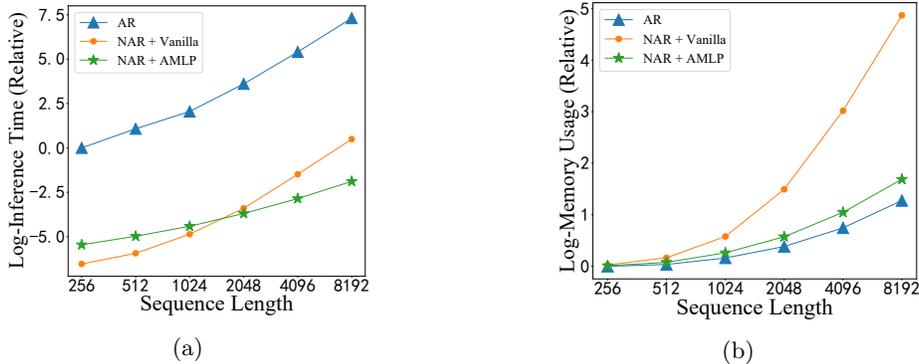


Fig. 4: (a) Empirical running time and (b) empirical memory cost with sequence length. Logarithms of relative measurement to the AR model are reported.

are performed with batch size 12 on a single A100 GPU, and the results are repeated with 100 runs. We remain running latency data ranging from the first quartile and the third quartile among the 100 runs to remove noise. Finally, the remaining figures are averaged to serve as the final time consumption.

Fig. 4a shows that NAR-AMLP extremely speeds up the inference process. To generate a long sequence with 8,192 tokens, vanilla NAR is $116\times$ faster than AR while NAR-AMLP is even $590\times$ faster. For sequences with more than 1500 tokens, both variants of AMLP are more efficient than vanilla attention; otherwise, the vanilla attention is faster. Fig. 4b shows that NAR-AMLP significantly reduces memory consumption in NAR generation. It saves 89% memory usage of NAR model when generating a sequence with 8,192 tokens. Note that AR models cost fewer memory resources because of incremental decoding, which caches previous states and processes only one token at each step. But AR models still suffer from huge memory usage as NAR models in training, since they are usually implemented with a causal mask on the attention matrix. Thus it is reasonable to infer that NAR-AMLP is more efficient than AR and NAR models in training.

4 Related Work

Non-Autoregressive Generation [17] first proposes a non-autoregressive model to generate all the tokens within a sequence in parallel, which extremely speeds up the inference process but is inferior in generation quality. To mitigate the quality degradation, many researchers devote to improve the model performance with iterative decoding [29,16,18,20,22], curriculum learning [19,34,42,43,4], latent variable modeling [38,45,3,4], imitation learning [31,56] and learning objective [48,15,32,12]. These previous works focus on pursuing the high efficacy of non-autoregressive generation, but few works are presented to improve NAR’s efficiency in long sequence modeling. We target to further improve its efficiency and scale non-autoregressive models to long sequences.

MLP Architecture Multi-layer perceptron [14] is a classic neural network architecture and has been widely used. Recently, novel variants of MLP architectures are proposed for text and image processing, achieving impressive results on image classification [52,33], text classification [50], multilingual parsing [13], and intent classification [13]. MLP-Mixer [52] is proposed by leveraging a token-mixing and a channel-mixing MLP to enable token-wise and channel-wise communication. MLP-Mixer is further improved to pNLP-Mixer with locality sensitive hashing [24] projection at the bottom calculating non-trainable fingerprints [13]. [33] propose gMLP by introducing a spatial gating unit to enhance the communication between neighboring tokens. CycleMLP [8] leverages a local window to achieve linear time complexity on dense prediction. Besides, previous studies focus on encoding text/image features with MLP, but we explore the possibility to leverage an MLP architecture for sequence generation.

Attention Mechanism Attention is first proposed to align the target and source sequence in neural machine translation [2], and is further improved to multi-head self/cross/causal attention [53]. Due to its quadratic time complexity and memory cost with sequence length, a surge of efficient attention is proposed to improve the efficiency of softmax attention. Due to the sparsity of attention matrix, many researchers propose to explicitly model a sparse attention mechanism to obtain fast computation without harming performance [21,51,28,5,58,47]. The low-rank property of attention matrix also brings out matrix decomposition-based methods [57,9]. The softmax attention can also be linearized via exponential kernel decomposition [10,40,41,60,44]. These attention variants are exploring an efficient way to approximate softmax attention, but we focus on MLP architecture, which is naturally an efficient architecture.

5 Conclusions

In this work, we introduced Attentive Multi-Layer Perceptron (AMLP), an efficient plugin alternative to vanilla attention for non-autoregressive generation tasks. AMLP uses adaptive weights to learn inter-token interactions as done in attention. And we also put forward two methods adopting different philosophies to parameterize the adaptive weight matrices in AMLP. Substantial experiments on generation tasks verify that AMLP surpasses attention in most tasks and achieves similar performances with other strong efficient models in other tasks. Besides, efficiency analysis indicates that AMLP combined NAR model could save time compared to AR models, and save space compared to vanilla NAR models in long sequence settings.

6 Ethical Issues

AMLP is designed to speed up the generation of non-autoregressive models, by replacing the inefficient softmax attention with our AMLP module to achieve

linear complexity. The potential positive implications imply lower difficulty in deploying NAR models on resource-limited devices, thus increasing the accessibility of NAR models. AMLP also makes positive impacts on extending NAR models to various domains, since it can do both self-attention and cross-attention. Moreover, the high efficiency of AMLP reduces the carbon footprint of training a model and thus brings positive environmental benefits. As such, we do not foresee any immediate negative ethical or societal consequences stemming from our work that are different from those that apply to other fundamental components of the transformer architecture and NAR models.

References

1. Ali, A., Touvron, H., Caron, M., Bojanowski, P., Douze, M., Joulin, A., Laptev, I., Neverova, N., Synnaeve, G., Verbeek, J., et al.: Xcit: Cross-covariance image transformers. *Advances in neural information processing systems* **34**, 20014–20027 (2021)
2. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. In: *ICLR* (2015)
3. Bao, Y., Zhou, H., Feng, J., Wang, M., Huang, S., Chen, J., Li, L.: Non-autoregressive transformer by position learning. *arXiv preprint arXiv:1911.10677* (2019)
4. Bao, Y., Zhou, H., Huang, S., Wang, D., Qian, L., Dai, X., Chen, J., Li, L.: latent-GLAT: Glancing at latent variables for parallel text generation. In: *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. pp. 8398–8409. Association for Computational Linguistics, Dublin, Ireland (May 2022). <https://doi.org/10.18653/v1/2022.acl-long.575>, <https://aclanthology.org/2022.acl-long.575>
5. Beltagy, I., Peters, M.E., Cohan, A.: Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150* (2020)
6. Bojar, O., Buck, C., Federmann, C., Haddow, B., Koehn, P., Leveling, J., Monz, C., Pecina, P., Post, M., Saint-Amand, H., Soricut, R., Specia, L., Tamchyna, A.: Findings of the 2014 workshop on statistical machine translation. In: *Proceedings of the Ninth Workshop on Statistical Machine Translation*. pp. 12–58. Association for Computational Linguistics, Baltimore, Maryland, USA (Jun 2014). <https://doi.org/10.3115/v1/W14-3302>
7. Chang, H., Zhang, H., Jiang, L., Liu, C., Freeman, W.T.: Maskgit: Masked generative image transformer. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 11315–11325 (2022)
8. Chen, S., Xie, E., GE, C., Chen, R., Liang, D., Luo, P.: CycleMLP: A MLP-like architecture for dense prediction. In: *International Conference on Learning Representations* (2022), <https://openreview.net/forum?id=NMEceG4v69Y>
9. Chen, Z., Gong, M., Ge, L., Du, B.: Compressed self-attention for deep metric learning with low-rank approximation. In: *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*. pp. 2058–2064 (2021)
10. Choromanski, K.M., Likhoshervostov, V., Dohan, D., Song, X., Gane, A., Sarlos, T., Hawkins, P., Davis, J.Q., Mohiuddin, A., Kaiser, L., Belanger, D.B., Colwell, L.J., Weller, A.: Rethinking attention with performers. In: *International Conference on Learning Representations* (2021), <https://openreview.net/forum?id=Ua6zuk0WRH>

11. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: International Conference on Learning Representations (2021)
12. Du, C., Tu, Z., Jiang, J.: Order-agnostic cross entropy for non-autoregressive machine translation. In: International Conference on Machine Learning. pp. 2849–2859. PMLR (2021)
13. Fusco, F., Pascual, D., Staar, P.: pnlp-mixer: an efficient all-mlp architecture for language. arXiv preprint arXiv:2202.04350 (2022)
14. Gardner, M.W., Dorling, S.: Artificial neural networks (the multilayer perceptron)—a review of applications in the atmospheric sciences. *Atmospheric environment* **32**(14-15), 2627–2636 (1998)
15. Ghazvininejad, M., Karpukhin, V., Zettlemoyer, L., Levy, O.: Aligned cross entropy for non-autoregressive machine translation. In: International Conference on Machine Learning. pp. 3515–3523. PMLR (2020)
16. Ghazvininejad, M., Levy, O., Liu, Y., Zettlemoyer, L.: Mask-predict: Parallel decoding of conditional masked language models. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 6112–6121 (2019)
17. Gu, J., Bradbury, J., Xiong, C., Li, V.O., Socher, R.: Non-autoregressive neural machine translation. In: International Conference on Learning Representations (2018)
18. Gu, J., Wang, C., Zhao, J.: Levenshtein transformer. *Advances in Neural Information Processing Systems* **32** (2019)
19. Guo, J., Tan, X., Xu, L., Qin, T., Chen, E., Liu, T.Y.: Fine-tuning by curriculum learning for non-autoregressive neural machine translation. *Proceedings of the AAAI Conference on Artificial Intelligence* **34**(05), 7839–7846 (Apr 2020). <https://doi.org/10.1609/aaai.v34i05.6289>, <https://ojs.aaai.org/index.php/AAAI/article/view/6289>
20. Guo, J., Xu, L., Chen, E.: Jointly masked sequence-to-sequence model for non-autoregressive neural machine translation. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 376–385 (2020)
21. Ho, J., Kalchbrenner, N., Weissenborn, D., Salimans, T.: Axial attention in multi-dimensional transformers. arXiv preprint arXiv:1912.12180 (2019)
22. Huang, C., Zhou, H., Zaïane, O.R., Mou, L., Li, L.: Non-autoregressive translation with layer-wise prediction and deep supervision. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 10776–10784 (2022)
23. Huang, X.S., Perez, F., Volkovs, M.: Improving non-autoregressive translation models without distillation. In: International Conference on Learning Representations (2022)
24. Indyk, P., Motwani, R.: Approximate nearest neighbors: towards removing the curse of dimensionality. In: Proceedings of the thirtieth annual ACM symposium on Theory of computing. pp. 604–613 (1998)
25. Ito, K., Johnson, L.: The lj speech dataset. <https://keithito.com/LJ-Speech-Dataset/> (2017)
26. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. In: International Conference on Learning Representations (2018)

27. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
28. Kitaev, N., Kaiser, L., Levskaya, A.: Reformer: The efficient transformer. In: International Conference on Learning Representations (2020), <https://openreview.net/forum?id=rkgNkKHtVb>
29. Lee, J., Mansimov, E., Cho, K.: Deterministic non-autoregressive neural sequence modeling by iterative refinement. In: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. pp. 1173–1182 (2018)
30. Li, N., Liu, S., Liu, Y., Zhao, S., Liu, M.: Neural speech synthesis with transformer network. Proceedings of the AAAI Conference on Artificial Intelligence **33**(01), 6706–6713 (Jul 2019). <https://doi.org/10.1609/aaai.v33i01.33016706>, <https://ojs.aaai.org/index.php/AAAI/article/view/4642>
31. Li, Z., Lin, Z., He, D., Tian, F., Qin, T., Wang, L., Liu, T.Y.: Hint-based training for non-autoregressive machine translation. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 5708–5713. Association for Computational Linguistics, Hong Kong, China (Nov 2019). <https://doi.org/10.18653/v1/D19-1573>, <https://aclanthology.org/D19-1573>
32. Liu, G., Yang, Z., Tao, T., Liang, X., Bao, J., Li, Z., He, X., Cui, S., Hu, Z.: Don't take it literally: An edit-invariant sequence loss for text generation. In: Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 2055–2078. Association for Computational Linguistics, Seattle, United States (Jul 2022). <https://doi.org/10.18653/v1/2022.naacl-main.150>, <https://aclanthology.org/2022.naacl-main.150>
33. Liu, H., Dai, Z., So, D., Le, Q.V.: Pay attention to mlps. Advances in Neural Information Processing Systems **34**, 9204–9215 (2021)
34. Liu, J., Ren, Y., Tan, X., Zhang, C., Qin, T., Zhao, Z., Liu, T.Y.: Task-level curriculum learning for non-autoregressive neural machine translation. In: IJCAI. pp. 3861–3867 (2020)
35. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10012–10022 (2021)
36. Luong, T., Pham, H., Manning, C.D.: Effective approaches to attention-based neural machine translation. In: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. pp. 1412–1421. Association for Computational Linguistics, Lisbon, Portugal (Sep 2015). <https://doi.org/10.18653/v1/D15-1166>, <https://aclanthology.org/D15-1166>
37. Ma, X., Zhou, C., Kong, X., He, J., Gui, L., Neubig, G., May, J., Zettlemoyer, L.: Mega: moving average equipped gated attention. arXiv preprint arXiv:2209.10655 (2022)
38. Ma, X., Zhou, C., Li, X., Neubig, G., Hovy, E.: FlowSeq: Non-autoregressive conditional sequence generation with generative flow. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 4282–4292. Association for Computational Linguistics, Hong Kong, China (Nov 2019). <https://doi.org/10.18653/v1/D19-1437>

39. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting of the Association for Computational Linguistics. pp. 311–318 (2002)
40. Peng, H., Kasai, J., Pappas, N., Yogatama, D., Wu, Z., Kong, L., Schwartz, R., Smith, N.A.: ABC: Attention with bounded-memory control. In: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 7469–7483. Association for Computational Linguistics, Dublin, Ireland (May 2022). <https://doi.org/10.18653/v1/2022.acl-long.515>, <https://aclanthology.org/2022.acl-long.515>
41. Peng, H., Pappas, N., Yogatama, D., Schwartz, R., Smith, N., Kong, L.: Random feature attention. In: International Conference on Learning Representations (2021)
42. Qian, L., Zhou, H., Bao, Y., Wang, M., Qiu, L., Zhang, W., Yu, Y., Li, L.: Glancing transformer for non-autoregressive neural machine translation. In: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 1993–2003. Association for Computational Linguistics, Online (Aug 2021). <https://doi.org/10.18653/v1/2021.acl-long.155>
43. Qian, L., Zhou, Y., Zheng, Z., Zhu, Y., Lin, Z., Feng, J., Cheng, S., Li, L., Wang, M., Zhou, H.: The volctrans GLAT system: Non-autoregressive translation meets WMT21. In: Proceedings of the Sixth Conference on Machine Translation. pp. 187–196. Association for Computational Linguistics, Online (Nov 2021)
44. Qin, Z., Sun, W., Deng, H., Li, D., Wei, Y., Lv, B., Yan, J., Kong, L., Zhong, Y.: cosformer: Rethinking softmax in attention. In: International Conference on Learning Representations (2022), <https://openreview.net/forum?id=B18CQrx2Up4>
45. Ran, Q., Lin, Y., Li, P., Zhou, J.: Guiding non-autoregressive neural machine translation decoding with reordering information. In: AAAI. pp. 13727–13735 (2021)
46. Ren, Y., Hu, C., Tan, X., Qin, T., Zhao, S., Zhao, Z., Liu, T.Y.: Fastspeech 2: Fast and high-quality end-to-end text to speech. In: International Conference on Learning Representations (2021)
47. Roy, A., Saffar, M., Vaswani, A., Grangier, D.: Efficient content-based sparse attention with routing transformers. *Transactions of the Association for Computational Linguistics* **9**, 53–68 (2021)
48. Saharia, C., Chan, W., Saxena, S., Norouzi, M.: Non-autoregressive machine translation with latent alignments. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 1098–1108. Association for Computational Linguistics, Online (Nov 2020). <https://doi.org/10.18653/v1/2020.emnlp-main.83>
49. Saharia, C., Ho, J., Chan, W., Salimans, T., Fleet, D.J., Norouzi, M.: Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022)
50. Tay, Y., Bahri, D., Metzler, D., Juan, D.C., Zhao, Z., Zheng, C.: Synthesizer: Rethinking self-attention for transformer models. In: International Conference on Machine Learning. pp. 10183–10192. PMLR (2021)
51. Tay, Y., Bahri, D., Yang, L., Metzler, D., Juan, D.C.: Sparse sinkhorn attention. In: International Conference on Machine Learning. pp. 9438–9447. PMLR (2020)
52. Tolstikhin, I.O., Houlsby, N., Kolesnikov, A., Beyer, L., Zhai, X., Unterthiner, T., Yung, J., Steiner, A., Keysers, D., Uszkoreit, J., et al.: Mlp-mixer: An all-mlp architecture for vision. *Advances in Neural Information Processing Systems* **34** (2021)

53. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
54. Wang, C., Hsu, W.N., Adi, Y., Polyak, A., Lee, A., Chen, P.J., Gu, J., Pino, J.: fairseq s²: A scalable and integrable speech synthesis toolkit. In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. pp. 143–152. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic (Nov 2021). <https://doi.org/10.18653/v1/2021.emnlp-demo.17>
55. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* **13**(4), 600–612 (2004)
56. Wei, B., Wang, M., Zhou, H., Lin, J., Sun, X.: Imitation learning for non-autoregressive neural machine translation. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. pp. 1304–1312. Association for Computational Linguistics, Florence, Italy (Jul 2019). <https://doi.org/10.18653/v1/P19-1125>
57. Xiong, Y., Zeng, Z., Chakraborty, R., Tan, M., Fung, G., Li, Y., Singh, V.: Nyströmformer: A nyström-based algorithm for approximating self-attention. *Proceedings of the AAAI Conference on Artificial Intelligence* **35**(16), 14138–14148 (May 2021). <https://doi.org/10.1609/aaai.v35i16.17664>, <https://ojs.aaai.org/index.php/AAAI/article/view/17664>
58. Zaheer, M., Guruganesh, G., Dubey, K.A., Ainslie, J., Alberti, C., Ontanon, S., Pham, P., Ravula, A., Wang, Q., Yang, L., et al.: Big bird: Transformers for longer sequences. *Advances in Neural Information Processing Systems* **33**, 17283–17297 (2020)
59. Zhang, J., Jiang, S., Feng, J., Zheng, L., Kong, L.: Cab: Comprehensive attention benchmarking on long sequence modeling. *arXiv preprint arXiv:2210.07661* (2022)
60. Zheng, L., Wang, C., Kong, L.: Linear complexity randomized self-attention mechanism. *arXiv preprint arXiv:2204.04667* (2022)
61. Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., Zhang, W.: Informer: Beyond efficient transformer for long sequence time-series forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence* **35**(12), 11106–11115 (May 2021), <https://ojs.aaai.org/index.php/AAAI/article/view/17325>