
FourierFormer: Transformer Meets Generalized Fourier Integral Theorem

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Multi-head attention empowers the recent success of transformers, the state-of-the-
2 art models that have achieved remarkable success in sequence modeling and beyond.
3 These attention mechanisms compute the pairwise dot products between the queries
4 and keys, which results from the use of unnormalized Gaussian kernels with the
5 assumption that the queries follow a mixture of Gaussian distribution. There is no
6 guarantee that this assumption is valid in practice. In response, we first interpret
7 attention in transformers as a nonparametric kernel regression. We then propose
8 the FourierFormer, a new class of transformers in which the dot-product kernels
9 are replaced by the novel generalized Fourier integral kernels. Different from the
10 dot-product kernels, where we need to choose a good covariance matrix to capture
11 the dependency of the features of data, the generalized Fourier integral kernels can
12 automatically capture such dependency and remove the need to tune the covariance
13 matrix. We theoretically prove that our proposed Fourier integral kernels can effi-
14 ciently approximate any key and query distributions. Compared to the conventional
15 transformers with dot-product attention, FourierFormers attain better accuracy
16 and reduce the redundancy between attention heads. We empirically corroborate
17 the advantages of FourierFormers over the baseline transformers in a variety of
18 practical applications including language modeling and image classification.

19 1 Introduction

20 Transformers [76] are powerful neural networks that have achieved tremendous success in many
21 areas of machine learning [38, 69, 34] and become the state-of-the-art model on a wide range
22 of applications across different data modalities, from language [22, 1, 17, 12, 55, 4, 8, 20] to
23 images [23, 41, 71, 56, 52, 26], videos [3, 42], point clouds [90, 29], and protein sequence [58, 32].
24 In addition to their excellent performance on supervised learning tasks, transformers can also
25 effectively transfer the learned knowledge from a pretraining task to new tasks with limited or no
26 supervision [53, 54, 22, 87, 40]. At the core of transformers is the dot-product self-attention, which
27 mainly accounts for the success of transformer models [13, 49, 39]. This dot-product self-attention
28 learn self-alignment between tokens in an input sequence by estimating the relative importance of a
29 given token with respect to all other tokens. It then transform each token into a weighted average of
30 the feature representations of other tokens where the weight is proportional to a importance score
31 between each pair of tokens. The importance scores in self-attention enable a token to attend to other
32 tokens in the sequence, thus capturing the contextual representation [6, 76, 36].

33 1.1 Self-Attention

34 Given an input sequence $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times D_x}$ of N feature vectors, self-attention
35 computes the output sequence \mathbf{H} from \mathbf{X} as follows:

Step 1: Projecting the input sequence into different subspaces. The input sequence \mathbf{X} is transformed into the query matrix \mathbf{Q} , the key matrix \mathbf{K} , and the value matrix \mathbf{V} via three linear
Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

transformations

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q^\top; \mathbf{K} = \mathbf{X}\mathbf{W}_K^\top; \mathbf{V} = \mathbf{X}\mathbf{W}_V^\top,$$

36 where $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{D \times D_x}$, and $\mathbf{W}_V \in \mathbb{R}^{D_v \times D_x}$ are the weight matrices. We denote $\mathbf{Q} :=$
 37 $[\mathbf{q}_1, \dots, \mathbf{q}_N]^\top$, $\mathbf{K} := [\mathbf{k}_1, \dots, \mathbf{k}_N]^\top$, and $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_N]^\top$, where the vectors $\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i$ for
 38 $i = 1, \dots, N$ are the query, key, and value vectors, respectively.

39 **Step 2: Computing the output as a weighted average.** The output sequence $\mathbf{H} := [\mathbf{h}_1, \dots, \mathbf{h}_N]^\top$
 40 is then given by

$$\mathbf{H} = \text{softmax}\left(\mathbf{Q}\mathbf{K}^\top / \sqrt{D}\right)\mathbf{V} := \mathbf{A}\mathbf{V}, \quad (1)$$

41 where the softmax function is applied to each row of the matrix $(\mathbf{Q}\mathbf{K}^\top) / \sqrt{D}$. For each query vector
 42 $\mathbf{q}_i, i = 1, \dots, N$, Eqn. (1) can be written in the vector form to compute the output vector \mathbf{h}_i as
 43 follows

$$\mathbf{h}_i = \sum_{j=1}^N \text{softmax}\left(\mathbf{q}_i^\top \mathbf{k}_j / \sqrt{D}\right) \mathbf{v}_j := \sum_{j=1}^N a_{ij} \mathbf{v}_j. \quad (2)$$

44 The matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ and its component a_{ij} for $i, j = 1, \dots, N$ are the attention matrix and
 45 attention scores, respectively. The self-attention computed by equations (1) and (2) is called the dot-
 46 product attention or softmax attention. In our paper, we refer a transformer that uses this attention as
 47 the baseline transformer with the dot-product attention or the dot-product transformer. The structure
 48 of the attention matrix \mathbf{A} after training governs the ability of the self-attention to capture contextual
 49 representation for each token.

50 **Multi-head Attention** Each output sequence \mathbf{H} forms an attention head. Multi-head attention
 51 concatenates multiple heads to compute the final output. Let H be the number of heads and
 52 $\mathbf{W}^O \in \mathbb{R}^{HD_v \times HD_v}$ be the projection matrix for the output. The multi-head attention is defined as

$$\text{MultiHead}(\{\mathbf{Q}, \mathbf{K}, \mathbf{V}\}_{i=1}^H) = \text{Concat}(\mathbf{H}_1, \dots, \mathbf{H}_H) \mathbf{W}^O.$$

53 The capacity of the attention mechanism and its ability to learn diverse syntactic and semantic
 54 relationships determine the success of transformers [70, 77, 16, 78, 30]. However, equations (1)
 55 and (2) implies that the dot-product attention assumes the features (q_{i1}, \dots, q_{iD}) in \mathbf{q}_i , as well as
 56 the features (k_{j1}, \dots, k_{jD}) in \mathbf{k}_j , are independent. Thus, the dot-product attention fail to capture the
 57 correlations between these features, limiting its representation capacity and inhibit the performance
 58 of transformers on practical tasks where there is no guarantee that independent features can learned
 59 from complex data. One solution to capture correlations between features \mathbf{q}_i and \mathbf{k}_j is to introduce
 60 covariance matrices into the formulation of the dot-product attention with the cost of significantly
 61 increasing of the computational complexity. Also, choosing good covariance matrices is difficult.

62 1.2 Contribution

63 In this paper, we first establish a correspondence between self-attention and nonparametric kernel
 64 regression. Under this new perspective of self-attention, we explain the limitation of the dot-product
 65 self-attention that it may fail to capture correlations between the features in the query and key
 66 vectors. We then leverage the generalized Fourier integral theorems, which can automatically capture
 67 these correlations, and derive the generalized Fourier integral estimators for the nonparametric
 68 regression problem. Using this new density estimator, we propose the FourierFormer, a novel
 69 class of transformers that can capture correlations between features in the query and key vectors of
 70 self-attention. In summary, our contribution is three-fold:

- 71 1. We derive the formula of self-attention from solving a nonparametric kernel regression
 72 problem, thus providing a nonparametric regression interpretation to study and further
 73 develop self-attention.
- 74 2. We develop the generalized Fourier integral estimators for the nonparametric regression
 75 problem and provide theoretical guarantees for these estimator.
- 76 3. We propose the FourierFormer whose attentions use the generalized Fourier integral es-
 77 timators to capture more efficiently correlations between features in the query and key
 78 vectors.

79 Finally, we empirically show that the FourierFormer attains significantly better accuracy than the
 80 baseline transformer with the dot-product attention on a variety of tasks including the WikiText
 81 language modeling and ImageNet image classification. We also demonstrate in our experiments that
 82 FourierFormer helps reduce the redundancy between attention heads.

83 **Organization** We structure this paper as follows: In Section 2, we present the correspondence
 84 between self-attention and nonparametric kernel regression. In Section 3, we discuss the generalized
 85 Fourier integral estimators and define the FourierFormer. We validate and empirically analyze the
 86 advantages of FourierFormer in Section 4. We discuss related works in Section 5. The paper ends with
 87 concluding remarks. Technical proofs and more experimental details are provided in the Appendix.

88 **Notation** For any $N \in \mathbb{N}$, we denote $[N] = \{1, 2, \dots, N\}$. For any $D \geq 1$, $\mathbb{L}_1(\mathbb{R}^D)$ denotes the
 89 space of real-valued functions on \mathbb{R}^D that are integrable. For any two sequences $\{a_N\}_{N \geq 1}$, $\{b_N\}_{N \geq 1}$,
 90 we denote $a_N = \mathcal{O}(b_N)$ to mean that $a_N \leq Cb_N$ for all $N \geq 1$ where C is some universal constant.

91 2 A Nonparametric Regression Interpretation of Self-attention

92 In this section, we establish the connection between self-attention and nonparametric kernel regression.
 93 In particular, we derive the self-attention in equation (2) as a nonparametric kernel regression in
 94 which the key vectors \mathbf{k}_j and value vectors \mathbf{v}_j are training inputs and training targets, respectively,
 95 while the query vectors \mathbf{q}_i and the output vectors \mathbf{h}_i form a set of new inputs and their corresponding
 96 targets that need to be estimated, respectively, for $i, j = 1, \dots, N$. In general, we can view the
 97 training set $\{\mathbf{k}_j, \mathbf{v}_j\}$ for $j \in [N]$ to come from the following *nonparametric regression model*:

$$\mathbf{v}_j = f(\mathbf{k}_j) + \varepsilon_j, \quad (3)$$

98 where $\varepsilon_1, \dots, \varepsilon_N$ are independent noises such that $\mathbb{E}(\varepsilon_j) = 0$. Furthermore, we consider a random
 99 design setting where the key vectors $\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_N$ are i.i.d. samples from the distribution that
 100 admits p as density function. By an abuse of notation, we also denote p as the joint density where the
 101 key and value vectors $(\mathbf{v}_1, \mathbf{k}_1), \dots, (\mathbf{v}_N, \mathbf{k}_N)$ are i.i.d. samples from. Here, f is a true but unknown
 102 function and we would like to estimate it.

103 **Nadaraya–Watson estimator** Our approach to estimate the function f is based on
 104 Nadaraya–Watson’s nonparametric kernel regression approach [48]. In particular, from the nonpara-
 105 metric regression model (3), we have $\mathbb{E}[\mathbf{v}_j | \mathbf{k}_j] = f(\mathbf{k}_j)$ for all $j \in [N]$. Therefore, it is sufficient to
 106 estimate the conditional distribution of the value vectors given the key vectors. Given the density
 107 function p of the key vectors and the joint density p of the key and value vectors, for any pair of
 108 vectors (\mathbf{v}, \mathbf{k}) generate from model (3) we have

$$\mathbb{E}[\mathbf{v} | \mathbf{k}] = \int_{\mathbb{R}^D} \mathbf{v} \cdot p(\mathbf{v} | \mathbf{k}) d\mathbf{v} = \int \frac{\mathbf{v} \cdot p(\mathbf{v}, \mathbf{k})}{p(\mathbf{k})} d\mathbf{v}. \quad (4)$$

109 The formulation (4) of the conditional expectation indicates that as long as we can estimate the joint
 110 density function $p(\mathbf{v}, \mathbf{k})$ and the marginal density function $p(\mathbf{v})$, we are able to obtain an estimation
 111 for the conditional expectation and thus for the function f . This approach is widely known as
 112 Nadaraya–Watson’s nonparametric kernel regression approach.

113 **Kernel density estimator** To estimate $p(\mathbf{v}, \mathbf{k})$ and $p(\mathbf{k})$, we employ the kernel density estimation
 114 approach [59, 50]. In particular, by using the isotropic Gaussian kernel with bandwidth σ , we have
 115 the following estimators of $p(\mathbf{v}, \mathbf{k})$ and $p(\mathbf{k})$:

$$\hat{p}_\sigma(\mathbf{v}, \mathbf{k}) = \frac{1}{N} \sum_{j=1}^N \varphi_\sigma(\mathbf{v} - \mathbf{v}_j) \varphi_\sigma(\mathbf{k} - \mathbf{k}_j), \quad \hat{p}_\sigma(\mathbf{k}) = \frac{1}{N} \sum_{j=1}^N \varphi_\sigma(\mathbf{k} - \mathbf{k}_j), \quad (5)$$

116 where $\varphi_\sigma(\cdot)$ is the isotropic multivariate Gaussian density function with diagonal covariance matrix
 117 $\sigma^2 \mathbf{I}_D$. Given the kernel density estimators (5), we obtain the following estimation of the function f :

$$\begin{aligned} \hat{f}_\sigma(\mathbf{k}) &= \int_{\mathbb{R}^D} \frac{\mathbf{v} \cdot \hat{p}_\sigma(\mathbf{v}, \mathbf{k})}{\hat{p}_\sigma(\mathbf{k})} d\mathbf{v} = \int_{\mathbb{R}^D} \frac{\mathbf{v} \cdot \sum_{j=1}^N \varphi_\sigma(\mathbf{v} - \mathbf{v}_j) \varphi_\sigma(\mathbf{k} - \mathbf{k}_j)}{\sum_{j=1}^N \varphi_\sigma(\mathbf{k} - \mathbf{k}_j)} d\mathbf{v} \\ &= \frac{\sum_{j=1}^N \varphi_\sigma(\mathbf{k} - \mathbf{k}_j) \int \mathbf{v} \cdot \varphi_\sigma(\mathbf{v} - \mathbf{v}_j) d\mathbf{v}}{\sum_{j=1}^N \varphi_\sigma(\mathbf{k} - \mathbf{k}_j)} = \frac{\sum_{j=1}^N v_j \varphi_\sigma(\mathbf{k} - \mathbf{k}_j)}{\sum_{j=1}^N \varphi_\sigma(\mathbf{k} - \mathbf{k}_j)}. \end{aligned} \quad (6)$$

118 **Connection between Self-Attention and nonparametric regression** By plugging the query vectors
 119 \mathbf{q}_i into the function \hat{f}_σ in equation (6), we obtain that

$$\begin{aligned}\hat{f}_\sigma(\mathbf{q}_i) &= \frac{\sum_j^N \mathbf{v}_j \exp(-\|\mathbf{q}_i - \mathbf{k}_j\|^2/2\sigma^2)}{\sum_j^N \exp(-\|\mathbf{q}_i - \mathbf{k}_j\|^2/2\sigma^2)} \\ &= \frac{\sum_j^N \mathbf{v}_j \exp[-(\|\mathbf{q}_i\|^2 + \|\mathbf{k}_j\|^2)/2\sigma^2] \exp(\mathbf{q}_i \mathbf{k}_j^\top / \sigma^2)}{\sum_j^N \exp[-(\|\mathbf{q}_i\|^2 + \|\mathbf{k}_j\|^2)/2\sigma^2] \exp(\mathbf{q}_i \mathbf{k}_j^\top / \sigma^2)}.\end{aligned}\quad (7)$$

120 If we further assume that the keys \mathbf{k}_j are normalized, which is usually done in practice to stabilize
 121 the training of transformers [64], the value of $\hat{f}_\sigma(\mathbf{q}_i)$ in equation (6) then becomes

$$\hat{f}_\sigma(\mathbf{q}_i) = \frac{\sum_j^N \mathbf{v}_j \exp(\mathbf{q}_i \mathbf{k}_j^\top / \sigma^2)}{\sum_{j'}^N \exp(\mathbf{q}_i \mathbf{k}_{j'}^\top / \sigma^2)} = \sum_{j=1}^N \text{softmax}(\mathbf{q}_i^\top \mathbf{k}_j / \sigma^2) \mathbf{v}_j. \quad (8)$$

122 When we choose $\sigma^2 = \sqrt{D}$ where D is the dimension of \mathbf{q}_i and \mathbf{k}_j , equation (8) matches equa-
 123 tion (2) of self-attention, namely, $\hat{f}_\sigma(\mathbf{q}_i) = \mathbf{h}_i$. Thus, we have shown that self-attention performs
 124 nonparametric regression using isotropic Gaussian kernels.

125 **Remark 1** *The assumption that \mathbf{k}_j is normalized is to recover the pairwise dot-product attention in*
 126 *transformers. In general, this assumption is not necessary. In fact, the isotropic Gaussian kernel in*
 127 *equation (7) is more desirable than the dot-product kernel in equation (8) of the pairwise dot-product*
 128 *attention since the former is Lipschitz while the later is not Lipschitz [35]. The Lipschitz constraint*
 129 *helps improve the robustness of the model [15, 74, 2] and stabilize the model training [46].*

130 **Limitation of Self-Attention** From our nonparametric regression interpretation, self-attention is
 131 derived from the use of isotropic Gaussian kernels for kernel density estimation and nonparametric
 132 regression estimation, which may fail to capture the complex correlations between D features
 133 in \mathbf{q}_i and \mathbf{k}_j [81, 31]. Using multivariate Gaussian kernels with dense covariance matrices can
 134 help capture such correlations; however, choosing good covariance matrices is challenging and
 135 inefficient [80, 66, 11]. In the following section, we discuss the Fourier integral estimator and its use
 136 as a kernel for computing self-attention in order to overcome these limitations.

137 3 FourierFormer: Transformer via Generalized Fourier Integral Theorem

138 In the following, we introduce generalized integral theorems that are able to capture the complex
 139 interactions among the features of the queries and keys. We then apply these theorems to density
 140 estimation and nonparametric regression problems. We also establish the convergence rates of these
 141 estimators. Given these density estimators, we introduce a novel family of transformers, named
 142 *FourierFormer*, that integrates the generalized Fourier integral theorem into the dot-product attention
 143 step of the standard transformer.

144 3.1 Generalized Fourier Integral Theorems and Their Applications

145 The Fourier integral theorem is a beautiful result in mathematics [85, 7] and has been recently used
 146 in nonparametric mode clustering, deconvolution problem, and generative modeling [31]. It is a
 147 combination of Fourier transform and Fourier inverse transform. In particular, for any function
 148 $p \in \mathbb{L}_1(\mathbb{R}^D)$, the *Fourier integral theorem* is given by

$$\begin{aligned}p(\mathbf{k}) &= \frac{1}{(2\pi)^D} \int_{\mathbb{R}^D} \int_{\mathbb{R}^D} \cos(\mathbf{s}^\top (\mathbf{k} - \mathbf{y})) p(\mathbf{y}) d\mathbf{y} d\mathbf{s} \\ &= \frac{1}{\pi^D} \lim_{R \rightarrow \infty} \int_{\mathbb{R}^D} \prod_{j=1}^D \frac{\sin(R(k_j - y_j))}{(k_j - y_j)} p(\mathbf{y}) d\mathbf{y},\end{aligned}\quad (9)$$

149 where $\mathbf{k} = (k_1, \dots, k_D)$, $\mathbf{y} = (y_1, \dots, y_D)$, $\mathbf{s} = (s_1, \dots, s_D)$, and R is the radius. The de-
 150 tailed derivation of Equation (9) is in Appendix A.3. Equation (9) suggests that $p_R(\mathbf{k}) :=$

151 $\frac{1}{\pi^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \frac{\sin(R(y_j - k_j))}{(y_j - k_j)} p(\mathbf{y}) d\mathbf{y}$ can be used as an estimator of the function p .

152 **Benefits of the Fourier integral over Gaussian kernel** There are two important benefits of the
 153 estimator p_R : (i) it can automatically preserve the correlated structure lying within p even when p is

154 very complex and high dimensional function. It is in stark contrast to the standard kernel estimator
 155 built based on multivariate Gaussian kernel where we need to choose good covariance matrix in the
 156 multivariate Gaussian kernel to guarantee such estimator to work well. We note that as the standard
 157 soft-max Transformer is constructed based on the multivariate Gaussian kernel, the issue of choosing
 158 good covariance matrix in dot-product transformer is inevitable; (ii) The product of sinc kernels in
 159 the estimator p_R does not decay to a point mass when $R \rightarrow \infty$. It is in stark difference from the
 160 multivariate Gaussian kernel estimator, which converges to a point mass when the covariance matrix
 161 goes to 0. It indicates that p_R is a non-trivial estimator of the function p . Finally, detailed illustrations
 162 of these benefits of the Fourier integral over Gaussian kernel in density estimation and nonparametric
 163 regression problems, which we have just shown to have connection to the self-attention in transformer,
 164 can be found in Section 8 in [31].

165 **Generalized Fourier integral estimator** Borrowing the above benefits of Fourier integral estimator
 166 p_R , in the paper we would like to consider a generalization of that estimator, named *generalized*
 167 *Fourier integral estimator*, which is given by:

$$p_R^\phi(\mathbf{k}) := \frac{R^D}{A^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(R(y_j - k_j))}{R(y_j - k_j)}\right) p(\mathbf{y}) d\mathbf{y}, \quad (10)$$

168 where $A := \int_{\mathbb{R}} \phi\left(\frac{\sin(z)}{z}\right) dz$ and $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is a given function. When $\phi(\mathbf{k}) = \mathbf{k}$ for all
 169 $\mathbf{k} \in \mathbb{R}^D$, the generalized Fourier integral estimator p_R^ϕ becomes the Fourier integral estimator p_R .
 170 Under appropriate conditions on the function ϕ (see Theorem 1 in Section 3.1.1 and Theorem 3 in
 171 Appendix B.1), the estimator p_R^ϕ converges to the true function p , namely,

$$p(\mathbf{k}) = \lim_{R \rightarrow \infty} p_R^\phi(\mathbf{k}) = \lim_{R \rightarrow \infty} \frac{R^D}{A^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(R(y_j - k_j))}{R(y_j - k_j)}\right) p(\mathbf{y}) d\mathbf{y}. \quad (11)$$

172 We name the above limit as *generalized Fourier integral theorem*. Furthermore, the estimator p_R^ϕ also
 173 inherits similar aforementioned benefits of the Fourier integral estimator p_R . Therefore, we will use
 174 the generalized Fourier integral theorem as a building block for constructing density estimators and
 175 nonparametric regression estimators, which are crucial to develop the FourierFormer in Section 3.2.

176 3.1.1 Density Estimation via Generalized Fourier Integral Theorems

177 We first apply the generalized Fourier integral theorem to the density estimation problem. To ease the
 178 presentation, we assume that $\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_N \in \mathbb{R}^D$ are i.i.d. samples from a distribution admitting
 179 density function p where $D \geq 1$ is the dimension. Inspired by the generalized Fourier integral
 180 theorem, we obtain the following *generalized Fourier density estimator* $p_{N,R}^\phi$ of p as follows:

$$p_{N,R}^\phi(\mathbf{k}) := \frac{R^D}{NA^D} \sum_{i=1}^N \prod_{j=1}^D \phi\left(\frac{\sin(R(k_j - k_{ij}))}{R(k_j - k_{ij})}\right), + \quad (12)$$

181 where $A = \int_{\mathbb{R}} \phi\left(\frac{\sin(z)}{z}\right) dz$ and $\mathbf{k}_i = (k_{i1}, \dots, k_{iD})$ for all $i \in [N]$. To quantify the error between
 182 the generalized Fourier density estimator $p_{n,R}^\phi$ and the true density p , we utilize mean integrated
 183 squared errors (MISE) [84], which is given by:

$$\text{MISE}(p_{N,R}^\phi, p) := \int_{\mathbb{R}^D} (p_{N,R}^\phi(\mathbf{k}) - p(\mathbf{k}))^2 d\mathbf{k}. \quad (13)$$

184 We start with the following bound on the MISE between $p_{n,R}^\phi$ and p .

185 **Theorem 1** Assume that $\int_{\mathbb{R}} \phi(\sin(z)/z) z^j dz = 0$ for all $j \in [m]$ and $\int_{\mathbb{R}} |\phi(\sin(z)/z)| |z|^{m+1} dz <$
 186 ∞ for some $m \in \mathbb{N}$. Then, there exist universal constants C and C' depending on d and A such that

$$\text{MISE}(p_{N,R}^\phi, p) \leq \frac{C}{R^{m+1}} + \frac{C' R^D}{N}.$$

187 Proof of Theorem 1 is in Appendix C.1. A few comments are in order. First, by choosing R
 188 to balance the bias and variance in the bound of MISE in Theorem 1, we have the optimal R as

189 $R = \mathcal{O}(N^{1/(D+m+1)})$. With that choice of R , the MISE rate of $p_{N,R}^\phi$ is $\mathcal{O}(N^{-(m+1)/(D+m+1)})$.
 190 Second, when $\phi(z) = z^l$ for $l \geq 4$ and $z \in \mathbb{R}$, the assumptions in Theorem 1 are satisfied when
 191 $m = 1$. Under this case, the MISE rate of $p_{N,R}^\phi$ is $\mathcal{O}(N^{-2/(D+2)})$. However, these assumptions
 192 do not satisfy when $\phi(z) = z^l$ and $l \in \{1, 2, 3\}$, which is due to the limitation of the current proof
 193 technique of Theorem 1 that is based on Taylor expansion of the estimator $p_{n,R}^\phi$.

194 To address the limitation of the Taylor expansion technique, we utilize the Plancherel theorem in
 195 Fourier analysis to establish the MISE rate of $p_{N,R}^\phi$ when $\phi(z) = z^l$ and $l \in \{1, 2, 3\}$. The details of
 196 the theoretical analyses for such setting are in Appendix B.

197 3.2 FourierFormer: Transformers with Fourier Attentions

198 Motivated by the preservation of the correlated structure of the function from the generalized Fourier
 199 integral theorem as well as the theoretical guarantees of density estimators, in this section we adapt
 200 the nonparametric regression interpretation of self-attention in Section 2 and propose the generalized
 201 Fourier nonparametric regression estimator in Section 3.2.1. We also establish the convergence
 202 properties of that estimator. Then, based on generalized Fourier nonparametric regression estimator,
 203 we develop the Fourier Attention and its corresponding FourierFormer in Section 3.2.2.

204 3.2.1 Nonparametric Regression via Generalized Fourier Integral Theorem

205 We now discuss an application of the generalized Fourier integral theorems to the nonparametric
 206 regression setting (3), namely, we assume that $(\mathbf{v}_1, \mathbf{k}_1), \dots, (\mathbf{v}_N, \mathbf{k}_N)$ are i.i.d. samples from the
 207 following nonparametric regression model:

$$\mathbf{v}_j = f(\mathbf{k}_j) + \varepsilon_j,$$

208 where $\varepsilon_1, \dots, \varepsilon_N$ are independent noises such that $\mathbb{E}(\varepsilon_j) = 0$ and the key vectors $\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_N$ are
 209 i.i.d. samples from p . Given the generalized Fourier density estimator (12), following the argument in
 210 Section 2, the Nadaraya–Watson estimator of the function f based on the generalized Fourier density
 211 estimator is given by:

$$f_{N,R}(\mathbf{k}) := \frac{\sum_{i=1}^N \mathbf{v}_i \prod_{j=1}^D \phi\left(\frac{\sin(R(k_j - k_{ij}))}{R(k_j - k_{ij})}\right)}{\sum_{i=1}^N \prod_{j=1}^D \phi\left(\frac{\sin(R(k_j - k_{ij}))}{R(k_j - k_{ij})}\right)}. \quad (14)$$

212 The main difference between the generalized Fourier nonparametric regression estimator $f_{N,R}$ in
 213 equation (14) and the estimator \hat{f}_σ in equation (6) is that the estimator $f_{N,R}$ utilizes the generalized
 214 Fourier density estimator to estimate the conditional distribution of the value vectors given the key
 215 vectors instead of the isotropic Gaussian kernel density estimator as in \hat{f}_σ . As we highlighted in
 216 Section 3, an important benefit of the generalized Fourier density estimator is that it can capture the
 217 complex dependencies of the features of the value vectors and the key vectors while the Gaussian
 218 kernel needs to have good covariance matrix to do that, which is computationally expensive in
 219 practice.

220 We now have the following result establishing the mean square error (MSE) of $f_{N,R}$.

221 **Theorem 2** Assume that $\int_{\mathbb{R}} \phi\left(\frac{\sin(z)}{z}\right) z^j dz = 0$ for all $1 \leq j \leq m$ and $\int_{\mathbb{R}} \left|\phi\left(\frac{\sin(z)}{z}\right)\right| |z|^j dz < \infty$
 222 for any $m+1 \leq j \leq 2m+2$ for some $m \in \mathbb{N}$. Then, for any $\mathbf{k} \in \mathbb{R}^D$, there exist universal constants
 223 C_1, C_2, C_3, C_4 such that the following holds:

$$\mathbb{E}[(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2] \leq \left(\frac{C_1}{R^{2(m+1)}} + \frac{(f(\mathbf{k}) + C_2)R^D}{N} \right) / (p^2(\mathbf{k})J(R)),$$

224 where $J(R) = 1 - \frac{1}{p^2(\mathbf{k})} \left(\frac{C_3}{R^{2(m+1)}} + \frac{C_4 R^d \log(NR)}{N} \right)$. Here, the outer expectation is taken with
 225 respect to the key vectors $\mathbf{k}_1, \dots, \mathbf{k}_N$ and the noises $\varepsilon_1, \dots, \varepsilon_N$.

226 Proof of Theorem 2 is in Appendix C.3. A few comments with Theorem 2 are in order. First, by
 227 choosing R to balance the bias and variance in the bound of the MSE of the nonparametric generalized
 228 Fourier estimator $f_{N,R}$, we have the optimal radius R as $R = \mathcal{O}(N^{\frac{1}{2(m+1)+D}})$. With that choice of
 229 the optimal radius R , the rate of $f_{N,R}$ is $\mathcal{O}(N^{-\frac{2(m+1)}{D+2(m+1)}})$. Second, when $\phi(z) = z^l$ for $l \geq 6$, the

230 assumption on the function ϕ of Theorem 2 is satisfied with $m = 1$. Under this case, the rate of $f_{N,R}$
 231 becomes $\mathcal{O}(N^{-\frac{4}{D+4}})$. In Appendix B, we also provide the rate of $f_{N,R}$ when $\phi(z) = z^l$ for some
 232 $l \leq 5$, which includes the original Fourier integral theorem.

233 3.2.2 FourierFormer

234 Given the generalized Fourier nonparametric regression estimator $f_{N,R}$ in equation (14), by plugging
 235 the query values $\mathbf{q}_1, \dots, \mathbf{q}_N$ into that function, we obtain the following definition of the Fourier
 236 attention:

237 **Definition 1 (Fourier Attention)** A Fourier attention is a multi-head attention that does nonpara-
 238 metric regression using the generalized Fourier nonparametric regression estimator $f_{N,R}$. The output
 239 $\hat{\mathbf{h}}_i$ of the Fourier attention is then computed as

$$\hat{\mathbf{h}}_i := f_{N,R}(\mathbf{q}_i) = \frac{\sum_{i=1}^N \mathbf{v}_i \prod_{j=1}^D \phi\left(\frac{\sin(R(\mathbf{q}_{ij} - \mathbf{k}_{ij}))}{R(\mathbf{q}_{ij} - \mathbf{k}_{ij})}\right)}{\sum_{i=1}^N \prod_{j=1}^D \phi\left(\frac{\sin(R(\mathbf{q}_{ij} - \mathbf{k}_{ij}))}{R(\mathbf{q}_{ij} - \mathbf{k}_{ij})}\right)} \quad \forall i \in [N]. \quad (15)$$

240 Given the Fourier Attention in Definition 1, we then give the definition of FourierFormer as follows.

241 **Definition 2 (FourierFormer)** A FourierFormer is a transformer that uses Fourier attention to
 242 capture dependency between tokens in the input sequence and the correlation between features in
 243 each token.

244 **Remark 2 (The Nonnegativity of the Fourier Kernel)** The density estimation via generalized
 245 Fourier integral theorem in Section 3.1.1 does not require the generalized Fourier density esti-
 246 mator to be nonnegative. However, empirically, we observe that negative density estimator can cause
 247 instability in training the FourierFormer. Thus, in FourierFormer, we choose the function ϕ to be a
 248 nonnegative function to enforce the density estimator to be nonnegative. In particular, we choose ϕ to
 249 be power functions of the form $\phi(x) = x^{2m}$, where m is an positive integer. Note that when $m = 2$
 250 and $m = 4$, the kernels in our generalized Fourier integral estimators are the well-known Fejer-de la
 251 Vallee Poussin and Jackson-de la Vallee Poussin kernels [19].

252 3.3 An Efficient Implementation of the Fourier Attention

The Fourier kernel is implemented efficiently in the C++/CUDA extension developed by Pytorch
 [51]. The idea is similar to the function `cdist` [51], which computes the p-norm distance between
 each pair of the two collections of row vectors. In our case, we aim to compute kernel functions that
 represent a Fourier attention in Definition 1. The core of this implementation is the following Fourier
 metric function d_f :

$$d_f(\mathbf{q}_i, \mathbf{k}_j) = \prod_{d=1}^D \phi\left(\frac{\sin(R(\mathbf{q}_{id} - \mathbf{k}_{jd}))}{R(\mathbf{q}_{id} - \mathbf{k}_{jd})}\right)$$

253 We directly implement d_f as a `torch.autograd.Function` [51] in which we provide an efficient
 254 way to compute forward and backward function (d_f and gradient of d_f). While the implementation
 255 of the forward function is straight forward, the backward function is more tricky since we need to
 256 optimize the code to compute the gradient of d_f w.r.t to variables \mathbf{q} , \mathbf{k} , and R all at once. We can
 257 develop the backward function with highly parallel computation by exploiting GPU architecture and
 258 utilizing the reduction technique. The computational time is comparable to function `cdist`; thus, our
 259 FourierFormer implementation is as computationally time-efficient.

260 4 Experimental Results

261 In this section, we numerically justify the advantage of FourierFormer over the baseline dot-product
 262 transformer on two large-scale tasks: language modeling on WikiText-103 [44] (Section 4.1) and
 263 image classification on ImageNet [21, 60] (Section 4.2). We aim to show that: (i) FourierFormer
 264 achieves better accuracy than the baseline transformer on a variety of practical tasks with different
 265 data modalities, and (ii) FourierFormer helps reduce head redundancy compared to the baseline
 266 transformer (Section 4.3).

267 Throughout the section, we compare FourierFormers with the baseline dot-product transformers
 268 of the same configuration. In all experiments, we made the constant R in Fourier attention (see

Table 1. Perplexity (PPL) on WikiText-103 of FourierFormers compared to the baselines. FourierFormers achieve much better PPL than the baselines.

Method	Valid PPL	Test PPL
<i>Baseline dot-product (small)</i>	33.15	34.29
FourierFormer (small)	31.86	32.85
<i>Baseline dot-product (medium)</i>	27.90	29.60
FourierFormer (medium)	26.51	28.01

269 equation (58)) to be a learnable scalar and set choose the function $\phi(x) = x^4$ (see Remark 2). All of
 270 our results are averaged over 5 runs with different seeds. More details on the models and training are
 271 provided in Appendix D. We also provide additional experimental results in Appendix E.

272 4.1 Language Modeling on WikiText-103

273 **Datasets and metrics** WikiText-103 is a collection of articles from Wikipedia, which have long
 274 contextual dependencies. The training set consists of about $28K$ articles containing $103M$ running
 275 words; this corresponds to text blocks of about 3600 words. The validation and test sets have $218K$
 276 and $246K$ running words, respectively. Each of them contains 60 articles and about $268K$ words. Our
 277 experiment follows the standard setting [44, 64] and splits the training data into L -word independent
 278 long segments. For evaluation, we use a batch size of 1, and process the text sequence with a sliding
 279 window of size L . The last position is used for computing perplexity (PPL) except in the first segment,
 280 where all positions are evaluated as in [1, 64].

281 **Models and baselines** Our implementation is based on the public code by [64].¹ We use their
 282 small and medium models in our experiments. In particular, for small models, the key, value, and
 283 query dimension are set to 128, and the training and evaluation context length are set to 256. For
 284 medium models, the key, value, and query dimension are set to 256, and the training and evaluation
 285 context length are set to 384. In both configurations, the number of heads is 8, the feed-forward layer
 286 dimension is 2048, and the number of layers is 16.

287 **Results** We report the validation and test perplexity (PPL) of FourierFormer versus the baseline
 288 transformer with the dot-product attention in Table 1. FourierFormers attain much better PPL than the
 289 baselines in both small and medium configurations. For the small configuration, the improvements of
 290 FourierFormer over the baseline are 1.29 PPL in validation and 1.44 PPL in test. For the medium
 291 configuration, these improvements are 1.39 PPL in validation and 1.59 PPL in test. These results
 292 suggest that the advantage of FourierFormer over the baseline dot-product transformer grows with the
 293 model’s size. This meets our expectation because larger models has larger query and key dimensions,
 294 e.g. the language model with medium configuration in this experiment has the query and key
 295 dimension of 256 versus 128 as in the language model with small configuration. Since the advantage
 296 of FourierFormer results from the property that FourierFormer can capture correlation between
 297 features in query and key vectors, the larger the query and key dimensions are, the more advantage
 298 FourierFormer has.

299 4.2 Image Classification on ImageNet

300 **Datasets and metrics** The ImageNet dataset [21, 60] consists of $1.28M$ training images and $50K$
 301 validation images. For this benchmark, the model learns to predict the category of the input image
 302 among 1000 categories. Top-1 and top-5 classification accuracies are reported.

303 **Models and baselines** We use the DeiT-tiny model [72] with 12 transformer layers, 4 attention heads
 304 per layer, and the model dimension of 192. To train the models, we follow the same setting and
 305 configuration as for the baseline [72].²

306 **Results** We summarize our results in Table 2. Same as in the language modeling experiment, for this
 307 image classification task, the Deit model equipped with FourierFormer significantly outperforms the
 308 baseline Deit dot-product transformer in both top-1 and top-5 accuracy. This result suggests that the
 309 advantage of FourierFormer over the baseline dot-product transformer holds across different data
 310 modalities.

¹Implementation available at <https://github.com/IDSIA/lmtool-fwp>.

²Implementation available at <https://github.com/facebookresearch/deit>.

Table 2. Top-1 and top-5 accuracy (%) of FourierFormer Deit vs. the baseline Deit with dot-product attention. FourierFormer Deit outperforms the baseline in both top-1 and top-5 accuracy.

Method	Top-1 Acc	Top-5 Acc
<i>Baseline DeiT</i>	72.23	91.13
FourierFormer DeiT	73.25	91.66

Table 3. Layer-average mean and standard deviation of \mathcal{L}_2 distances between heads of FourierFormer versus the baseline transformer with dot-product attention trained for the WikiText-103 language modeling task. FourierFormer has greater \mathcal{L}_2 distance between heads than the baseline and thus captures more diverse attention patterns.

Method	Train	Test
<i>Baseline dot-product</i>	6.20 \pm 2.30	6.17 \pm 2.30
FourierFormer	7.45 \pm 2.50	7.37 \pm 2.44

311 4.3 FourierFormer Helps Reducing Head Redundancy

312 To study the diversity between attention heads, given the model trained for the WikiText-103 language
 313 modeling task, we compute the average \mathcal{L}_2 distance between heads in each layer. We show the
 314 layer-average mean and variance of distances between heads in Table 3. Results in Table 3 shows
 315 that FourierFormer obtains greater \mathcal{L}_2 distance between attention heads than the baseline transformer
 316 with the dot-product attention and thus helps reduce the head redundancy. Note that we use the small
 317 configuration as specified in Section 4.1 for both models.

318 5 Related Work

319 **Interpretation of Attention Mechanism in Transformers** Recent works have tried to gain an
 320 understanding of transformer’s attention from different perspectives. [73] considers attention as
 321 applying kernel smoother over the inputs. Extending this kernel approach, [33, 14, 82] linearize the
 322 softmax kernel in dot-product attention and propose a family of efficient transformers with linear
 323 computational and memory complexity. [9] then shows that these linear transformers are comparable
 324 to a Petrov-Galerkin projection [57], suggesting that the softmax normalization in the dot-product
 325 attention is sufficient but not necessary. Other works provide an understanding of attention in
 326 transformers via ordinary/partial differential equation include [43, 62]. In addition, [68, 28, 89] relate
 327 attentions in transformers to a Gaussian mixture models. Several works also connect the attention
 328 mechanism to graph-structured learning and message passing in graphical models [83, 65, 37]. Our
 329 work focuses on deriving the connection between self-attention and nonparametric kernel regression
 330 and exploring better regression estimator, such as the generalized Fourier nonparametric regression
 331 estimator, to improve the performance of transformers.

332 **Redundancy in Transformers** [18, 45, 24] show that neurons and attention heads in the pre-trained
 333 transformer are redundant and can be removed when applied on a downstream task. By studying
 334 the contextualized embeddings in pre-trained networks, it has been demonstrated that the learned
 335 representations from these redundant models are highly anisotropic [47, 25]. Furthermore, [63, 67, 79,
 336 61] employ knowledge distillation and sparse approximation to enhance the efficiency of transformers.
 337 Our FourierFormer is complementary to these methods and can be combined with them.

338 6 Concluding Remarks

339 In this paper, we establish the correspondence between the nonparametric kernel regression and the
 340 self-attention in transformer. We then develop the generalized Fourier integral estimators and propose
 341 the FourierFormer, a novel class of transformers that use the generalized Fourier integral estimators to
 342 construct their attentions for efficiently capturing the correlations between features in the query and
 343 key vectors. We theoretically prove the approximation guarantees of the generalized Fourier integral
 344 estimators and empirically validate the advantage of FourierFormer over the baseline transformer
 345 with the dot-product attention in terms of accuracy and head redundancy reduction. It is interesting
 346 to incorporate robust kernels into the nonparametric regression framework of FourierFormer to
 347 enhance the robustness of the model under data perturbation and adversarial attacks. A limitation of
 348 FourierFormer is that it still has the same quadratic computational and memory complexity as the
 349 baseline transformer with the dot-product attention. We leave the development of the linear version
 350 of FourierFormer that achieves linear computational and memory complexity as future work. It is
 351 worth noting that there is no potential negative societal impacts of FourierFormer.

References

- 352
- 353 [1] Rami Al-Rfou, Dokook Choe, Noah Constant, Mandy Guo, and Llion Jones. Character-level
354 language modeling with deeper self-attention. In *Proceedings of the AAAI Conference on*
355 *Artificial Intelligence*, volume 33, pages 3159–3166, 2019.
- 356 [2] Cem Anil, James Lucas, and Roger Grosse. Sorting out lipschitz function approximation. In
357 *International Conference on Machine Learning*, pages 291–301. PMLR, 2019.
- 358 [3] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia
359 Schmid. Vivit: A video vision transformer. In *2021 IEEE/CVF International Conference on*
360 *Computer Vision (ICCV)*, pages 6816–6826, 2021.
- 361 [4] Alexei Baevski and Michael Auli. Adaptive input representations for neural language modeling.
362 In *International Conference on Learning Representations*, 2019.
- 363 [5] Anthony Bagnall, Hoang Anh Dau, Jason Lines, Michael Flynn, James Large, Aaron Bostrom,
364 Paul Southam, and Eamonn Keogh. The uea multivariate time series classification archive, 2018.
365 *arXiv preprint arXiv:1811.00075*, 2018.
- 366 [6] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
367 learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- 368 [7] S. Bochner. *Lectures on Fourier Integrals*. Princeton University Press, 1959.
- 369 [8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
370 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
371 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 372 [9] Shuhao Cao. Choose a transformer: Fourier or galerkin. *Advances in Neural Information*
373 *Processing Systems*, 34, 2021.
- 374 [10] Mauro Cettolo, Jan Niehues, Sebastian Stüker, Luisa Bentivogli, and Marcello Federico. Report
375 on the 11th iwslt evaluation campaign, iwslt 2014. In *Proceedings of the International Workshop*
376 *on Spoken Language Translation, Hanoi, Vietnam*, volume 57, 2014.
- 377 [11] J.E. Chacón and T. Duong. *Multivariate Kernel Smoothing and its Applications*. CRC Press,
378 2018.
- 379 [12] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with
380 sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- 381 [13] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares,
382 Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-
383 decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical*
384 *Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar, October
385 2014. Association for Computational Linguistics.
- 386 [14] Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea
387 Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser,
388 David Benjamin Belanger, Lucy J Colwell, and Adrian Weller. Rethinking attention with
389 performers. In *International Conference on Learning Representations*, 2021.
- 390 [15] Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier.
391 Parseval networks: Improving robustness to adversarial examples. In *International Conference*
392 *on Machine Learning*, pages 854–863. PMLR, 2017.
- 393 [16] Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What does
394 BERT look at? an analysis of BERT’s attention. In *Proceedings of the 2019 ACL Workshop*
395 *BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence,
396 Italy, August 2019. Association for Computational Linguistics.
- 397 [17] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov.
398 Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint*
399 *arXiv:1901.02860*, 2019.

- 400 [18] Fahim Dalvi, Hassan Sajjad, Nadir Durrani, and Yonatan Belinkov. Analyzing redundancy in
401 pretrained transformer models. *arXiv preprint arXiv:2004.04010*, 2020.
- 402 [19] Kathryn Bullock Davis. Mean square error properties of density estimates. *The Annals of*
403 *Statistics*, 3(4):1025–1030, 1975.
- 404 [20] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Uni-
405 versal transformers. In *International Conference on Learning Representations*, 2019.
- 406 [21] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
407 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*
408 *recognition*, pages 248–255. Ieee, 2009.
- 409 [22] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
410 deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Confer-*
411 *ence of the North American Chapter of the Association for Computational Linguistics: Human*
412 *Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis,
413 Minnesota, June 2019. Association for Computational Linguistics.
- 414 [23] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
415 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
416 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
417 recognition at scale. In *International Conference on Learning Representations*, 2021.
- 418 [24] Nadir Durrani, Hassan Sajjad, Fahim Dalvi, and Yonatan Belinkov. Analyzing individual
419 neurons in pre-trained language models. *arXiv preprint arXiv:2010.02695*, 2020.
- 420 [25] Kawin Ethayarajh. How contextual are contextualized word representations? comparing the
421 geometry of bert, elmo, and gpt-2 embeddings. *arXiv preprint arXiv:1909.00512*, 2019.
- 422 [26] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and
423 Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF*
424 *International Conference on Computer Vision*, pages 6824–6835, 2021.
- 425 [27] J. Fan. On the optimal rates of convergence for nonparametric deconvolution problems. *Annals*
426 *of Statistics*, 19(3):1257–1272, 1991.
- 427 [28] Prasad Gabbur, Manjot Bilkhu, and Javier Movellan. Probabilistic attention for interactive
428 segmentation. *Advances in Neural Information Processing Systems*, 34, 2021.
- 429 [29] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min
430 Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7(2):187–199, 2021.
- 431 [30] John Hewitt and Percy Liang. Designing and interpreting probes with control tasks. In *Proceed-*
432 *ings of the 2019 Conference on Empirical Methods in Natural Language Processing and the*
433 *9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages
434 2733–2743, Hong Kong, China, November 2019. Association for Computational Linguistics.
- 435 [31] N. Ho and S.G. Walker. Multivariate smoothing via the Fourier integral theorem and Fourier
436 kernel. *Arxiv preprint Arxiv:2012.14482*, 2021.
- 437 [32] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ron-
438 neberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al.
439 Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- 440 [33] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers
441 are rnns: Fast autoregressive transformers with linear attention. In *International Conference on*
442 *Machine Learning*, pages 5156–5165. PMLR, 2020.
- 443 [34] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan,
444 and Mubarak Shah. Transformers in vision: A survey. *ACM Computing Surveys (CSUR)*, 2021.
- 445 [35] Hyunjik Kim, George Papamakarios, and Andriy Mnih. The lipschitz constant of self-attention.
446 In *International Conference on Machine Learning*, pages 5562–5571. PMLR, 2021.

- 447 [36] Yoon Kim, Carl Denton, Luong Hoang, and Alexander M Rush. Structured attention networks.
448 *arXiv preprint arXiv:1702.00887*, 2017.
- 449 [37] Devin Kreuzer, Dominique Beaini, Will Hamilton, Vincent Létourneau, and Prudencio Tossou.
450 Rethinking graph transformers with spectral attention. *Advances in Neural Information Pro-*
451 *cessing Systems*, 34, 2021.
- 452 [38] Tianyang Lin, Yuxin Wang, Xiangyang Liu, and Xipeng Qiu. A survey of transformers. *arXiv*
453 *preprint arXiv:2106.04554*, 2021.
- 454 [39] Zhouhan Lin, Minwei Feng, Cícero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou,
455 and Yoshua Bengio. A structured self-attentive sentence embedding. *CoRR*, abs/1703.03130,
456 2017.
- 457 [40] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
458 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
459 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 460 [41] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining
461 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings*
462 *of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021.
- 463 [42] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin
464 transformer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- 465 [43] Yiping Lu, Zhuohan Li, Di He, Zhiqing Sun, Bin Dong, Tao Qin, Liwei Wang, and Tie-Yan Liu.
466 Understanding and improving transformer from a multi-particle dynamic system point of view.
467 *arXiv preprint arXiv:1906.02762*, 2019.
- 468 [44] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
469 models. In *5th International Conference on Learning Representations, ICLR 2017, Toulon,*
470 *France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- 471 [45] Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In
472 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors,
473 *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 474 [46] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization
475 for generative adversarial networks. In *International Conference on Learning Representations*,
476 2018.
- 477 [47] Jiaqi Mu and Pramod Viswanath. All-but-the-top: Simple and effective postprocessing for word
478 representations. In *International Conference on Learning Representations*, 2018.
- 479 [48] E.A. Nadaraya. On estimating regression. *Theory of Probability and its Applications*, 9:141–142,
480 1964.
- 481 [49] Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention
482 model for natural language inference. In *Proceedings of the 2016 Conference on Empirical*
483 *Methods in Natural Language Processing*, pages 2249–2255, Austin, Texas, November 2016.
484 Association for Computational Linguistics.
- 485 [50] E. Parzen. On estimation of a probability density function and mode. *Annals of Mathematical*
486 *Statistics*, 33:1065–1076, 1962.
- 487 [51] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
488 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
489 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
490 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style,
491 high-performance deep learning library. In *Advances in Neural Information Processing Systems*
492 32, pages 8024–8035. Curran Associates, Inc., 2019.

- 493 [52] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
494 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
495 models from natural language supervision. In *International Conference on Machine Learning*,
496 pages 8748–8763. PMLR, 2021.
- 497 [53] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language
498 understanding by generative pre-training. *OpenAI report*, 2018.
- 499 [54] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever.
500 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 501 [55] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
502 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified
503 text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- 504 [56] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark
505 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on
506 Machine Learning*, pages 8821–8831. PMLR, 2021.
- 507 [57] JN Reddy. *An introduction to the finite element method*, volume 1221. McGraw-Hill New York,
508 2004.
- 509 [58] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo,
510 Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from
511 scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National
512 Academy of Sciences*, 118(15), 2021.
- 513 [59] M. Rosenblatt. Remarks on some nonparametric estimates of a density function. *Annals of
514 Mathematical Statistics*, 27:832–837, 1956.
- 515 [60] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
516 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual
517 recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- 518 [61] Hassan Sajjad, Fahim Dalvi, Nadir Durrani, and Preslav Nakov. Poor man’s bert: Smaller and
519 faster transformer models. *arXiv e-prints*, pages arXiv–2004, 2020.
- 520 [62] Michael E Sander, Pierre Ablin, Mathieu Blondel, and Gabriel Peyré. Sinkformers: Transform-
521 ers with doubly stochastic attention. In *International Conference on Artificial Intelligence and
522 Statistics*, pages 3515–3530. PMLR, 2022.
- 523 [63] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version
524 of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- 525 [64] Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast
526 weight programmers. In *International Conference on Machine Learning*, pages 9355–9366.
527 PMLR, 2021.
- 528 [65] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position rep-
529 resentations. In *Proceedings of the 2018 Conference of the North American Chapter of the
530 Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short
531 Papers)*, pages 464–468, New Orleans, Louisiana, June 2018. Association for Computational
532 Linguistics.
- 533 [66] J.G. Staniswalis, K. Messer, and D.R. Finston. Kernel estimators for multivariate regression.
534 *Journal of Nonparametric Statistics*, 3:103–121, 1993.
- 535 [67] Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. Patient knowledge distillation for bert model
536 compression. *arXiv preprint arXiv:1908.09355*, 2019.
- 537 [68] Binh Tang and David S. Matteson. Probabilistic transformer for time series analysis. In
538 A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural
539 Information Processing Systems*, 2021.

- 540 [69] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey.
541 *arXiv preprint arXiv:2009.06732*, 2020.
- 542 [70] Ian Tenney, Dipanjan Das, and Ellie Pavlick. BERT rediscovers the classical NLP pipeline.
543 In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*,
544 pages 4593–4601, Florence, Italy, July 2019. Association for Computational Linguistics.
- 545 [71] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and
546 Hervé Jégou. Training data-efficient image transformers & distillation through attention. In
547 *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.
- 548 [72] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and
549 Hervé Jégou. Training data-efficient image transformers & distillation through attention. In
550 *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.
- 551 [73] Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan
552 Salakhutdinov. Transformer dissection: An unified understanding for transformer’s attention
553 via the lens of kernel. In *Proceedings of the 2019 Conference on Empirical Methods in Natural
554 Language Processing and the 9th International Joint Conference on Natural Language Process-
555 ing (EMNLP-IJCNLP)*, pages 4344–4353, Hong Kong, China, November 2019. Association for
556 Computational Linguistics.
- 557 [74] Yusuke Tsuzuku, Issei Sato, and Masashi Sugiyama. Lipschitz-margin training: Scalable
558 certification of perturbation invariance for deep neural networks. *Advances in neural information
559 processing systems*, 31, 2018.
- 560 [75] A. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, 2009.
- 561 [76] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
562 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information
563 processing systems*, pages 5998–6008, 2017.
- 564 [77] Jesse Vig and Yonatan Belinkov. Analyzing the structure of attention in a transformer language
565 model. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpret-
566 ing Neural Networks for NLP*, pages 63–76, Florence, Italy, August 2019. Association for
567 Computational Linguistics.
- 568 [78] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head
569 self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of
570 the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5797–5808,
571 Florence, Italy, July 2019. Association for Computational Linguistics.
- 572 [79] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head
573 self-attention: Specialized heads do the heavy lifting, the rest can be pruned. *arXiv preprint
574 arXiv:1905.09418*, 2019.
- 575 [80] M.P. Wand. Error analysis for general multivariate kernel estimators. *Journal of Nonparametric
576 Statistics*, 2:1–15, 1992.
- 577 [81] M.P. Wand and M.C. Jones. Comparison of smoothing parameterizations in bivariate kernel
578 density estimation. *Journal of the American Statistical Association*, 88:520–528, 1993.
- 579 [82] Sinong Wang, Belinda Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention
580 with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- 581 [83] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks.
582 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
583 7794–7803, 2018.
- 584 [84] L. Wasserman. *All of Nonparametric Statistics*. Springer, 2006.
- 585 [85] N. Wiener. *The Fourier Integral and Certain of its Applications*. Cambridge University Press,
586 1933.

- 587 [86] Haixu Wu, Jialong Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Flowformer: Linearizing
588 transformers with conservation flows. In *International Conference on Machine Learning*,
589 2022.
- 590 [87] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V
591 Le. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint*
592 *arXiv:1906.08237*, 2019.
- 593 [88] George Zerveas, Srideepika Jayaraman, Dhaval Patel, Anuradha Bhamidipaty, and Carsten
594 Eickhoff. A transformer-based framework for multivariate time series representation learning.
595 In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*,
596 pages 2114–2124, 2021.
- 597 [89] Shaolei Zhang and Yang Feng. Modeling concentrated cross-attention for neural machine
598 translation with Gaussian mixture model. In *Findings of the Association for Computational*
599 *Linguistics: EMNLP 2021*, pages 1401–1411, Punta Cana, Dominican Republic, November
600 2021. Association for Computational Linguistics.
- 601 [90] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer.
602 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16259–
603 16268, 2021.

604 Checklist

605 The checklist follows the references. Please read the checklist guidelines carefully for information on
606 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
607 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
608 the appropriate section of your paper or providing a brief inline description. For example:

- 609 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 610 • Did you include the license to the code and datasets? **[No]** The code and the data are
611 proprietary.
- 612 • Did you include the license to the code and datasets? **[N/A]**

613 Please do not modify the questions and only use the provided macros for your answers. Note that the
614 Checklist section does not count towards the page limit. In your paper, please delete this instructions
615 block and only keep the Checklist section heading above along with the questions/answers below.

- 616 1. For all authors...
 - 617 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
618 contributions and scope? **[Yes]**
 - 619 (b) Did you describe the limitations of your work? **[Yes]** See Section 6
 - 620 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See
621 Section 6
 - 622 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
623 them? **[Yes]**
- 624 2. If you are including theoretical results...
 - 625 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - 626 (b) Did you include complete proofs of all theoretical results? **[Yes]**
- 627 3. If you ran experiments...
 - 628 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
629 mental results (either in the supplemental material or as a URL)? **[Yes]**
 - 630 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
631 were chosen)? **[Yes]**
 - 632 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
633 ments multiple times)? **[Yes]**
 - 634 (d) Did you include the total amount of compute and the type of resources used (e.g., type
635 of GPUs, internal cluster, or cloud provider)? **[Yes]**

- 636 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 637 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 638 (b) Did you mention the license of the assets? [N/A]
- 639 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 640
- 641 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 642 using/curating? [N/A]
- 643 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 644 information or offensive content? [N/A]
- 645 5. If you used crowdsourcing or conducted research with human subjects...
- 646 (a) Did you include the full text of instructions given to participants and screenshots, if
- 647 applicable? [N/A]
- 648 (b) Did you describe any potential participant risks, with links to Institutional Review
- 649 Board (IRB) approvals, if applicable? [N/A]
- 650 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 651 spent on participant compensation? [N/A]