

Transformer with Fourier Integral Attentions

Tan Nguyen^{b,*} Minh Pham^{b,*} Tam Nguyen[†] Khai Nguyen[†] Stanley J. Osher^b Nhat Ho[†]

University of Texas at Austin[†],
University of California, Los Angeles^b,
May 31, 2022

Abstract

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

1 Introduction

Transformers [75] are powerful neural networks that have achieved tremendous success in many areas of machine learning [36, 68, 32] and become the state-of-the-art model on a wide range of applications across different data modalities, from language [20, 1, 15, 10, 54, 4, 7, 18] to images [21, 39, 70, 55, 51, 24], videos [3, 40], point clouds [87, 27], and protein sequence [57, 30]. In addition to their excellent performance on supervised learning tasks, transformers can also effectively transfer the learned knowledge from a pretraining task to new tasks with limited or no supervision [52, 53, 20, 85, 38]. At the core of transformers is the dot-product self-attention, which mainly accounts for the success of transformer models [11, 48, 37]. This dot-product self-attention learn self-alignment between tokens in an input sequence by estimating the relative importance of a given token with respect to all other tokens. It then transform each token into a weighted average of the feature representations of other tokens where the weight is proportional to a importance score between each pair of tokens. The importance scores in self-attention enable a token to attend to other tokens in the sequence, thus capturing the contextual representation [5, 75, 34].

* Tan Nguyen and Minh Pham contributed equally to this work. Correspondence to: Nhat Ho (minhnhat@utexas.edu) and Tan Nguyen (tanmnguyen89@ucla.edu).

1.1 Self-Attention

Given an input sequence $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D_x}$ of N feature vectors, self-attention 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 transformations

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

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 $\mathbf{Q} := [\mathbf{q}_1, \dots, \mathbf{q}_N]^\top$, $\mathbf{K} := [\mathbf{k}_1, \dots, \mathbf{k}_N]^\top$, and $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_N]^\top$, where the vectors $\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i$ for $i = 1, \dots, N$ are the query, key, and value vectors, respectively.

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

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

where the softmax function is applied to each row of the matrix $(\mathbf{Q}\mathbf{K}^\top)/\sqrt{D}$. For each query vector \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 follows

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

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 attention scores, respectively. The self-attention computed by equations (1) and (2) is called the dot-product attention or softmax attention. In our paper, we refer a transformer that uses this attention as the baseline transformer with the dot-product attention or the dot-product transformer. The structure of the attention matrix \mathbf{A} after training governs the ability of the self-attention to capture contextual representation for each token.

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

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

The capacity of the attention mechanism and its ability to learn diverse syntactic and semantic relationships determine the success of transformers [69, 76, 14, 77, 28]. However, equations (1) and (2) implies that the dot-product attention assumes the features (q_{i1}, \dots, q_{iD}) in \mathbf{q}_i , as well as the features (k_{j1}, \dots, k_{jD}) in \mathbf{k}_j , are independent. Thus, the dot-product attention fail to capture the correlations between these features, limiting its representation capacity and inhibit the performance of transformers on practical tasks where there is no guarantee that independent features can learned from complex data. One solution to capture correlations between features \mathbf{q}_i and \mathbf{k}_j is to introduce covariance matrices into the formulation of the dot-product attention with the cost of significantly increasing of the computational complexity. Also, choosing good covariance matrices is difficult.

1.2 Contribution

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

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

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

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

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 space of real-valued functions on \mathbb{R}^D that are integrable. For any two sequences $\{a_N\}_{N \geq 1}, \{b_N\}_{N \geq 1}$, 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.

2 A Nonparametric Regression Interpretation of Self-attention

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

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

where $\varepsilon_1, \dots, \varepsilon_N$ are independent noises such that $\mathbb{E}(\varepsilon_j) = 0$. Furthermore, we consider a random 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 admits p as density function. By an abuse of notation, we also denote p as the joint density where the 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 function and we would like to estimate it.

Nadaraya–Watson estimator: Our approach to estimate the function f is based on Nadaraya–Watson’s nonparametric kernel regression approach [46]. In particular, from the nonparametric 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 estimate the conditional distribution of the value vectors given the key vectors. Given the density function p of the key vectors and the joint density p of the key and value vectors, for any pair of vectors (\mathbf{v}, \mathbf{k}) generate from model (3) we have

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

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

Kernel density estimator: To estimate $p(\mathbf{v}, \mathbf{k})$ and $p(\mathbf{k})$, we employ the kernel density estimation approach [58, 49]. In particular, by using the isotropic Gaussian kernel with bandwidth σ , we have the following estimators of $p(\mathbf{v}, \mathbf{k})$ and $p(\mathbf{k})$:

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

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

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

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

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

If we further assume that the keys \mathbf{k}_j are normalized, which is usually done in practice to stabilize the training of transformers [63], the value of $\hat{f}_\sigma(\mathbf{q}_i)$ in equation (6) then becomes

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

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

Remark 1. *The assumption that \mathbf{k}_j is normalized is to recover the pairwise dot-product attention in transformers. In general, this assumption is not necessary. In fact, the isotropic Gaussian kernel in equation (7) is more desirable than the dot-product kernel in equation (8) of the pairwise dot-product attention since the former is Lipschitz while the later is not Lipschitz [33]. The Lipschitz constraint helps improve the robustness of the model [13, 73, 2] and stabilize the model training [44].*

Limitation of Self-Attention: From our nonparametric regression interpretation, self-attention is derived from the use of isotropic Gaussian kernels for kernel density estimation and nonparametric regression estimation, which may fail to capture the complex correlations between D features in \mathbf{q}_i and \mathbf{k}_j [80, 29]. Using multivariate Gaussian kernels with dense covariance matrices can help capture such correlations; however, choosing good covariance matrices is challenging and inefficient [79, 65, 9]. In the following section, we discuss the Fourier integral estimator and its use as a kernel for computing self-attention in order to overcome these limitations.

3 FourierFormer: Transformer via Generalized Fourier Integral Theorem

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

3.1 (Generalized) Fourier Integral Theorems and Their Applications

The Fourier integral theorem is a beautiful result in mathematics [84, 6] and has been recently used in nonparametric mode clustering, deconvolution problem, and generative modeling [29]. It is a combination of Fourier transform and Fourier inverse transform. In particular, for any function $p \in \mathbb{L}_1(\mathbb{R}^D)$, the *Fourier integral theorem* is given by

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

where $\mathbf{k} = (k_1, \dots, k_D)$ and $\mathbf{y} = (y_1, \dots, y_D)$. Equation (9) suggests that

$$p_R(\mathbf{k}) := \frac{1}{\pi^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \frac{\sin(R(y_j - k_j))}{(y_j - k_j)} p(\mathbf{y}) d\mathbf{y}$$

can be used as an estimator of the function p .

Benefits of the Fourier integral over Gaussian kernel: There are two important benefits of the estimator p_R : (i) it can automatically preserve the correlated structure lying within p even when p is very complex and high dimensional function. It is in stark contrast to the standard kernel estimator built based on multivariate Gaussian kernel where we need to choose good covariance matrix in the multivariate Gaussian kernel to guarantee such estimator to work well. We note that as the standard soft-max Transformer is constructed based on the multivariate Gaussian kernel, the issue of choosing good covariance matrix in dot-product transformer is inevitable; (ii) The product of sinc kernels in the estimator p_R does not decay to a point mass when $R \rightarrow \infty$. It is in stark difference from the multivariate Gaussian kernel estimator, which converges to a point mass when the covariance matrix goes to 0. It indicates that p_R is a non-trivial estimator of the function p . Finally, detailed illustrations of these benefits of the Fourier integral over Gaussian kernel in density estimation and nonparametric regression problems, which we have just shown to have connection to the self-attention in transformer, can be found in Section 8 in [29].

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

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

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 $\mathbf{k} \in \mathbb{R}^D$, the generalized Fourier integral estimator p_R^ϕ becomes the Fourier integral estimator p_R . Under appropriate conditions on the function ϕ (see Theorem 1 in Section 3.1.1 and Theorem 3 in Appendix A.1), the estimator p_R^ϕ converges to the true function p , namely,

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

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

3.1.1 Density Estimation via Generalized Fourier Integral Theorems

We first apply the generalized Fourier integral theorem to the density estimation problem. To ease the 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 density function p where $D \geq 1$ is the dimension. Inspired by the generalized Fourier integral theorem, we obtain the following *generalized Fourier density estimator* $p_{N,R}^\phi$ of p as follows:

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

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 the generalized Fourier density estimator $p_{n,R}^\phi$ and the true density p , we utilize mean integrated squared errors (MISE) [83], which is given by:

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

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

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

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

Proof of Theorem 1 is in Appendix B.1. A few comments are in order. First, by choosing R to balance the bias and variance in the bound of MISE in Theorem 1, we have the optimal R as $R = \mathcal{O}(N^{1/(D+m+1)})$. With that choice of R , the MISE rate of $p_{N,R}^\phi$ is $\mathcal{O}(N^{-(m+1)/(D+m+1)})$. Second, when $\phi(z) = z^l$ for $l \geq 4$ and $z \in \mathbb{R}$, the assumptions in Theorem 1 are satisfied when $m = 1$. Under this case, the MISE rate of $p_{N,R}^\phi$ is $\mathcal{O}(N^{-2/(D+2)})$. However, these assumptions do not satisfy when $\phi(z) = z^l$ and $l \in \{1, 2, 3\}$, which is due to the limitation of the current proof technique of Theorem 1 that is based on Taylor expansion of the estimator $p_{n,R}^\phi$.

To address the limitation of the Taylor expansion technique, we utilize the Plancherel theorem in 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 the theoretical analyses for such setting are in Appendix A.

3.2 FourierFormer: Transformers with Fourier Attentions

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

3.2.1 Nonparametric Regression via Generalized Fourier Integral Theorem

We now discuss an application of the generalized Fourier integral theorems to the nonparametric 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 following nonparametric regression model:

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

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 i.i.d. samples from p . Given the generalized Fourier density estimator (12), following the argument in Section 2, the Nadaraya–Watson estimator of the function f based on the generalized Fourier density estimator is given by:

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

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

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

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 $\int_{\mathbb{R}} \left|\phi\left(\frac{\sin(z)}{z}\right)\right| |z|^j dz < \infty$ 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 C_1, C_2, C_3, C_4 such that the following holds:*

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

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 respect to the key vectors $\mathbf{k}_1, \dots, \mathbf{k}_N$ and the noises $\varepsilon_1, \dots, \varepsilon_N$.

Proof of Theorem 2 is in Appendix B.3. A few comments with Theorem 2 are in order. First, by choosing R to balance the bias and variance in the bound of the MSE of the nonparametric generalized 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 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 assumption on the function ϕ of Theorem 2 is satisfied with $m = 1$. Under this case, the rate of $f_{N,R}$ becomes $\mathcal{O}(N^{-\frac{4}{D+4}})$. In Appendix A, we also provide the rate of $f_{N,R}$ when $\phi(z) = z^l$ for some $l \leq 5$, which includes the original Fourier integral theorem.

3.2.2 FourierFormer

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

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

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

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

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

Remark 2 (The Nonnegativity of the Fourier Kernel). *The density estimation via generalized Fourier integral theorem in Section 3.1.1 does not require the generalized Fourier density estimator to be nonnegative. However, empirically, we observe that negative density estimator can cause instability in training the FourierFormer. Thus, in FourierFormer, we choose the function ϕ to be a nonnegative function to enforce the density estimator to be nonnegative. In particular, we choose ϕ to be power functions of the form $\phi(x) = x^{2m}$, where m is an positive integer. Note that when $m = 2$ and $m = 4$, the kernels in our generalized Fourier integral estimators are the well-known Fejer-de la Vallee Poussin and Jackson-de la Vallee Poussin kernels [17].*

3.3 An Efficient Implementation of the Fourier Attention

The Fourier kernel is implemented efficiently in the C++/CUDA extension developed by Pytorch [50]. The idea is similar to the function `cdist` [50], 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 :

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

We directly implement d_f as a `torch.autograd.Function` [50] in which we provide an efficient way to compute forward and backward function (d_f and gradient of d_f). While the implementation of the forward function is straight forward, the backward function is more tricky since we need to 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 develop the backward function with highly parallel computation by exploiting GPU architecture and utilizing the reduction technique. The computational time is comparable to function `cdist`; thus, our FourierFormer implementation is as computationally time-efficient.

4 Experimental Results

In this section, we numerically justify the advantage of FourierFormer over the baseline dot-product transformer on two large-scale tasks: language modeling on WikiText-103 [42] (Section 4.1) and image classification on ImageNet [19, 59] (Section 4.2). We aim to show that: (i) FourierFormer achieves better accuracy than the baseline transformer on a variety of practical tasks with different data modalities, and (ii) FourierFormer helps reduce head redundancy compared to the baseline transformer (Section 4.3).

Throughout the section, we compare FourierFormers with the baseline dot-product transformers of the same configuration. In all experiments, we made the constant R in Fourier attention (see equation (54)) to be a learnable scalar and set choose the function $\phi(x) = x^4$ (see Remark 2). All of our results are averaged over 5 runs with different seeds. More details on the models and training are provided in Appendix C. We also provide additional experimental results in Appendix D.

4.1 Language Modeling on WikiText-103

Datasets and metrics WikiText-103 is a collection of articles from Wikipedia, which have long contextual dependencies. The training set consists of about 28K articles containing

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

103M running words; this corresponds to text blocks of about 3600 words. The validation and test sets have 218K and 246K running words, respectively. Each of them contains 60 articles and about 268K words. Our experiment follows the standard setting [42, 63] and splits the training data into L -word independent long segments. For evaluation, we use a batch size of 1, and process the text sequence with a sliding window of size L . The last position is used for computing perplexity (PPL) except in the first segment, where all positions are evaluated as in [1, 63].

Models and baselines: Our implementation is based on the public code by [63]. We use their small and medium models in our experiments. In particular, for small models, the key, value, and query dimension are set to 128, and the training and evaluation context length are set to 256. For medium models, the key, value, and query dimension are set to 256, and the training and evaluation context length are set to 384. In both configurations, the number of heads is 8, the feed-forward layer dimension is 2048, and the number of layers is 16.

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

4.2 Image Classification on ImageNet

Datasets and metrics The ImageNet dataset [19, 59] consists of 1.28M training images and 50K validation images. For this benchmark, the model learns to predict the category of the input image among 1000 categories. Top-1 and top-5 classification accuracies are reported.

Models and baselines: We use the DeiT-tiny model [71] with 12 transformer layers, 4 attention heads per layer, and the model dimension of 192. To train the models, we follow the same setting and configuration as for the baseline [71].

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

Implementation 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

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

4.3 FourierFormer Helps Reducing Head Redundancy

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

5 Related Work

Interpretation of Attention Mechanism in Transformers: Recent works have tried to gain an understanding of transformer’s attention from different perspectives. [72] considers attention as applying kernel smoother over the inputs. Extending this kernel approach, [31, 12, 81] linearize the softmax kernel in dot-product attention and propose a family of efficient transformers with linear computational and memory complexity. [8] then shows that these linear transformers are comparable to a Petrov-Galerkin projection [56], suggesting that the softmax normalization in the dot-product attention is sufficient but not necessary. Other works provide an understanding of attention in transformers via ordinary/partial differential equation include [41, 61]. In addition, [67, 26, 86, 47] relate attentions in transformers to a Gaussian mixture models. Several works also connect the attention mechanism to graph-structured learning and message passing in graphical models [82, 64, 35]. Our work focuses on deriving the connection between self-attention and nonparametric kernel regression and exploring better regression estimator, such as the generalized Fourier nonparametric regression

estimator, to improve the performance of transformers.

Redundancy in Transformers: [16, 43, 22] show that neurons and attention heads in the pre-trained transformer are redundant and can be removed when applied on a downstream task. By studying the contextualized embeddings in pre-trained networks, it has been demonstrated that the learned representations from these redundant models are highly anisotropic [45, 23]. Furthermore, [62, 66, 78, 60] employ knowledge distillation and sparse approximation to enhance the efficiency of transformers. Our FourierFormer is complementary to these methods and can be combined with them.

6 Concluding Remarks

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

Supplement to “FourierFormer: Transformer Meets Generalized Fourier Integral Theorem”

In the supplementary material, we collect proofs, additional theories, and experiment results deferred from the main text. In Appendix A, we provide additional theoretical results for generalized Fourier density estimator and for generalized Fourier nonparametric regression estimator. We provide proofs of key results in the main text and additional theories in Appendix B. We present experiment details in Appendix C while including additional experimental results in Appendix D.

A Additional Theoretical Results

In this section, we provide additional theoretical results for generalized Fourier density estimator in Appendix A.1 and for generalized Fourier nonparametric regression estimator in Appendix A.2.

A.1 Generalized Fourier density estimator

We now establish the MISE rate of $p_{N,R}^\phi$ in equation (12) when $\phi(z) = z^l$ and $l \in \{1, 2\}$. We consider the following tail bounds on the Fourier transform of the true density function p as follows.

Definition 3. (1) We say that p is supersmooth of order α if we have universal constants C_1 and C_2 such that the following inequalities hold for almost surely $x \in \mathbb{R}^D$:

$$|\widehat{p}(x)| \leq C_1 \exp \left(-C_2 \left(\sum_{j=1}^D |x_j|^\alpha \right) \right).$$

Here, \widehat{p} denotes the Fourier transform of the function p .

(2) The function p is ordinary smooth of order β if there exists universal constant c such that the following inequality holds for almost surely $x \in \mathbb{R}^D$:

$$|\widehat{p}(x)| \leq c \cdot \prod_{j=1}^D \frac{1}{(1 + |x_j|^\beta)}.$$

The notions of supersmoothness and ordinary smoothness had been used widely in deconvolution problems [25] and density estimation problems [17, 74, 29]. The supersmooth condition is satisfied when the function p is Gaussian distribution or Cauchy distribution while the ordinary smooth condition is satisfied when the function p is Laplace distribution and Beta distribution.

Based on the smoothness conditions in Definition 3, we have the following result regarding the mean-square integrated error (MISE) of the function generalized Fourier density estimator (12) (see equation (13) for a definition of MISE) when $\phi(z) = z^l$ and $l \in \{1, 2\}$.

Theorem 3. (a) When $\phi(z) = z$, the following holds:

- (Supersmooth setting) If the true density function p is supersmooth function of order α for some $\alpha > 0$, then there exists universal constants \bar{C}_1, \bar{C}_2 , and \bar{C}_3 such that as long as $R \geq \bar{C}_1$ we have

$$MISE(p_{N,R}^\phi) \leq \bar{C}_2 \left(R^{\max\{1-\alpha, 0\}} \exp(-\bar{C}_3 R^\alpha) + \frac{R^D}{N} \right).$$

- (Ordinary smooth setting) If the true density function p is ordinary smooth function of order β for some $\beta > 1$, then there exists universal constants \bar{c} such that

$$MISE(p_{N,R}^\phi) \leq \bar{c} \left(R^{-\beta+1} + \frac{R^D}{N} \right).$$

(b) When $\phi(z) = z^2$, the following holds

- (Supersmooth setting) If the true density function p is supersmooth function of order α for some $\alpha > 0$, then there exists universal constants C'_1 and C'_2 such that as long as $R \geq C'_1$ we have

$$MISE(p_{N,R}^\phi) \leq C'_2 \left(\frac{1}{R^2} + \frac{R^D}{N} \right).$$

- (Ordinary smooth setting) If the true density function p is ordinary smooth function of order β for some $\beta > 3$, then there exists universal constants c' such that

$$MISE(p_{N,R}^\phi) \leq c' \left(\frac{1}{R^2} + \frac{R^D}{N} \right).$$

Proof of Theorem 3 is in Appendix B.2. A few comments with the results of Theorem 3 are in order.

When $\phi(z) = z$: As part (a) of Theorem 3 indicates, when the function p is supersmooth, by choosing the radius R to balance the bias and variance, we have the optimal R as $R = \left(\frac{\log(N)}{C_3}\right)^{1/\alpha}$ and the MISE rate of the generalized Fourier density estimator $p_{N,R}^\phi$ becomes $\mathcal{O}\left(\frac{\log(N)^{D/\alpha}}{N}\right)$. It indicates that, the MISE rate of $p_{N,R}^\phi$ is parametric when the function p is supersmooth. On the other hand, when the function p is ordinary smooth, the optimal R becomes $\mathcal{R} = \mathcal{O}(N^{\frac{1}{D+\beta-1}})$ and the MISE rate becomes $\mathcal{O}\left(N^{-\frac{\beta-1}{D+\beta-1}}\right)$. It is slower than the MISE rate when the function p is supersmooth.

When $\phi(z) = z^2$: The results of part (b) of Theorem 3 demonstrate that the upper bounds for the MISE rate of the generalized Fourier density estimator $p_{N,R}^\phi$ is similar for both the supersmooth and ordinary smooth settings. The optimal radius $R = \mathcal{O}\left(N^{\frac{1}{D+2}}\right)$ and the MISE rate of the estimator is $\mathcal{O}\left(N^{-\frac{2}{D+2}}\right)$.

A.2 Generalized Fourier nonparametric regression estimator

In this appendix, we provide additional result for the mean square error (MSE) rate of the generalized Fourier nonparametric regression estimator $f_{N,R}$ in equation (14) when $\phi(z) = z$, namely, the setting of the Fourier integral theorem. The results when $\phi(z) = z^l$ for $l \in \{2, 3, 4, 5\}$ are left for the future work.

When $\phi(z) = z$, the MSE rate of $f_{N,R}$ had been established in Theorem 9 of Ho et al. [29] when the function p is supersmooth function. Here, we restate that result for the completeness.

Theorem 4. *Assume that the function p is supersmooth function of order α for some $\alpha > 0$ and $\sup_{\mathbf{k} \in \mathbb{R}^D} |p(\mathbf{k})| < \infty$. Furthermore, we assume that the function f in the nonparametric regression model (3) is such that $\sup_{\mathbf{k} \in \mathbb{R}^D} |f^2(\mathbf{k})p(\mathbf{k})| < \infty$ and*

$$|\widehat{f \cdot p}(\mathbf{t})| \leq C_1 Q(|t_1|, |t_2|, \dots, |t_D|) \exp\left(-C_2 \left(\sum_{j=1}^D |t_j|^\alpha\right)\right),$$

where $\widehat{f \cdot p}(\mathbf{t})$ is the Fourier transform of the function $f \cdot p$, C_1 and C_2 are some universal constants, and $Q(|t_1|, |t_2|, \dots, |t_D|)$ is some polynomial function of $|t_1|, \dots, |t_D|$ with non-negative coefficients. Then, we can find universal constants C_3, C_4, C_5 such that as long as $R \geq C_3$ we have

$$\mathbb{E} [(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2] \leq C_4 \frac{R^{\max\{2\deg(Q)+2-2\alpha, 0\}} \exp(-2C_2 R^\alpha) + \frac{(f(\mathbf{k})+C_5)R^D}{N}}{p^2(\mathbf{k})\bar{J}(R)},$$

where $\deg(Q)$ denotes the degree of the polynomial function Q , and we define $\bar{J}(R) = 1 - \frac{R^{\max\{2-2\alpha, 0\}} \exp(-2C_2 R^\alpha) + \frac{R^D \log(NR)}{N}}{p^2(\mathbf{k})}$.

Proof of Theorem 4 is similar to the proof of Theorem 9 of Ho et al. [29]; therefore, it is omitted. The result of Theorem 4 indicates that the optimal radius $R = \left(\frac{\log(N)}{2C_2}\right)^{1/\alpha}$ and the MSE rate of the generalized Fourier nonparametric regression estimator $f_{N,R}$ is $\mathcal{O}\left(\frac{\log(N)^{D/\alpha}}{N}\right)$.

B Proofs

In this Appendix, we provide proofs for key results in the paper and in Appendix A.

B.1 Proof of Theorem 1

Recall that, $\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_N \in \mathbb{R}^D$ are i.i.d. samples from the density function p . In equation (12), the generalized Fourier density estimator of p_0 is given by:

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

where $A = \int_{\mathbb{R}} \phi\left(\frac{\sin(z)}{z}\right) dz$, $\mathbf{k}_i = (k_{i1}, \dots, k_{iD})$, and $\mathbf{k} = (k_1, \dots, k_D)$. Direct calculation demonstrates that

$$\begin{aligned} \mathbb{E}[p_{N,R}^\phi(\mathbf{k})] &= \frac{R^D}{A^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(R(k_j - y_j))}{R(k_j - y_j)}\right) p(\mathbf{y}) d\mathbf{y} \\ &= \frac{1}{A^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(y_j)}{y_j}\right) p\left(\mathbf{k} - \frac{\mathbf{y}}{R}\right) d\mathbf{y}. \end{aligned} \quad (16)$$

An application of Taylor expansion up to the m -th order indicates that

$$p\left(\mathbf{k} - \frac{\mathbf{y}}{R}\right) = \sum_{0 \leq |\alpha| \leq m} \frac{1}{R^{|\alpha|} \alpha!} \prod_{j=1}^D (-y_j)^{\alpha_j} \frac{\partial^{|\alpha|} p}{\partial \mathbf{k}^\alpha}(\mathbf{k}) + \bar{R}(\mathbf{k}, \mathbf{y}), \quad (17)$$

where $\alpha = (\alpha_1, \dots, \alpha_d)$, $|\alpha| = \sum_{j=1}^d \alpha_j$, and $\bar{R}(\mathbf{k}, \mathbf{y})$ is Taylor remainder admitting the following form:

$$\bar{R}(\mathbf{k}, \mathbf{y}) = \sum_{|\beta|=m+1} \frac{m+1}{R^{m+1} \beta!} \prod_{j=1}^D (-y_j)^{\beta_j} \int_0^1 (1-t)^m \frac{\partial^{m+1} p}{\partial \mathbf{k}^\beta} \left(\mathbf{k} - \frac{t\mathbf{y}}{R}\right) dt. \quad (18)$$

Plugging equations (17) and (18) into equation (16), we find that

$$\begin{aligned} \mathbb{E}[p_{N,R}^\phi(\mathbf{k})] &= p(\mathbf{k}) + \frac{1}{A^D} \sum_{1 \leq |\alpha| \leq m} \frac{1}{R^{|\alpha|} \alpha!} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(y_j)}{y_j}\right) \prod_{j=1}^D (-y_j)^{\alpha_j} \frac{\partial^{|\alpha|} p}{\partial \mathbf{k}^\alpha}(\mathbf{k}) d\mathbf{y} \\ &+ \frac{1}{A^D} \sum_{|\beta|=m+1} \frac{m+1}{R^{m+1} \beta!} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(y_j)}{y_j}\right) \prod_{j=1}^D (-y_j)^{\beta_j} \int_0^1 (1-t)^m \frac{\partial^{m+1} p_0}{\partial \mathbf{k}^\beta} \left(\mathbf{k} - \frac{t\mathbf{y}}{R}\right) d\mathbf{y} dt. \end{aligned}$$

According to the hypothesis that $\int_{\mathbb{R}} \phi\left(\frac{\sin(z)}{z}\right) z^j dz = 0$ for all $1 \leq j \leq m$, we obtain that

$$\int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(y_j)}{y_j}\right) \prod_{j=1}^D (-y_j)^{\alpha_j} \frac{\partial^{|\alpha|} p}{\partial \mathbf{k}^\alpha}(\mathbf{k}) d\mathbf{y} = 0$$

for any $\alpha = (\alpha_1, \dots, \alpha_d)$ such that $1 \leq |\alpha| \leq m$. Collecting the above results, we arrive at

$$\begin{aligned} & |\mathbb{E}[p_{N,R}^\phi(\mathbf{k})] - p(\mathbf{k})| \\ &= \left| \frac{1}{A^D} \sum_{|\beta|=m+1} \frac{m+1}{R^{m+1}\beta!} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(y_j)}{y_j}\right) \prod_{j=1}^D (-y_j)^{\beta_j} \int_0^1 (1-t)^m \frac{\partial^{m+1} p}{\partial \mathbf{k}^\beta} \left(\mathbf{k} - \frac{t\mathbf{y}}{R}\right) d\mathbf{y} dt \right| \\ &\leq \frac{1}{A^D} \sum_{|\beta|=m+1} \frac{m+1}{R^{m+1}\beta!} \int_{\mathbb{R}^D} \prod_{j=1}^D \left| \phi\left(\frac{\sin(y_j)}{y_j}\right) \right| \prod_{j=1}^D |y_j|^{\beta_j} \int_0^1 (1-t)^m \left| \frac{\partial^{m+1} p}{\partial \mathbf{k}^\beta} \left(\mathbf{k} - \frac{t\mathbf{y}}{R}\right) \right| d\mathbf{y} dt. \end{aligned}$$

Since the function $p \in \mathcal{C}^{m+1}(\mathbb{R}^D)$, we can find positive constant M such that $\|\frac{\partial^{m+1} p}{\partial \mathbf{k}^\beta}(\mathbf{k})\|_\infty \leq M$ for all $\beta = (\beta_1, \dots, \beta_d)$ such that $|\beta| = m+1$. Therefore, we find that

$$\begin{aligned} |\mathbb{E}[p_{N,R}^\phi(\mathbf{k})] - p(\mathbf{k})| &\leq \frac{M}{A^D} \sum_{|\beta|=m+1} \frac{m+1}{R^{m+1}\beta!} \int_{\mathbb{R}^D} \prod_{j=1}^D \left| \phi\left(\frac{\sin(y_j)}{y_j}\right) \right| \prod_{j=1}^D |y_j|^{\beta_j} d\mathbf{y} \int_0^1 (1-t)^m dt \\ &= \frac{M}{A^D} \sum_{|\beta|=m+1} \frac{1}{R^{m+1}\beta!} \int_{\mathbb{R}^D} \prod_{j=1}^D \left| \phi\left(\frac{\sin(y_j)}{y_j}\right) \right| \prod_{j=1}^D |y_j|^{\beta_j} d\mathbf{y}. \end{aligned}$$

For any $\beta = (\beta_1, \dots, \beta_d)$ such that $|\beta| = m+1$, an application of the AM-GM inequality indicates that $\prod_{j=1}^D |y_j|^{\beta_j} \leq m(\sum_{j=1}^D |y_j|^{m+1})$. Hence, putting these results together leads to

$$|\mathbb{E}[p_{N,R}^\phi(\mathbf{k})] - p(\mathbf{k})| \leq \frac{Mm}{A^D R^{m+1}} \sum_{|\beta|=m+1} \frac{1}{\beta!} \int_{\mathbb{R}^D} \prod_{j=1}^D \left| \phi\left(\frac{\sin(y_j)}{y_j}\right) \right| \left(\sum_{j=1}^D |y_j|^{m+1} \right) d\mathbf{y}.$$

From the hypothesis, we have $\int_{\mathbb{R}} \left| \phi\left(\frac{\sin(z)}{z}\right) \right| |z|^{m+1} dz < \infty$. As a consequence, we can find a universal constant C depending on A and d such that

$$|\mathbb{E}[p_{n,R}^\phi(\mathbf{k})] - p(\mathbf{k})| \leq \frac{C}{R^{m+1}}$$

for all $\mathbf{k} \in \mathbb{R}^D$.

Bounding the variance: We now move to bound the variance of $p_{N,R}^\phi(\mathbf{k})$. Indeed, direct computation indicates that

$$\begin{aligned} \text{Var}[p_{N,R}^\phi(\mathbf{k})] &= \frac{R^{2D}}{nA^{2D}} \text{Var} \left[\prod_{j=1}^D \phi\left(\frac{\sin(R(k_j - K_j))}{R(x_j - K_j)}\right) \right] \\ &\leq \frac{R^{2D}}{nA^{2D}} \mathbb{E} \left[\prod_{j=1}^D \phi^2\left(\frac{\sin(R(k_j - K_j))}{R(k_j - K_j)}\right) \right] \\ &= \frac{R^D}{nA^{2D}} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi^2\left(\frac{\sin(y_j)}{y_j}\right) p\left(\mathbf{k} - \frac{\mathbf{y}}{R}\right) d\mathbf{y} \leq \frac{R^D \|p\|_\infty}{NA^{2D}} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi^2\left(\frac{\sin(y_j)}{y_j}\right) d\mathbf{y} \end{aligned}$$

where the variance and the expectation are taken with respect to $K = (K_1, \dots, K_d) \sim p$. As $\int_{\mathbb{R}} \phi^2\left(\frac{\sin(z)}{z}\right) dz < \infty$, there exists a universal constant C' depending on A and D such that

$$\text{Var}[p_{N,R}^\phi(\mathbf{k})] \leq \frac{C' R^D}{N}.$$

As a consequence, we obtain the conclusion of the theorem.

B.2 Proof of Theorem 3

From the Plancherel theorem, we obtain that

$$\int_{\mathbb{R}^D} \left[(p_{N,R}^\phi(\mathbf{k}) - p(\mathbf{k})) \right]^2 d\mathbf{k} = \frac{1}{(2\pi)^D} \int_{\mathbb{R}^D} \left[\widehat{p}_{N,R}^\phi(\mathbf{t}) - \widehat{p}(\mathbf{t}) \right]^2 dt, \quad (19)$$

where $\widehat{p}_{N,R}^\phi$ and \widehat{p} are respectively the Fourier transforms of $p_{N,R}$ and p . From the definition of generalized Fourier density estimator $p_{N,R}^\phi$ in equation (12), it is clear that

$$\widehat{p}_{N,R}^\phi(\mathbf{t}) = \frac{1}{N} \sum_{i=1}^N \exp(i\mathbf{t}^\top \mathbf{k}_i) \prod_{j=1}^D K_R(t_j),$$

for any $\mathbf{t} = (t_1, \dots, t_D) \in \mathbb{R}^D$ where we define $K_R(y) := \frac{1}{\pi} \int_{\mathbb{R}} R \phi\left(\frac{\sin(Rx)}{Rx}\right) \exp(iyx) dx$ for any $y \in \mathbb{R}$. To ease the presentation, we denote $\bar{K}_R(\mathbf{t}) := \prod_{j=1}^D K_R(t_j)$ and $\varphi_N(\mathbf{t}) = \frac{1}{N} \sum_{i=1}^N \exp(i\mathbf{t}^\top \mathbf{k}_i)$ for any $\mathbf{t} = (t_1, t_2, \dots, t_D) \in \mathbb{R}^D$. Based on these notations, we can rewrite

$$\widehat{p}_{N,R}^\phi(\mathbf{t}) = \varphi_N(\mathbf{t}) \bar{K}_R(\mathbf{t})$$

Direct calculation shows that $\mathbb{E}_{\mathbf{k}_1^N}[\varphi_N(\mathbf{t})] = \widehat{p}(\mathbf{t})$ for any $\mathbf{t} \in \mathbb{R}^D$ where $\mathbf{k}_1^N := (\mathbf{k}_1, \dots, \mathbf{k}_n)$. Furthermore, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{k}_1^N} [|\varphi_N(\mathbf{t})|^2] &= \mathbb{E}[\varphi_N(\mathbf{t})\varphi_N(-\mathbf{t})] = \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N \exp(i\mathbf{t}^\top \mathbf{k}_i) \right) \left(\frac{1}{N} \sum_{i=1}^N \exp(-i\mathbf{t}^\top \mathbf{k}_i) \right) \right] \\ &= \frac{1}{N} + \frac{(N-1)}{N} \mathbb{E} \left[\exp(i\mathbf{t}^\top \mathbf{k}) \exp(-i\mathbf{t}^\top \mathbf{k}) \right] \\ &= \frac{1}{N} + \frac{(N-1)}{N} |\widehat{p}(\mathbf{t})|^2. \end{aligned}$$

Collecting the above results, we have the following equations:

$$\begin{aligned} \mathbb{E}_{\mathbf{k}_1^N} \left[\int_{\mathbb{R}^D} \left[\widehat{p}_{N,R}^\phi(\mathbf{t}) - \widehat{p}(\mathbf{t}) \right]^2 dt \right] &= \mathbb{E}_{\mathbf{k}_1^N} \left[\int_{\mathbb{R}^D} [\varphi_N(\mathbf{t}) \bar{K}_R(\mathbf{t}) - \widehat{p}(\mathbf{t})]^2 dt \right] \\ &= \mathbb{E}_{\mathbf{k}_1^N} \left[\int_{\mathbb{R}^D} [(\varphi_N(\mathbf{t}) - \widehat{p}(\mathbf{t})) \bar{K}_R(\mathbf{t}) - \widehat{p}(\mathbf{t})(1 - \bar{K}_R(\mathbf{t}))]^2 dt \right] \\ &= \int_{\mathbb{R}^D} \mathbb{E}_{\mathbf{k}_1^N} [(\varphi_N(\mathbf{t}) - \widehat{p}(\mathbf{t}))^2] \bar{K}_R^2(\mathbf{t}) + \widehat{p}^2(\mathbf{t})(1 - \bar{K}_R(\mathbf{t}))^2 dt \\ &= \int_{\mathbb{R}^D} \widehat{p}^2(\mathbf{t})(1 - \bar{K}_R(\mathbf{t}))^2 dt + \frac{1}{N} \int_{\mathbb{R}^D} (1 - |\widehat{p}(\mathbf{t})|^2) \bar{K}_R^2(\mathbf{t}) dt. \end{aligned} \quad (20)$$

Combining the results from equations (19) and (20), we find that

$$\begin{aligned} \text{MISE}(p_{N,R}^\phi) &= \mathbb{E}_{\mathbf{k}_1^N} \left[\int_{\mathbb{R}^D} \left[(p_{N,R}^\phi(\mathbf{k}) - p(\mathbf{k})) \right]^2 d\mathbf{k} \right] \\ &= \frac{1}{(2\pi)^D} \left(\int_{\mathbb{R}^D} \widehat{p}^2(\mathbf{t})(1 - \bar{K}_R(\mathbf{t}))^2 dt + \frac{1}{N} \int_{\mathbb{R}^D} (1 - |\widehat{p}(\mathbf{t})|^2) \bar{K}_R^2(\mathbf{t}) dt \right). \end{aligned} \quad (21)$$

B.2.1 When $\phi(z) = z$

We first consider the setting when $\phi(z) = z$, namely, the setting of the Fourier integral theorem. Under this setting, direct computation indicates that

$$\bar{K}_R(\mathbf{t}) = \prod_{i=1}^d \mathbf{1}_{\{|t_i| \leq R\}}.$$

Given the smoothness assumptions on the function p , we have two settings on that function.

Supersmooth setting of the function p : When the function p is supersmooth density, we have

$$|\hat{p}(\mathbf{t})| \leq C_1 \exp \left(-C_2 \left(\sum_{j=1}^D |t_j|^\alpha \right) \right),$$

where C_1 and C_2 are some universal constants. Therefore, we find that

$$\begin{aligned} \int_{\mathbb{R}^D} \hat{p}^2(\mathbf{t})(1 - \bar{K}_R(\mathbf{t}))^2 d\mathbf{t} &= \int_{\mathbb{R}^D \setminus [-R, R]^D} \hat{p}^2(\mathbf{t}) d\mathbf{t} \leq C_1 \int_{\mathbb{R}^D \setminus [-R, R]^D} \exp \left(-C_2 \left(\sum_{j=1}^D |t_j|^\alpha \right) \right) d\mathbf{t} \\ &\leq C_1 \sum_{i=1}^D \int_{B_i} \exp \left(-C_2 \left(\sum_{j=1}^D |t_j|^\alpha \right) \right) d\mathbf{t}, \end{aligned} \tag{22}$$

where $B_i := \{\mathbf{t} \in \mathbb{R}^D : |t_i| \geq R\}$. We now proceed to bound $\int_{B_i} \exp \left(-C_2 \left(\sum_{j=1}^D |t_j|^\alpha \right) \right) d\mathbf{t}$ for all $i \in [D]$. Indeed, we have that

$$\begin{aligned} \int_{B_i} \exp \left(-C_2 \left(\sum_{j=1}^D |t_j|^\alpha \right) \right) d\mathbf{t} &= \left(\int_{\mathbb{R}} \exp(-C_2|x|^\alpha) dx \right)^{D-1} \cdot \int_{|x| \geq R} \exp(-C_2|x|^\alpha) dx \\ &= \frac{C_2 \alpha^{D-1}}{(2C_2 \Gamma(1/\alpha))^{D-1}} \cdot \int_{|x| \geq R} \exp(-C_2|x|^\alpha) dx. \end{aligned}$$

When $\alpha \geq 1$, we have that

$$\int_R^\infty \exp(-C_2 x^\alpha) dx \leq \int_R^\infty x^{\alpha-1} \exp(-C_2 x^\alpha) dx = \exp(-C_2 R^\alpha) / (C_2 \alpha).$$

When $\alpha \in (0, 1)$, then we find that

$$\begin{aligned} \int_R^\infty \exp(-C_2 x^\alpha) dx &= \int_R^\infty x^{1-\alpha} x^{\alpha-1} \exp(-C_2 x^\alpha) dx \\ &\leq \frac{R^{1-\alpha} \exp(-C_2 R^\alpha)}{C_2 \alpha} + \frac{1-\alpha}{C_2 \alpha R^\alpha} \int_R^\infty \exp(-C_2 x^\alpha) dx, \end{aligned}$$

When the R is such that $R^\alpha \geq \frac{2(1-\alpha)}{C_2 \alpha}$, the above inequality becomes

$$\int_R^\infty \exp(-C_2 x^\alpha) dx \leq \frac{2R^{1-\alpha} \exp(-C_2 R^\alpha)}{C_2 \alpha}.$$

Collecting the above results, we arrive at

$$\int_{|x| \geq R} \exp(-C_2|x|^\alpha) dx \leq \frac{4R^{\max\{1-\alpha, 0\}}}{C_2\alpha} \exp(-C_2R^\alpha). \quad (23)$$

Plugging the inequality (23) into the inequality (26), there exists universal constant C_3 depending on α and D such that

$$\int_{\mathbb{R}^D} \widehat{p}^2(\mathbf{t})(1 - \bar{K}_R(\mathbf{t}))^2 d\mathbf{t} \leq C_3 R^{\max\{1-\alpha, 0\}} \exp(-C_1R^\alpha). \quad (24)$$

On the other hand, we also have

$$\frac{1}{N} \int_{\mathbb{R}^D} (1 - |\widehat{p}(\mathbf{t})|^2) \bar{K}_R^2(\mathbf{t}) d\mathbf{t} \leq \frac{1}{N} \int_{\mathbb{R}^D} \bar{K}_R^2(\mathbf{t}) \leq \frac{R^D}{N}. \quad (25)$$

Combining the results from equations (24) and (25), we obtain that

$$\text{MISE}(p_{N,R}^\phi) \leq C_4 \left(R^{\max\{1-\alpha, 0\}} \exp(-C_1R^\alpha) + \frac{R^D}{N} \right).$$

As a consequence, we obtain the conclusion of Theorem 3 under the supersmooth setting of the function p and $\phi(z) = z$.

Ordinary smooth setting of the function p : The proof of Theorem 3 when the function p is ordinary smooth also proceeds in the similar fashion as that when p is supersmooth. In particular, we have

$$\int_{\mathbb{R}^D} \widehat{p}^2(\mathbf{t})(1 - \bar{K}_R(\mathbf{t}))^2 d\mathbf{t} \leq c \sum_{i=1}^D \int_{B_i} \prod_{j=1}^D \frac{1}{(1 + |t_j|^\beta)} d\mathbf{t}, \quad (26)$$

where $B_i := \{\mathbf{t} \in \mathbb{R}^D : |t_i| \geq R\}$. By simple algebra, we obtain that

$$\begin{aligned} \int_{B_i} \prod_{j=1}^D \frac{1}{(1 + |t_j|^\beta)} d\mathbf{t} &= \left(\int_{\mathbb{R}} \frac{1}{1 + |x|^\beta} dx \right)^{D-1} \cdot \int_{|x| \geq R} \frac{1}{1 + |x|^\beta} dx \\ &\leq \left(\int_{\mathbb{R}} \frac{1}{1 + |x|^\beta} dx \right)^{D-1} \frac{2}{\beta - 1} R^{-\beta+1}. \end{aligned}$$

Putting the above results together leads to

$$\int_{\mathbb{R}^D} \widehat{p}^2(\mathbf{t})(1 - \bar{K}_R(\mathbf{t}))^2 d\mathbf{t} \leq c_1 R^{-\beta+1}, \quad (27)$$

where c_1 is some universal constant.

Similar to the supersmooth setting, we also can bound the variance $\frac{1}{N} \int_{\mathbb{R}^D} (1 - |\widehat{p}(\mathbf{t})|^2) \bar{K}_R^2(\mathbf{t}) d\mathbf{t}$ under the ordinary smooth setting as follows:

$$\frac{1}{N} \int_{\mathbb{R}^D} (1 - |\widehat{p}(\mathbf{t})|^2) \bar{K}_R^2(\mathbf{t}) d\mathbf{t} \leq \frac{R^D}{N}. \quad (28)$$

Combining the results from equations (27) and (18), we obtain that

$$\text{MISE}(p_{N,R}^\phi) \leq c_2 \left(R^{-\beta+1} + \frac{R^D}{N} \right),$$

where c_2 is a universal constant. As a consequence, we obtain the conclusion of Theorem 3 under the ordinary smooth setting of the function p and $\phi(z) = z$.

B.2.2 When $\phi(z) = z^2$

When $\phi(z) = z^2$, which corresponds to the Féjer integral setting, we find that

$$\bar{K}_R(t) = \frac{1}{2^D} \prod_{i=1}^d \left(2 - \left| \frac{t_i}{R} \right| \right) \mathbf{1}_{\{|t_i| \leq 2R\}}.$$

Given the formulation of the function \bar{K}_R , we first bound $\frac{1}{N} \int_{\mathbb{R}^D} (1 - |\hat{p}(\mathbf{t})|^2) \bar{K}_R^2(\mathbf{t}) d\mathbf{t}$. Indeed, direct calculation shows that

$$\begin{aligned} \frac{1}{N} \int_{\mathbb{R}^D} (1 - |\hat{p}(\mathbf{t})|^2) \bar{K}_R^2(\mathbf{t}) d\mathbf{t} &\leq \frac{1}{N} \int_{\mathbb{R}^D} \bar{K}_R^2(\mathbf{t}) d\mathbf{t} = \frac{1}{N 2^D} \left(\int_{|x| \leq 2R} \left(2 - \frac{|x|}{R} \right) dx \right)^D \\ &= \frac{2^D R^D}{N}. \end{aligned} \quad (29)$$

Now, we proceed to upper bound $\int_{\mathbb{R}^D} \hat{p}^2(\mathbf{t}) (1 - \bar{K}_R(\mathbf{t}))^2 d\mathbf{t}$. We have two settings of the function p .

Supersmooth setting of the function p : Given the above formulation of the function \bar{K}_R , we have

$$\begin{aligned} \int_{\mathbb{R}^D} \hat{p}^2(\mathbf{t}) (1 - \bar{K}_R(\mathbf{t}))^2 d\mathbf{t} &= \int_{\mathbb{R}^D \setminus [-2R, 2R]^D} \hat{p}^2(\mathbf{t}) d\mathbf{t} \\ &\quad + \int_{[-2R, 2R]^D} \hat{p}^2(\mathbf{t}) \left(1 - \prod_{i=1}^D \left(1 - \frac{|t_i|}{2R} \right) \right)^2 d\mathbf{t}. \end{aligned} \quad (30)$$

By using the similar argument as when $\phi(x) = x$, when p is supersmooth function, we obtain that

$$\int_{\mathbb{R}^D \setminus [-2R, 2R]^D} \hat{p}^2(\mathbf{t}) d\mathbf{t} \leq C'_1 R^{\max\{1-\alpha, 0\}} \exp(-C'_2 R^\alpha), \quad (31)$$

where C'_1 and C'_2 are universal constants. On the other hand, we have

$$\begin{aligned} &\int_{[-2R, 2R]^D} \hat{p}^2(\mathbf{t}) \left(1 - \prod_{i=1}^D \left(1 - \frac{|t_i|}{2R} \right) \right)^2 d\mathbf{t} \\ &\leq C_1 \int_{[-2R, 2R]^D} \exp \left(-C_2 \left(\sum_{j=1}^D |t_j|^\alpha \right) \right) \left(1 - \prod_{i=1}^D \left(1 - \frac{|t_i|}{2R} \right) \right)^2 d\mathbf{t} \\ &\leq \bar{C}_1 \sum_{m=1}^D \sum_{i_1, \dots, i_m} \int_{[-2R, 2R]^D} \exp \left(-C_2 \left(\sum_{j=1}^D |t_j|^\alpha \right) \right) \frac{\prod_{l=1}^m t_{i_l}^2}{R^{2m}} d\mathbf{t}, \end{aligned} \quad (32)$$

where \bar{C}_1 is some universal constant. Here, i_1, \dots, i_m in the sum satisfy that they are pairwise different and $1 \leq i_1, \dots, i_m \leq D$. Now, simple calculations indicate that

$$\int_{[-2R, 2R]^D} \exp \left(-C_2 \left(\sum_{j=1}^D |t_j|^\alpha \right) \right) \frac{\prod_{l=1}^m t_{i_l}^2}{R^{2m}} d\mathbf{t} \leq$$

$$\frac{1}{R^{2m}} \int_{\mathbb{R}^D} \exp \left(-C_2 \left(\sum_{j=1}^D |t_j|^\alpha \right) \right) \prod_{l=1}^m t_{i_l}^2 dt \leq \frac{\bar{C}_2}{R^{2m}}, \quad (33)$$

where \bar{C}_2 is some universal constant. Combining the results from equations (32) and (33), there exists universal constant \bar{C}_3 depending on D such that

$$\int_{[-2R, 2R]^D} \hat{p}^2(\mathbf{t}) \left(1 - \prod_{i=1}^D \left(1 - \frac{|t_i|}{2R} \right) \right)^2 dt \leq \frac{\bar{C}_3}{R^2}. \quad (34)$$

Plugging the inequalities (31) and (34) to equation (30) leads to the following bound

$$\int_{\mathbb{R}^D} \hat{p}^2(\mathbf{t}) (1 - \bar{K}_R(\mathbf{t}))^2 dt \leq C'_1 R^{\max\{1-\alpha, 0\}} \exp(-C'_2 R^\alpha) + \frac{\bar{C}_3}{R^2} \leq \frac{\bar{C}_4}{R^2}. \quad (35)$$

Combining the results from equations (29) and (35), we have

$$\text{MISE}(p_{N,R}^\phi) \leq \bar{C}_5 \left(\frac{1}{R^2} + \frac{R^D}{N} \right).$$

As a consequence, we obtain the conclusion of Theorem 3 when $\phi(z) = z^2$ and the function p is supersmooth function.

Ordinary smooth setting of the function p : Using similar proof argument as that of the supersmooth setting of the function p , as $\beta > 3$, we find that

$$\begin{aligned} \int_{\mathbb{R}^D} \hat{p}^2(\mathbf{t}) (1 - \bar{K}_R(\mathbf{t}))^2 dt &\leq \frac{c}{R^{\beta-1}} + \int_{[-2R, 2R]^D} \hat{p}^2(\mathbf{t}) \left(1 - \prod_{i=1}^D \left(1 - \frac{|t_i|}{2R} \right) \right)^2 dt \\ &\leq \frac{c}{R^{\beta-1}} + \frac{c_1}{R^2} \leq \frac{c_2}{R^2}, \end{aligned} \quad (36)$$

where c, c_1, c_2 are universal constants. Combining the inequalities (29) and (36), we obtain the conclusion of Theorem 3 under the ordinary smooth setting of the function p and $\phi(z) = z^2$.

B.3 Proof of Theorem 2

Our proof strategy is to first bound the bias of $f_{N,R}(\mathbf{k})$ and then establish an upper bound for the variance of $f_{N,R}(\mathbf{k})$ for each $\mathbf{k} \in \mathbb{R}^D$.

B.3.1 Upper bound on the bias

Recall that in equation (14), we define $f_{N,R}(\mathbf{k})$ as follows:

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

where $p_{N,R}^\phi(\mathbf{k})$ is generalized Fourier density estimator in equation (12) while $a_{N,R}(\mathbf{k})$ is defined as follows:

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

Simple algebra leads to

$$f_{N,R}(\mathbf{k}) - f(\mathbf{k}) = \frac{a_{N,R}(\mathbf{k}) - f(\mathbf{k})p_{N,R}^\phi(\mathbf{k})}{p(\mathbf{k})} + \frac{(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))(p(\mathbf{k}) - p_{N,R}^\phi(\mathbf{k}))}{p(\mathbf{k})}. \quad (37)$$

Therefore, via an application of Cauchy-Schwarz inequality we obtain that

$$\begin{aligned} & (\mathbb{E}[f_{N,R}(\mathbf{k})] - f(\mathbf{k}))^2 \\ & \leq 2 \frac{\left(\mathbb{E}\left[a_{N,R}(\mathbf{k}) - f(\mathbf{k})p_{N,R}^\phi(\mathbf{k})\right]\right)^2}{p^2(\mathbf{k})} + 2 \frac{\left(\mathbb{E}\left[(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))(p(\mathbf{k}) - p_{N,R}^\phi(\mathbf{k}))\right]\right)^2}{p^2(\mathbf{k})} \\ & \leq 2 \frac{\left(\mathbb{E}\left[a_{N,R}(\mathbf{k}) - f(\mathbf{k})p_{N,R}^\phi(\mathbf{k})\right]\right)^2}{p^2(\mathbf{k})} + 2 \frac{\mathbb{E}\left[(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2\right] \mathbb{E}\left[(p(\mathbf{k}) - p_{N,R}^\phi(\mathbf{k}))^2\right]}{p^2(\mathbf{k})}, \end{aligned} \quad (38)$$

where the second inequality is due to the standard inequality $\mathbb{E}^2(XY) \leq \mathbb{E}(X^2)\mathbb{E}(Y^2)$ for all the random variables X, Y .

According to the assumptions of Theorem 2 and the result of Theorem 1, we have

$$\mathbb{E}\left[(p(\mathbf{k}) - p_{N,R}^\phi(\mathbf{k}))^2\right] \leq \frac{C_1}{R^{2(m+1)}} + \frac{C_2 R^D}{N}, \quad (39)$$

where C_1 and C_2 are some universal constants in Theorem 1.

Now, we proceed to bound $|\mathbb{E}[a_{N,R}(\mathbf{k}) - f(\mathbf{k})p_{N,R}^\phi(\mathbf{k})]|$. Direct calculation demonstrates that

$$\begin{aligned} \mathbb{E}[a_{N,R}(\mathbf{k})] &= \frac{R^D}{A^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(R(k_j - y_j))}{R(k_j - y_j)}\right) p(\mathbf{y}) f(\mathbf{y}) d\mathbf{y} \\ &= \frac{1}{A^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(y_j)}{y_j}\right) p\left(\mathbf{k} - \frac{\mathbf{y}}{R}\right) f\left(\mathbf{k} - \frac{\mathbf{y}}{R}\right) d\mathbf{y}. \end{aligned} \quad (40)$$

An application of Taylor expansion up to the m -th order indicates that

$$\begin{aligned} p\left(\mathbf{k} - \frac{\mathbf{y}}{R}\right) &= \sum_{0 \leq |\alpha| \leq m} \frac{1}{R^{|\alpha| \alpha!}} \prod_{j=1}^D (-y_j)^{\alpha_j} \frac{\partial^{|\alpha|} p}{\partial \mathbf{k}^\alpha}(\mathbf{k}) + \bar{R}_1(\mathbf{k}, \mathbf{y}), \\ f\left(\mathbf{k} - \frac{\mathbf{y}}{R}\right) &= \sum_{0 \leq |\alpha| \leq m} \frac{1}{R^{|\alpha| \alpha!}} \prod_{j=1}^D (-y_j)^{\alpha_j} \frac{\partial^{|\alpha|} f}{\partial \mathbf{k}^\alpha}(\mathbf{k}) + \bar{R}_2(\mathbf{k}, \mathbf{y}), \end{aligned} \quad (41)$$

where $\alpha = (\alpha_1, \dots, \alpha_d)$, $|\alpha| = \sum_{j=1}^d \alpha_j$, and $\bar{R}_1(\mathbf{k}, \mathbf{y})$, $\bar{R}_2(\mathbf{k}, \mathbf{y})$ are Taylor remainders admitting the following forms:

$$\begin{aligned} \bar{R}_1(\mathbf{k}, \mathbf{y}) &= \sum_{|\beta|=m+1} \frac{m+1}{R^{m+1} \beta!} \prod_{j=1}^D (-y_j)^{\beta_j} \int_0^1 (1-t)^m \frac{\partial^{m+1} p}{\partial \mathbf{k}^\beta} \left(\mathbf{k} - \frac{t\mathbf{y}}{R}\right) dt, \\ \bar{R}_2(\mathbf{k}, \mathbf{y}) &= \sum_{|\beta|=m+1} \frac{m+1}{R^{m+1} \beta!} \prod_{j=1}^D (-y_j)^{\beta_j} \int_0^1 (1-t)^m \frac{\partial^{m+1} f}{\partial \mathbf{k}^\beta} \left(\mathbf{k} - \frac{t\mathbf{y}}{R}\right) dt. \end{aligned} \quad (42)$$

Combining equations (41) and (42), we obtain that

$$\begin{aligned}
p\left(\mathbf{k} - \frac{\mathbf{y}}{R}\right) f\left(\mathbf{k} - \frac{\mathbf{y}}{R}\right) &= \sum_{0 \leq |\alpha|, |\beta| \leq m} \frac{1}{R^{|\alpha|+|\beta|} \alpha! \beta!} \prod_{j=1}^D (-y_j)^{\alpha_j + \beta_j} \frac{\partial^{|\alpha|} p}{\partial \mathbf{k}^\alpha}(\mathbf{k}) \frac{\partial^{|\beta|} f}{\partial \mathbf{k}^\beta}(\mathbf{k}) \\
&+ \left(\sum_{0 \leq |\alpha| \leq m} \frac{1}{R^{|\alpha|} \alpha!} \prod_{j=1}^D (-y_j)^{\alpha_j} \frac{\partial^{|\alpha|} p}{\partial \mathbf{k}^\alpha}(\mathbf{k}) \right) \bar{R}_2(\mathbf{k}, \mathbf{y}) \\
&+ \left(\sum_{0 \leq |\alpha| \leq m} \frac{1}{R^{|\alpha|} \alpha!} \prod_{j=1}^D (-y_j)^{\alpha_j} \frac{\partial^{|\alpha|} f}{\partial \mathbf{k}^\alpha}(\mathbf{k}) \right) \bar{R}_1(\mathbf{k}, \mathbf{y}) + \bar{R}_1(\mathbf{k}, \mathbf{y}) \bar{R}_2(\mathbf{k}, \mathbf{y}).
\end{aligned}$$

As we have $\int_{\mathbb{R}} \phi\left(\frac{\sin(z)}{z}\right) z^j dz = 0$ for all $1 \leq j \leq m$, plugging the equation in the above display to equation (40) leads to

$$\mathbb{E}[a_{n,R}(\mathbf{k})] = f(\mathbf{k}) \mathbb{E}\left[p_{N,R}^\phi(\mathbf{k})\right] + B_1 + B_2 + B_3 + B_4,$$

where B_1, B_2, B_3, B_4 are defined as follows:

$$\begin{aligned}
B_1 &= \frac{1}{A^D} \sum_{m+1 \leq |\alpha|+|\beta| \leq 2m} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(y_j)}{y_j}\right) \frac{1}{R^{|\alpha|+|\beta|} \alpha! \beta!} \prod_{j=1}^D (-y_j)^{\alpha_j + \beta_j} \frac{\partial^{|\alpha|} p}{\partial \mathbf{k}^\alpha}(\mathbf{k}) \frac{\partial^{|\beta|} f}{\partial \mathbf{k}^\beta}(\mathbf{k}) d\mathbf{y}, \\
B_2 &= \frac{1}{A^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(y_j)}{y_j}\right) \left(\sum_{0 \leq |\alpha| \leq m} \frac{1}{R^{|\alpha|} \alpha!} \prod_{j=1}^D (-y_j)^{\alpha_j} \frac{\partial^{|\alpha|} p_0}{\partial \mathbf{k}^\alpha}(\mathbf{k}) \right) \bar{R}_2(\mathbf{k}, \mathbf{y}) d\mathbf{y}, \\
B_3 &= \frac{1}{A^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(y_j)}{y_j}\right) \left(\sum_{0 \leq |\alpha| \leq m} \frac{1}{R^{|\alpha|} \alpha!} \prod_{j=1}^D (-y_j)^{\alpha_j} \frac{\partial^{|\alpha|} f}{\partial \mathbf{k}^\alpha}(\mathbf{k}) \right) \bar{R}_1(\mathbf{k}, \mathbf{y}) d\mathbf{y}, \\
B_4 &= \frac{1}{A^D} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi\left(\frac{\sin(y_j)}{y_j}\right) \bar{R}_1(\mathbf{k}, \mathbf{y}) \bar{R}_2(\mathbf{k}, \mathbf{y}) d\mathbf{y}.
\end{aligned}$$

Since we have $\int_{\mathbb{R}} \left| \phi\left(\frac{\sin(z)}{z}\right) \right| |z|^j dz < \infty$ for any $m+1 \leq j \leq 2m+2$ and $p_0, f \in \mathcal{C}^{m+1}(\mathbb{R}^d)$, we find that as long as $R \geq \bar{c}$ for some given constant \bar{c}

$$\begin{aligned}
|B_1| &\leq \frac{1}{A^D} \sum_{m+1 \leq |\alpha|+|\beta| \leq 2m} \frac{1}{R^{|\alpha|+|\beta|} \alpha! \beta!} \int_{\mathbb{R}^D} \prod_{j=1}^D \left| \phi\left(\frac{\sin(y_j)}{y_j}\right) \right| \prod_{j=1}^D |y_j|^{\alpha_j + \beta_j} \left\| \frac{\partial^{|\alpha|} p}{\partial \mathbf{k}^\alpha} \right\|_\infty \left\| \frac{\partial^{|\beta|} f}{\partial \mathbf{k}^\beta} \right\|_\infty \\
&\leq \frac{c_1}{R^{m+1}},
\end{aligned}$$

where c_1 is some universal constant depending on A, D , and \bar{c} . Furthermore, we find that

$$\begin{aligned}
|B_2| &\leq \frac{1}{A^D} \sum_{0 \leq |\alpha| \leq m, |\beta|=m+1} \frac{m+1}{R^{|\alpha|+m+1} \alpha! \beta!} \int_{\mathbb{R}^D} \prod_{j=1}^D \left| \phi\left(\frac{\sin(y_j)}{y_j}\right) \right| \prod_{j=1}^D |y_j|^{\alpha_j + \beta_j} \\
&\quad \times \int_0^1 (1-t)^m \left\| \frac{\partial^{m+1} f}{\partial \mathbf{k}^\beta} \right\|_\infty d\mathbf{y} dt \leq \frac{c_2}{R^{m+1}},
\end{aligned}$$

where c_2 is some universal constant depending on A , d , and \bar{c} . Similarly, we also can demonstrate that $B_3 \leq c_3/R^{m+1}$ and $B_4 \leq c_4/R^{2(m+1)}$ for some universal constants c_3 and c_4 . Putting the above results together, we arrive at the following bound:

$$\left| \mathbb{E} \left[a_{n,R}(\mathbf{k}) - f(\mathbf{k}) p_{N,R}^\phi(\mathbf{k}) \right] \right| \leq \frac{c'}{R^{m+1}}. \quad (43)$$

Plugging the results from equations (39) and (43) to equation (38), we obtain that

$$(\mathbb{E}[f_{N,R}(\mathbf{k})] - f(\mathbf{k}))^2 \leq \frac{2(c')^2}{p^2(\mathbf{k})R^{2(m+1)}} + \frac{2\mathbb{E}[(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2]}{p^2(\mathbf{k})} \left(\frac{C_1}{R^{2(m+1)}} + \frac{C_2 R^D}{N} \right). \quad (44)$$

B.3.2 Upper bound on the variance

Now, we study the variance of $f_{N,R}(\mathbf{k})$. By taking variance both sides of the equation (37), we obtain that

$$\begin{aligned} \text{var}(f_{N,R}(\mathbf{k})) &= \text{var} \left(\frac{a_{N,R}(\mathbf{k}) - f(\mathbf{k}) p_{N,R}^\phi(\mathbf{k})}{p(\mathbf{k})} + \frac{(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))(p(\mathbf{k}) - p_{N,R}^\phi(\mathbf{k}))}{p(\mathbf{k})} \right) \\ &\leq \frac{2}{p^2(\mathbf{k})} \left(\underbrace{\mathbb{E} \left[\left(a_{N,R}(\mathbf{k}) - f(\mathbf{k}) p_{N,R}^\phi(\mathbf{k}) \right)^2 \right]}_{T_1} + \underbrace{\mathbb{E} \left[\left(f_{N,R}(\mathbf{k}) - f(\mathbf{k}) \right)^2 \left(p(\mathbf{k}) - p_{N,R}^\phi(\mathbf{k}) \right)^2 \right]}_{T_2} \right). \end{aligned} \quad (45)$$

Upper bound of T_2 : To upper bound T_2 , we utilize the following lemma.

Lemma 1. *Assume that the function ϕ and p_0 satisfy the assumptions of Theorem 1. Furthermore, $\phi(z) \leq C$ as long as $|z| \leq 1$ for some universal constant C . Then, for almost all $\mathbf{k} \in \mathbb{R}^D$, there exist universal constants C' such that*

$$\mathbb{P} \left(\left| p_{N,R}^\phi(\mathbf{k}) - p(\mathbf{k}) \right| \geq C' \left(\frac{1}{R^{m+1}} + \sqrt{\frac{R^D \log(2/\delta)}{N}} \right) \right) \leq \delta.$$

Proof of Lemma 1 is given in Appendix B.4. Now given the result of Lemma 1, we denote B as the event such that

$$\left| p_{N,R}^\phi(\mathbf{k}) - p(\mathbf{k}) \right| \leq C' \left(\frac{1}{R^{m+1}} + \sqrt{\frac{R^D \log(2/\delta)}{N}} \right)$$

where C' is a universal constant in Lemma 1. Then, we obtain $\mathbb{P}(B) \geq 1 - \delta$. Hence, we have the following bound with T_2 :

$$\begin{aligned} T_2 &= \mathbb{E} \left[(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2 (p(\mathbf{k}) - p_{N,R}^\phi(\mathbf{k}))^2 | B \right] \mathbb{P}(B) \\ &\quad + \mathbb{E} \left[(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2 (p(\mathbf{k}) - p_{N,R}^\phi(\mathbf{k}))^2 | B^c \right] \mathbb{P}(B^c) \\ &\leq 2c' \mathbb{E} \left[(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2 \right] \left(\frac{1}{R^{2(m+1)}} + \frac{R^D \log(2/\delta)}{N} + \delta \left(p^2(\mathbf{k}) + \frac{C^D R^{2D}}{A^D} \right) \right), \end{aligned}$$

where c' is some universal constant and the final inequality is based on the inequalities: $\mathbb{P}(B^c) \leq \delta$ and $(p(\mathbf{k}) - p_{N,R}^\phi(\mathbf{k}))^2 \leq 2(p^2(\mathbf{k}) + (p_{N,R}^\phi)^2(\mathbf{k})) \leq 2\left(p^2(\mathbf{k}) + \frac{C^D R^{2D}}{A^D}\right)$ where C is a universal constant such that $\phi(z) \leq C$ when $|z| \leq 1$. By choosing δ such that $\delta = \frac{R^D}{N(p^2(\mathbf{k}) + C^D R^{2D}/A^D)}$, we obtain that

$$T_2 \leq c'' \mathbb{E} [(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2] \left(\frac{1}{R^{2(m+1)}} + \frac{R^D \log(NR)}{N} \right), \quad (46)$$

for some universal constant c'' when R is sufficiently large.

Upper bound of T_1 : As $\mathbf{v}_i = f(\mathbf{k}_i) + \epsilon_i$ for all $i \in [N]$, direct calculation shows that

$$T_1 = \mathbb{E} \left[\left(\frac{R^D}{NA^D} \sum_{i=1}^N (f(\mathbf{k}_i) - f(\mathbf{k})) \prod_{j=1}^D \phi \left(\frac{\sin(R(k_j - k_{ij}))}{R(k_j - k_{ij})} \right) + \frac{R^D}{NA^D} \sum_{i=1}^N \epsilon_i \prod_{j=1}^D \phi \left(\frac{\sin(R(k_j - k_{ij}))}{R(k_j - k_{ij})} \right) \right)^2 \right].$$

An application of Cauchy-Schwarz inequality leads to

$$T_1 \leq 2\mathbb{E} \left[\left(\frac{R^D}{NA^D} \sum_{i=1}^N (f(\mathbf{k}_i) - f(\mathbf{k})) \prod_{j=1}^D \phi \left(\frac{\sin(R(k_j - k_{ij}))}{R(k_j - k_{ij})} \right) \right)^2 \right] + 2\mathbb{E} \left[\left(\frac{1}{N\pi^D} \sum_{i=1}^N \epsilon_i \prod_{j=1}^D \phi \left(\frac{\sin(R(k_j - k_{ij}))}{R(k_j - k_{ij})} \right) \right)^2 \right] = 2(S_1 + S_2).$$

Since we have $\mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N Z_i \right)^2 \right] \leq \frac{1}{N} \mathbb{E} [Z_1^2] + \mathbb{E}^2 [Z_1]$ for any i.i.d. samples Z_1, \dots, Z_N , we obtain that

$$S_1 \leq \frac{R^{2D}}{NA^{2D}} \mathbb{E} \left[(f(X) - f(\mathbf{k}))^2 \prod_{j=1}^D \phi^2 \left(\frac{\sin(R(k_j - X_{.j}))}{R(k_j - X_{.j})} \right) \right] + \frac{R^{2D}}{A^{2D}} \mathbb{E}^2 \left[(f(X) - f(\mathbf{k})) \prod_{j=1}^D \phi \left(\frac{\sin(R(k_j - X_{.j}))}{R(k_j - X_{.j})} \right) \right],$$

where the outer expectation is taken with respect to $X = (X_{.1}, \dots, X_{.d}) \sim p$. From the result in equation (43), we have

$$\frac{R^{2D}}{A^{2D}} \mathbb{E}^2 \left[(f(X) - f(\mathbf{k})) \prod_{j=1}^D \phi \left(\frac{\sin(R(k_j - X_{.j}))}{R(k_j - X_{.j})} \right) \right] = \mathbb{E}^2 [a_{N,R}(\mathbf{k}) - f(\mathbf{k}) p_{N,R}^\phi(\mathbf{k})] \leq \frac{c'}{R^{2(m+1)}},$$

where c' is some universal constant. In addition, an application of Cauchy-Schwarz inequality leads to

$$\frac{R^{2D}}{NA^{2D}} \mathbb{E} \left[(f(X) - f(\mathbf{k}))^2 \prod_{j=1}^D \phi^2 \left(\frac{\sin(R(k_j - X_{.j}))}{R(k_j - X_{.j})} \right) \right]$$

$$\begin{aligned}
&\leq \frac{2R^{2D}}{NA^{2D}} \mathbb{E} \left[(f^2(X) + f^2(\mathbf{k})) \prod_{j=1}^D \phi^2 \left(\frac{\sin(R(k_j - X_{.j}))}{R(k_j - X_{.j})} \right) \right] \\
&= \frac{2R^D}{NA^{2D}} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi^2 \left(\frac{\sin(y_j)}{y_j} \right) \left(f^2 \left(\mathbf{k} - \frac{\mathbf{y}}{R} \right) p \left(\mathbf{k} - \frac{\mathbf{y}}{R} \right) + f^2(\mathbf{k}) \right) d\mathbf{y} \\
&\leq \frac{2R^D (\|f^2 \times p\|_\infty + f^2(\mathbf{k}))}{NA^{2D}} \int_{\mathbb{R}^D} \prod_{j=1}^D \phi^2 \left(\frac{\sin(y_j)}{y_j} \right) dy.
\end{aligned}$$

Since we have $\int_{\mathbb{R}} \phi^2(\sin(z)/z) dz < \infty$, it indicates that we can find a universal constant c'' such that

$$\frac{R^{2D}}{NA^{2D}} \mathbb{E} \left[(f(X) - f(\mathbf{k}))^2 \prod_{j=1}^D \phi^2 \left(\frac{\sin(R(k_j - X_{.j}))}{R(k_j - X_{.j})} \right) \right] \leq \frac{c'' R^D (\|f^2 \times p\|_\infty + f^2(\mathbf{k}))}{NA^{2D}}.$$

Putting the above results together, we obtain that

$$S_1 \leq \frac{c'}{R^{2(m+1)}} + \frac{c'' R^D (\|f^2 \times p\|_\infty + f^2(\mathbf{k}))}{NA^{2D}}. \quad (47)$$

Similarly, since $\mathbb{E}(\epsilon_i) = 0$ and $\text{var}(\epsilon_i) = \sigma^2$ for all $i \in [N]$, we have

$$S_2 = \frac{\sigma^2 R^{2D}}{NA^{2D}} \mathbb{E} \left[\prod_{j=1}^D \phi^2 \left(\frac{\sin(R(k_j - X_{.j}))}{R(k_j - X_{.j})} \right) \right] \leq \frac{c''' \sigma^2 R^D \|p\|_\infty R^D}{NA^{2D}}, \quad (48)$$

where c''' is some universal constant. Combining the results from equation (47) and equation (48), we find that

$$T_1 \leq C \left(\frac{(\|f^2 \times p\|_\infty + f^2(\mathbf{k}) + \sigma^2 \|p\|_\infty) R^D}{N} + \frac{1}{R^{2(m+1)}} \right), \quad (49)$$

where C is some universal constant. Plugging the bounds of T_1 and T_2 from equations (46) and (49) into equation (45), when $R \geq C'$ where C' is some universal constant, we have

$$\begin{aligned}
\text{var}(f_{N,R}(\mathbf{k})) &\leq \frac{C'_1}{p^2(\mathbf{k})} \mathbb{E} [(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2] \left(\frac{1}{R^{2(m+1)}} + \frac{R^D \log(NR)}{N} \right) \\
&\quad + \frac{C'_2}{p^2(\mathbf{k})} \left(\frac{(f(\mathbf{k}) + C'_3) R^D}{N} + \frac{1}{R^{2(m+1)}} \right), \quad (50)
\end{aligned}$$

where C'_1, C'_2, C'_3 are some universal constants. Combining the results with bias and variance in equations (44) and (50), we obtain the following bound:

$$\begin{aligned}
\mathbb{E} [(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2] &\leq \frac{2(c')^2}{p^2(\mathbf{k}) R^{2(m+1)}} + \frac{2\mathbb{E} [(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2]}{p^2(\mathbf{k})} \left(\frac{C_1}{R^{2(m+1)}} + \frac{C_2 R^D}{N} \right) \\
&\quad + \frac{C'_1}{p^2(\mathbf{k})} \mathbb{E} [(f_{N,R}(\mathbf{k}) - f(\mathbf{k}))^2] \left(\frac{1}{R^{2(m+1)}} + \frac{R^D \log(NR)}{N} \right) \\
&\quad + \frac{C'_2}{p^2(\mathbf{k})} \left(\frac{(f(\mathbf{k}) + C'_3) R^D}{N} + \frac{1}{R^{2(m+1)}} \right).
\end{aligned}$$

As a consequence, we obtain the conclusion of the theorem.

B.4 Proof of Lemma 1

Invoking triangle inequality, we obtain that

$$\left| p_{N,R}^\phi(\mathbf{k}) - p(\mathbf{k}) \right| \leq \left| p_{N,R}^\phi(\mathbf{k}) - \mathbb{E} \left[p_{N,R}^\phi(\mathbf{k}) \right] \right| + \left| \mathbb{E} \left[p_{N,R}^\phi(\mathbf{k}) \right] - p(\mathbf{k}) \right|. \quad (51)$$

If we denote $\mathbf{v}_i = \frac{R^D}{A^D} \prod_{j=1}^D \phi \left(\frac{\sin(R(k_j - k_{ij}))}{R(k_j - k_{ij})} \right)$ for all $i \in [N]$, then as $\sin(R(k_j - k_{ij})) / (R(k_j - k_{ij})) \leq 1$ for all $j \in [D]$ we have $|\mathbf{v}_i| \leq C^D R^D / A^D$ for all $i \in [N]$ where C is the constant such that $\phi(z) \leq C$ when $|z| \leq 1$. Furthermore, from the proof of Theorem 1 we have $\text{var}(\mathbf{v}_i) \leq C' R^D$ where $C' > 0$ is some universal constant. Given these bounds of \mathbf{v}_i and $\text{var}(\mathbf{v}_i)$, for any $t \in (0, C'']$ Bernstein's inequality shows that

$$\mathbb{P} \left(\left| \frac{1}{N} \sum_{i=1}^N \mathbf{v}_i - \mathbb{E}[\mathbf{v}_1] \right| \geq t \right) \leq 2 \exp \left(- \frac{Nt^2}{2C'R^D + 2C^D R^D t / (3A^D)} \right).$$

By choosing $t = \bar{C} \sqrt{R^D \log(2/\delta) / N}$, where \bar{C} is some universal constant, we find that

$$\mathbb{P} \left(\left| p_{N,R}^\phi(\mathbf{k}) - \mathbb{E} \left[p_{N,R}^\phi(\mathbf{k}) \right] \right| \geq t \right) = \mathbb{P} \left(\left| \frac{1}{N} \sum_{i=1}^N \mathbf{v}_i - \mathbb{E}[\mathbf{v}_1] \right| \geq t \right) \leq \delta. \quad (52)$$

From the result of Theorem 1, there exists universal constant c such that

$$\left| \mathbb{E} \left[p_{N,R}^\phi(\mathbf{k}) \right] - p(\mathbf{k}) \right| \leq c / R^{m+1}. \quad (53)$$

Plugging the bounds (52) and (53) into the triangle inequality (51), we obtain the conclusion of the lemma.

C Experiment Details

C.1 Language Modeling on WikiText-103

In our experiments on WikiText-103 in Section 4.1, we let R be a learnable scalar initialized to 2 and choose $\phi(x) = x^4$. The same setting is used for all attention units in the model; each unit has a different R . We observe that by setting R to be a learnable vector $[R_1, \dots, R_D]^\top$, the FourierFormer gains advantage in accuracy but with the cost of the increase in the number of parameters. When R is a vector $[R_1, \dots, R_D]^\top$, the equation of the Fourier Attention is given by

$$\hat{\mathbf{h}}_i := f_{N,R}(\mathbf{q}_i) = \frac{\sum_{i=1}^N \mathbf{v}_i \prod_{j=1}^D \phi \left(\frac{\sin(R_j(q_{ij} - k_{ij}))}{R_j(q_{ij} - k_{ij})} \right)}{\sum_{i=1}^N \prod_{j=1}^D \phi \left(\frac{\sin(R_j(q_{ij} - k_{ij}))}{R_j(q_{ij} - k_{ij})} \right)} \quad \forall i \in [N]. \quad (54)$$

We provide an ablation study for the effect of R and ϕ in Section D below.

C.2 Image Classification on ImageNet

Similar to setting for language modeling, in our experiments on ImageNet image classification, we set R to be a learnable scalar initialized to 1 and choose $\phi(x) = x^4$. Different attention units have different R .

Table 4. Ablation study on how the choice of $\phi(x) = x^k$ influences the performance of FourierFormer. Odd values of k cause training to diverge. For even values of k , greater k yields better perplexity (PPL), but the improvement is small for $k > 4$.

Method	Valid PPL	Test PPL
<i>Baseline dot-product (small)</i>	33.15	34.29
FourierFormer, $\phi(x) = x^2$ (small)	32.09	33.10
FourierFormer, $\phi(x) = x^4$ (small)	31.86	32.85
FourierFormer, $\phi(x) = x^6$ (small)	31.84	32.81
FourierFormer, $\phi(x) = x$ (small)	not converge	not converge
FourierFormer, $\phi(x) = x^3$ (small)	not converge	not converge

Table 5. Ablation study on how the initialization of R influences the performance of FourierFormer. When R is initialized to a too small or too big value, the PPL of the trained FourierFormer is reduced. $R_{\text{init}} = 1, 2, 3$ yield the best results. Fourierformer with learnable vectors R yields better results than Fourierformer of the same setting using learnable scalars R with the cost of increasing the number of parameters in the model.

Method	Valid PPL	Test PPL
<i>Baseline dot-product (small)</i>	33.15	34.29
FourierFormer, $R_{\text{init}} = 0.1$ (small)	32.04	33.01
FourierFormer, $R_{\text{init}} = 1.0$ (small)	31.89	32.87
FourierFormer, $R_{\text{init}} = 2.0$ (small)	31.86	32.85
FourierFormer, $R_{\text{init}} = 3.0$ (small)	31.90	32.88
FourierFormer, $R_{\text{init}} = 4.0$ (small)	32.58	33.65
FourierFormer, $R_{\text{init}} = 2.0$ (small, R is a vector)	31.82	32.80

D Additional Experimental Results

D.1 Effect of ϕ

Using the WikiText-103 language modeling as a case study, we analyze the effect of $\phi(x)$ on the performance of FourierFormer. In particular, we set $\phi(x) = x^k$ and compare the performance of FourierFormer for $k = 1, 2, 3, 4$ and 6 . We keep other settings the same as in our experiments in Section 4.1. We summarize our results in Table 4. We observe that for odd values of k such as $k = 1, 3$, the training diverges, confirming that negative density estimator cause instability in training FourierFormer (see Remark 3.1). For even values of k such as $k = 2, 4, 6$, we observe that the greater value of k results in better valid and test PPL. However, the gap between $k = 4$ and $k = 6$ is smaller compared to the gap between $k = 2$ and $k = 4$, suggesting that using $k > 4$ does not add much advantage in terms of accuracy.

D.2 Effect of the Initialization of R

In this section, we study the effect of the initialization value of R on the performance of FourierFormer when trained for the WikiText-103 language modeling and summarize our results in Table 5. Here we choose R to be learnable scalars as in experiments described in our main text. Other settings are also the same as in our experiments in Section 4.1. We observe that when R is initialized too small (e.g. $R_{\text{init}} = 0.1$) or too big (e.g. $R_{\text{init}} = 4$), the

PPL of the trained FourierFormer decreases. $R_{\text{init}} = 1, 2, 3$ yield best results.

We also study the performance of the FourierFormer when R is chosen to be a learnable vector, $R = [R_1, \dots, R_D]^\top$. We report our result in the last row of Table 5. FourierFormer with R be learnable vectors achieves better PPLs than FourierFormer with R be learnable scalars of the same setting. As we mentioned in Section C, this advantage comes with an increase in the number of parameters in the model.

Finally, from our experiments, we observe that making R a learnable parameter yields better PPLs than making R a constant and selecting its value via a careful search.

References

- [1] R. Al-Rfou, D. Choe, N. Constant, M. Guo, and L. Jones. Character-level language modeling with deeper self-attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3159–3166, 2019. (Cited on pages 1 and 10.)
- [2] C. Anil, J. Lucas, and R. Grosse. Sorting out lipschitz function approximation. In *International Conference on Machine Learning*, pages 291–301. PMLR, 2019. (Cited on page 5.)
- [3] A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lučić, and C. Schmid. Vivit: A video vision transformer. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6816–6826, 2021. (Cited on page 1.)
- [4] A. Baevski and M. Auli. Adaptive input representations for neural language modeling. In *International Conference on Learning Representations*, 2019. (Cited on page 1.)
- [5] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014. (Cited on page 1.)
- [6] S. Bochner. *Lectures on Fourier Integrals*. Princeton University Press, 1959. (Cited on page 5.)
- [7] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. (Cited on page 1.)
- [8] S. Cao. Choose a transformer: Fourier or galerkin. *Advances in Neural Information Processing Systems*, 34, 2021. (Cited on page 11.)
- [9] J. Chacón and T. Duong. *Multivariate Kernel Smoothing and its Applications*. CRC Press, 2018. (Cited on page 5.)
- [10] R. Child, S. Gray, A. Radford, and I. Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019. (Cited on page 1.)
- [11] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar, Oct. 2014. Association for Computational Linguistics. (Cited on page 1.)

- [12] K. M. Choromanski, V. Likhoshesterov, D. Dohan, X. Song, A. Gane, T. Sarlos, P. Hawkins, J. Q. Davis, A. Mohiuddin, L. Kaiser, D. B. Belanger, L. J. Colwell, and A. Weller. Rethinking attention with performers. In *International Conference on Learning Representations*, 2021. (Cited on page 11.)
- [13] M. Cisse, P. Bojanowski, E. Grave, Y. Dauphin, and N. Usunier. Parseval networks: Improving robustness to adversarial examples. In *International Conference on Machine Learning*, pages 854–863. PMLR, 2017. (Cited on page 5.)
- [14] K. Clark, U. Khandelwal, O. Levy, and C. D. Manning. What does BERT look at? an analysis of BERT’s attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy, Aug. 2019. Association for Computational Linguistics. (Cited on page 2.)
- [15] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019. (Cited on page 1.)
- [16] F. Dalvi, H. Sajjad, N. Durrani, and Y. Belinkov. Analyzing redundancy in pretrained transformer models. *arXiv preprint arXiv:2004.04010*, 2020. (Cited on page 12.)
- [17] K. B. Davis. Mean square error properties of density estimates. *The Annals of Statistics*, 3(4):1025–1030, 1975. (Cited on pages 9 and 13.)
- [18] M. Dehghani, S. Gouws, O. Vinyals, J. Uszkoreit, and L. Kaiser. Universal transformers. In *International Conference on Learning Representations*, 2019. (Cited on page 1.)
- [19] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. (Cited on pages 9 and 10.)
- [20] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. (Cited on page 1.)
- [21] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. (Cited on page 1.)
- [22] N. Durrani, H. Sajjad, F. Dalvi, and Y. Belinkov. Analyzing individual neurons in pretrained language models. *arXiv preprint arXiv:2010.02695*, 2020. (Cited on page 12.)
- [23] K. Ethayarajh. How contextual are contextualized word representations? comparing the geometry of bert, elmo, and gpt-2 embeddings. *arXiv preprint arXiv:1909.00512*, 2019. (Cited on page 12.)
- [24] H. Fan, B. Xiong, K. Mangalam, Y. Li, Z. Yan, J. Malik, and C. Feichtenhofer. Multi-scale vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6824–6835, 2021. (Cited on page 1.)

- [25] J. Fan. On the optimal rates of convergence for nonparametric deconvolution problems. *Annals of Statistics*, 19(3):1257–1272, 1991. (Cited on page 13.)
- [26] P. Gabbur, M. Bilkhu, and J. Movellan. Probabilistic attention for interactive segmentation. *Advances in Neural Information Processing Systems*, 34, 2021. (Cited on page 11.)
- [27] M.-H. Guo, J.-X. Cai, Z.-N. Liu, T.-J. Mu, R. R. Martin, and S.-M. Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7(2):187–199, 2021. (Cited on page 1.)
- [28] J. Hewitt and P. Liang. Designing and interpreting probes with control tasks. 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)*, pages 2733–2743, Hong Kong, China, Nov. 2019. Association for Computational Linguistics. (Cited on page 2.)
- [29] N. Ho and S. Walker. Multivariate smoothing via the Fourier integral theorem and Fourier kernel. *Arxiv preprint Arxiv:2012.14482*, 2021. (Cited on pages 5, 6, 13, and 14.)
- [30] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Židek, A. Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021. (Cited on page 1.)
- [31] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, pages 5156–5165. PMLR, 2020. (Cited on page 11.)
- [32] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan, and M. Shah. Transformers in vision: A survey. *ACM Computing Surveys (CSUR)*, 2021. (Cited on page 1.)
- [33] H. Kim, G. Papamakarios, and A. Mnih. The lipschitz constant of self-attention. In *International Conference on Machine Learning*, pages 5562–5571. PMLR, 2021. (Cited on page 5.)
- [34] Y. Kim, C. Denton, L. Hoang, and A. M. Rush. Structured attention networks. *arXiv preprint arXiv:1702.00887*, 2017. (Cited on page 1.)
- [35] D. Kreuzer, D. Beaini, W. Hamilton, V. Létourneau, and P. Tossou. Rethinking graph transformers with spectral attention. *Advances in Neural Information Processing Systems*, 34, 2021. (Cited on page 11.)
- [36] T. Lin, Y. Wang, X. Liu, and X. Qiu. A survey of transformers. *arXiv preprint arXiv:2106.04554*, 2021. (Cited on page 1.)
- [37] Z. Lin, M. Feng, C. N. dos Santos, M. Yu, B. Xiang, B. Zhou, and Y. Bengio. A structured self-attentive sentence embedding. *CoRR*, abs/1703.03130, 2017. (Cited on page 1.)
- [38] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. (Cited on page 1.)
- [39] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. (Cited on page 1.)

- [40] Z. Liu, J. Ning, Y. Cao, Y. Wei, Z. Zhang, S. Lin, and H. Hu. Video swin transformer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. (Cited on page 1.)
- [41] Y. Lu, Z. Li, D. He, Z. Sun, B. Dong, T. Qin, L. Wang, and T.-Y. Liu. Understanding and improving transformer from a multi-particle dynamic system point of view. *arXiv preprint arXiv:1906.02762*, 2019. (Cited on page 11.)
- [42] S. Merity, C. Xiong, J. Bradbury, and R. Socher. Pointer sentinel mixture models. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. (Cited on pages 9 and 10.)
- [43] P. Michel, O. Levy, and G. Neubig. Are sixteen heads really better than one? In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. (Cited on page 12.)
- [44] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*, 2018. (Cited on page 5.)
- [45] J. Mu and P. Viswanath. All-but-the-top: Simple and effective postprocessing for word representations. In *International Conference on Learning Representations*, 2018. (Cited on page 12.)
- [46] E. Nadaraya. On estimating regression. *Theory of Probability and its Applications*, 9:141–142, 1964. (Cited on page 4.)
- [47] T. Nguyen, T. Nguyen, D. Le, K. Nguyen, A. Tran, R. Baraniuk, N. Ho, and S. Osher. Improving transformers with probabilistic attention keys. In *ICML, 2022*. (Not cited.)
- [48] A. Parikh, O. Täckström, D. Das, and J. Uszkoreit. A decomposable attention model for natural language inference. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2249–2255, Austin, Texas, Nov. 2016. Association for Computational Linguistics. (Cited on page 1.)
- [49] E. Parzen. On estimation of a probability density function and mode. *Annals of Mathematical Statistics*, 33:1065–1076, 1962. (Cited on page 4.)
- [50] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. (Cited on page 9.)
- [51] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. (Cited on page 1.)

- [52] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever. Improving language understanding by generative pre-training. *OpenAI report*, 2018. (Cited on page 1.)
- [53] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. (Cited on page 1.)
- [54] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. (Cited on page 1.)
- [55] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. (Cited on page 1.)
- [56] J. Reddy. *An introduction to the finite element method*, volume 1221. McGraw-Hill New York, 2004. (Cited on page 11.)
- [57] A. Rives, J. Meier, T. Sercu, S. Goyal, Z. Lin, J. Liu, D. Guo, M. Ott, C. L. Zitnick, J. Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15), 2021. (Cited on page 1.)
- [58] M. Rosenblatt. Remarks on some nonparametric estimates of a density function. *Annals of Mathematical Statistics*, 27:832–837, 1956. (Cited on page 4.)
- [59] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. (Cited on pages 9 and 10.)
- [60] H. Sajjad, F. Dalvi, N. Durrani, and P. Nakov. Poor man’s bert: Smaller and faster transformer models. *arXiv e-prints*, pages arXiv–2004, 2020. (Cited on page 12.)
- [61] M. E. Sander, P. Ablin, M. Blondel, and G. Peyré. Sinkformers: Transformers with doubly stochastic attention. In *International Conference on Artificial Intelligence and Statistics*, pages 3515–3530. PMLR, 2022. (Cited on page 11.)
- [62] V. Sanh, L. Debut, J. Chaumond, and T. Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019. (Cited on page 12.)
- [63] I. Schlag, K. Irie, and J. Schmidhuber. Linear transformers are secretly fast weight programmers. In *International Conference on Machine Learning*, pages 9355–9366. PMLR, 2021. (Cited on pages 4 and 10.)
- [64] P. Shaw, J. Uszkoreit, and A. Vaswani. Self-attention with relative position representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 464–468, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. (Cited on page 11.)
- [65] J. Staniswalis, K. Messer, and D. Finston. Kernel estimators for multivariate regression. *Journal of Nonparametric Statistics*, 3:103–121, 1993. (Cited on page 5.)

- [66] S. Sun, Y. Cheng, Z. Gan, and J. Liu. Patient knowledge distillation for bert model compression. *arXiv preprint arXiv:1908.09355*, 2019. (Cited on page 12.)
- [67] B. Tang and D. S. Matteson. Probabilistic transformer for time series analysis. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. (Cited on page 11.)
- [68] Y. Tay, M. Dehghani, D. Bahri, and D. Metzler. Efficient transformers: A survey. *arXiv preprint arXiv:2009.06732*, 2020. (Cited on page 1.)
- [69] I. Tenney, D. Das, and E. Pavlick. BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy, July 2019. Association for Computational Linguistics. (Cited on page 2.)
- [70] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021. (Cited on page 1.)
- [71] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021. (Cited on page 10.)
- [72] Y.-H. H. Tsai, S. Bai, M. Yamada, L.-P. Morency, and R. Salakhutdinov. Transformer dissection: An unified understanding for transformer’s attention via the lens of kernel. 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)*, pages 4344–4353, Hong Kong, China, Nov. 2019. Association for Computational Linguistics. (Cited on page 11.)
- [73] Y. Tsuzuku, I. Sato, and M. Sugiyama. Lipschitz-margin training: Scalable certification of perturbation invariance for deep neural networks. *Advances in neural information processing systems*, 31, 2018. (Cited on page 5.)
- [74] A. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, 2009. (Cited on page 13.)
- [75] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. (Cited on page 1.)
- [76] J. Vig and Y. Belinkov. Analyzing the structure of attention in a transformer language model. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 63–76, Florence, Italy, Aug. 2019. Association for Computational Linguistics. (Cited on page 2.)
- [77] E. Voita, D. Talbot, F. Moiseev, R. Sennrich, and I. Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5797–5808, Florence, Italy, July 2019. Association for Computational Linguistics. (Cited on page 2.)

- [78] E. Voita, D. Talbot, F. Moiseev, R. Sennrich, and I. Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. *arXiv preprint arXiv:1905.09418*, 2019. (Cited on page 12.)
- [79] M. Wand. Error analysis for general multivariate kernel estimators. *Journal of Nonparametric Statistics*, 2:1–15, 1992. (Cited on page 5.)
- [80] M. Wand and M. Jones. Comparison of smoothing parameterizations in bivariate kernel density estimation. *Journal of the American Statistical Association*, 88:520–528, 1993. (Cited on page 5.)
- [81] S. Wang, B. Li, M. Khabsa, H. Fang, and H. Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020. (Cited on page 11.)
- [82] X. Wang, R. Girshick, A. Gupta, and K. He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. (Cited on page 11.)
- [83] L. Wasserman. *All of Nonparametric Statistics*. Springer, 2006. (Cited on page 7.)
- [84] N. Wiener. *The Fourier Integral and Certain of its Applications*. Cambridge University Press, 1933. (Cited on page 5.)
- [85] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. Salakhutdinov, and Q. V. Le. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*, 2019. (Cited on page 1.)
- [86] S. Zhang and Y. Feng. Modeling concentrated cross-attention for neural machine translation with Gaussian mixture model. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1401–1411, Punta Cana, Dominican Republic, Nov. 2021. Association for Computational Linguistics. (Cited on page 11.)
- [87] H. Zhao, L. Jiang, J. Jia, P. H. Torr, and V. Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16259–16268, 2021. (Cited on page 1.)