

# Manifold Learning and Artificial Intelligence

## Lecture 13

### Mixed Transformer

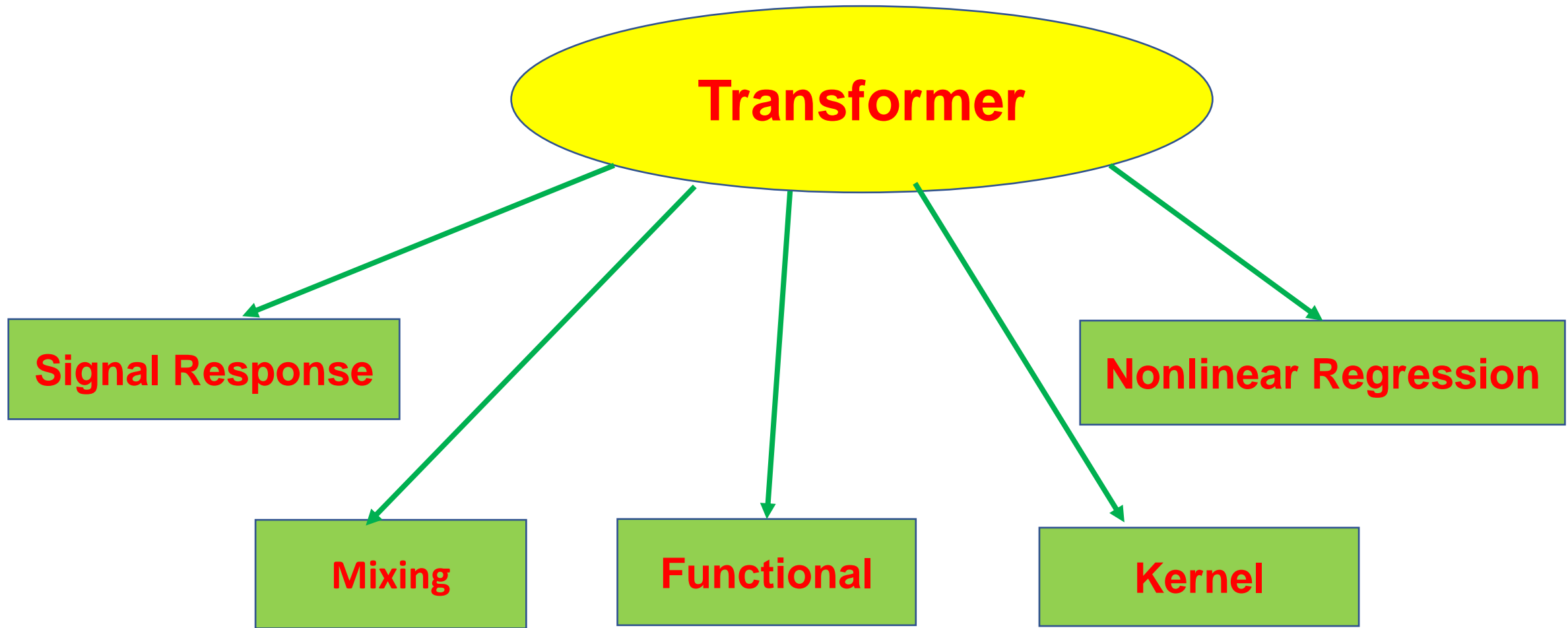
### Spectral and Space Transformer

Momiao Xiong, University of Texas School of Public Health

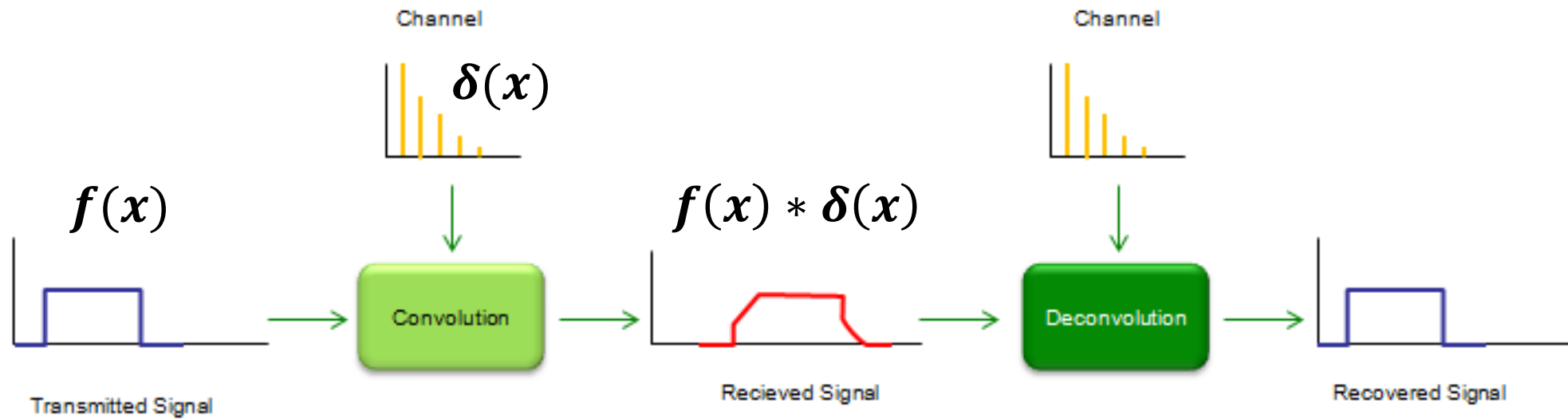
- Time: 10:00 pm, US East Time, 04/22/2023
- 10:00 am, Beijing Time. 04/23/2023

Github Address: <https://ai2healthcare.github.io/>

- **Theoretic foundation of transformer**
- **View transformer as response of system**
- **View transformer as nonlinear regression**
- **Kernel transformer**
- **Generalized Fourier Integral Theorems and their applications to transformer**
- **Functional Model**
- **Mixing MLP**



# 13.1 Signal Response



$$h = f * \delta + \epsilon, \quad f(x, y) * \delta(x, y) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} f(m, n) \delta(x-m, y-n) \quad \text{Causal Model}$$

$$f(x) * \delta(x) = \int_{-\infty}^{\infty} f(u) \delta(x-u) du = \sum_{m=0}^{M-1} f(m) \delta(x-m) \quad h(x) = f(x) * \delta(x) + \epsilon$$

$$x \perp \epsilon$$

$$f(x, y) * \delta(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(u, v) \delta(x-u, y-v) du dv$$

# 13.2. Fourier Transform

- Fourier Transform

## Continuous Function

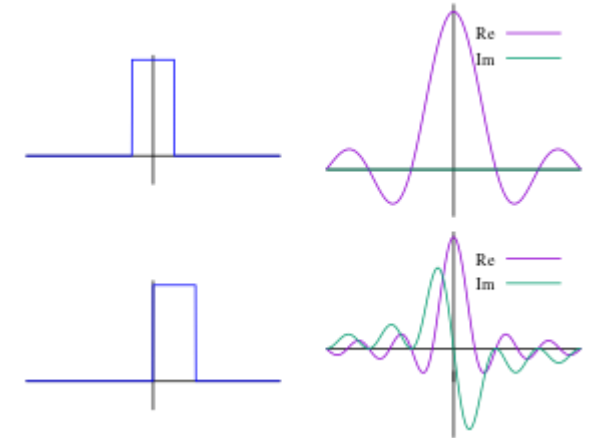
$$\hat{f}(s) = \int_{-R^D}^{R^D} f(x) e^{-i2\pi s^T x} dx, x \in R^D, s \in R^D \quad (F1)$$

$$\hat{f}(u, v) = \int_{-R^D}^{R^D} f(x, y) e^{-i2\pi(u^T x + v^T y)} dx dy \quad (F2)$$

## Discrete Function

$$\hat{f}(s) = \sum_{n=-\infty}^{n=\infty} f(x_n) e^{-i2\pi s^T x_n}$$

$$\hat{f}(u, v) = \sum_{n=-\infty}^{n=\infty} \sum_{m=-\infty}^{m=\infty} f(x_n, y_m) e^{-i2\pi(u^T x_n + v^T y_m)}$$



# 13.2. Fourier Transform

- Inverse Fourier Transform

## Continuous Function

$$f(x) = \int_{-R^D}^{R^D} \hat{f}(s) e^{i2\pi s^T x} ds \quad (\text{F3})$$

$$f(x, y) = \int_{-R^D}^{R^D} \hat{f}(u, v) e^{i2\pi(u^T x + v^T y)} du dv \quad (\text{F4})$$

## Discrete Function

$$f(x) = \sum_{n=-\infty}^{n=\infty} \hat{f}(s_n) e^{-i2\pi x^T s_n}$$

$$f(x, y) = \sum_{n=-\infty}^{n=\infty} \sum_{m=-\infty}^{m=\infty} \hat{f}(u_n, v_m) e^{-i2\pi(x^T u_n + y^T v_m)}$$

# 13.3. Convolution and Fourier Transform

$$\hat{h}(s) = \hat{f}(s)\hat{\delta}(s)$$

$$h = f * \delta + \epsilon$$

$$\hat{h}(x) = F^{-1}(\hat{f}(s)\hat{\delta}(s))$$

$F^{-1}$ : Inverse Foureier Transform

$$f(x, y) * \delta(x, y) < = > \hat{f}(u, v)\hat{\delta}(u, v)$$

$$f(x, y) * \delta(x, y) = F^{-1}\{\hat{f}(u, v)\hat{\delta}(u, v)\}$$

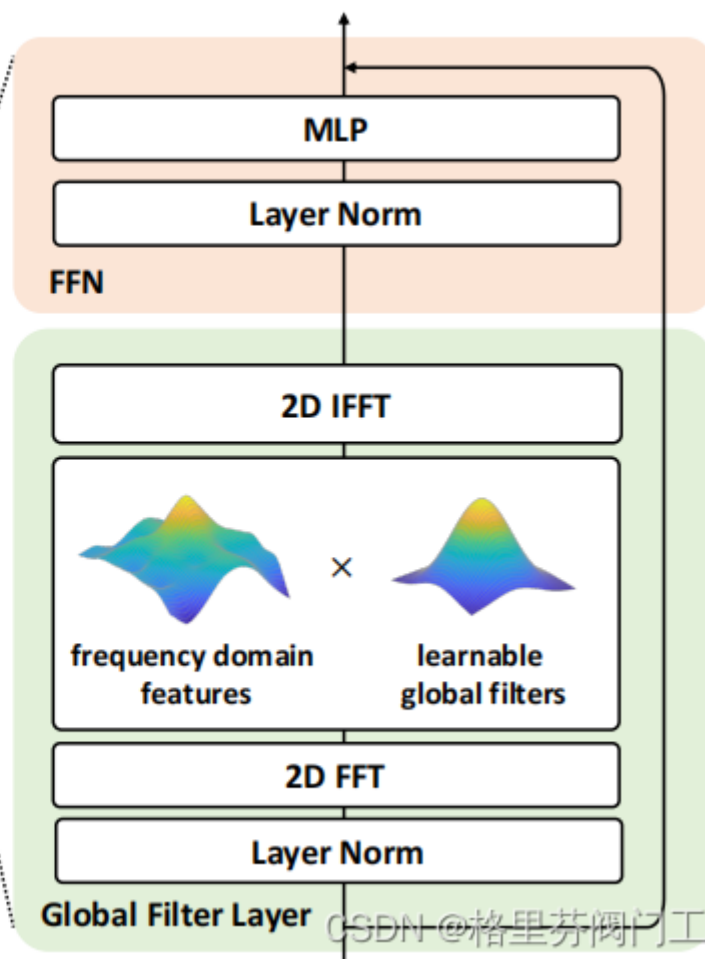
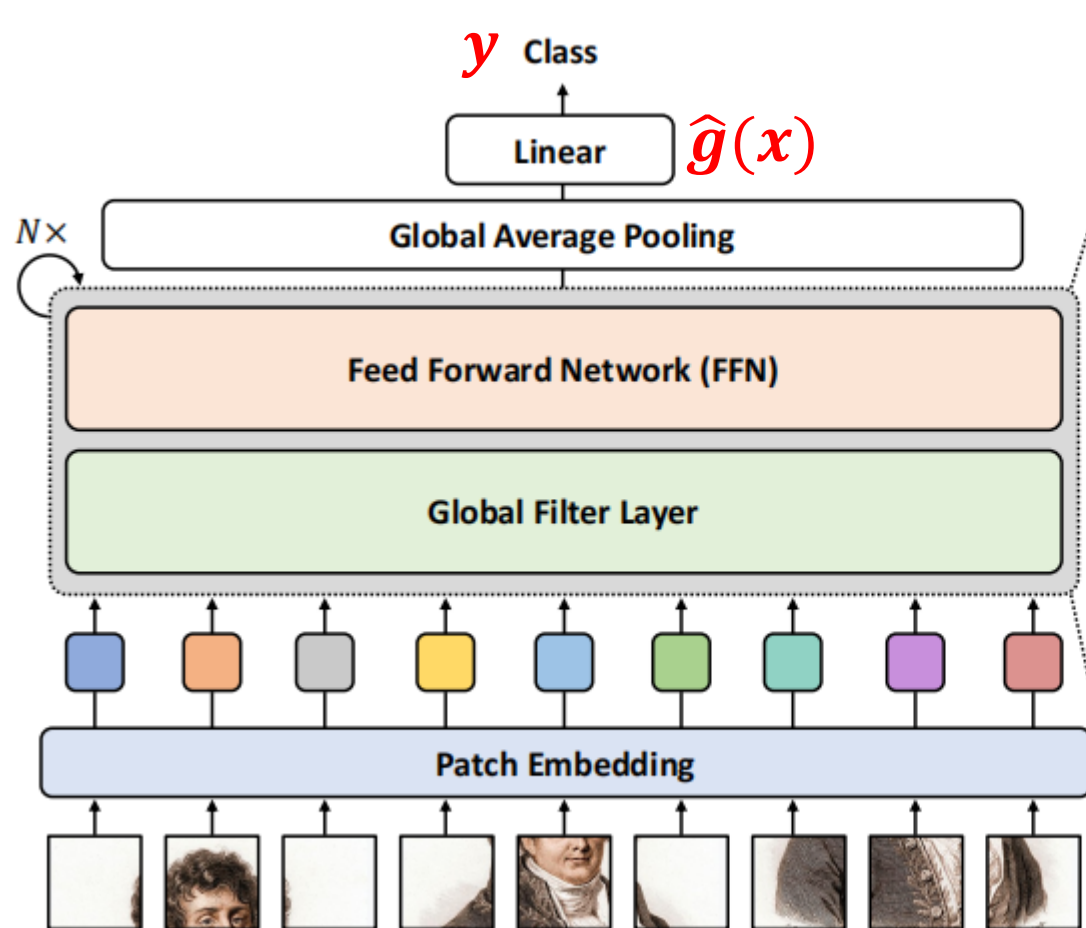
## Causal Test

$$\epsilon = h(x) - \hat{h}(x)$$

$$\epsilon \perp\!\!\!\perp x$$

# 13.4. Global Filter Networks for Image Classification

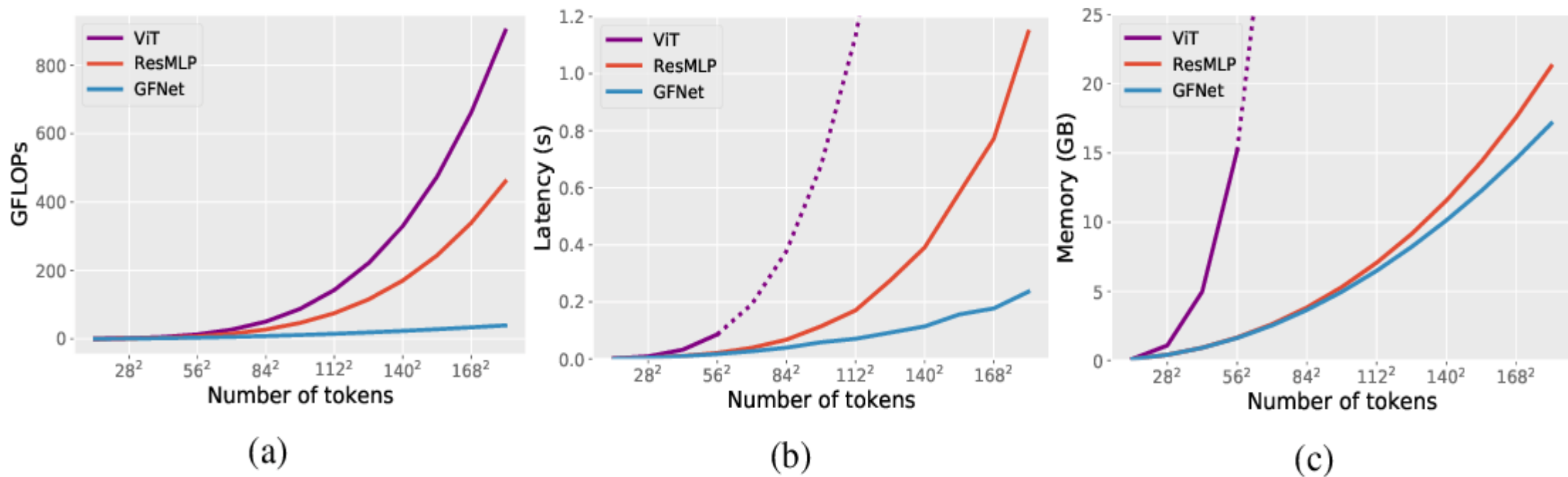
Code is available at <https://github.com/raoyongming/GFNet>.



**Causal Test**

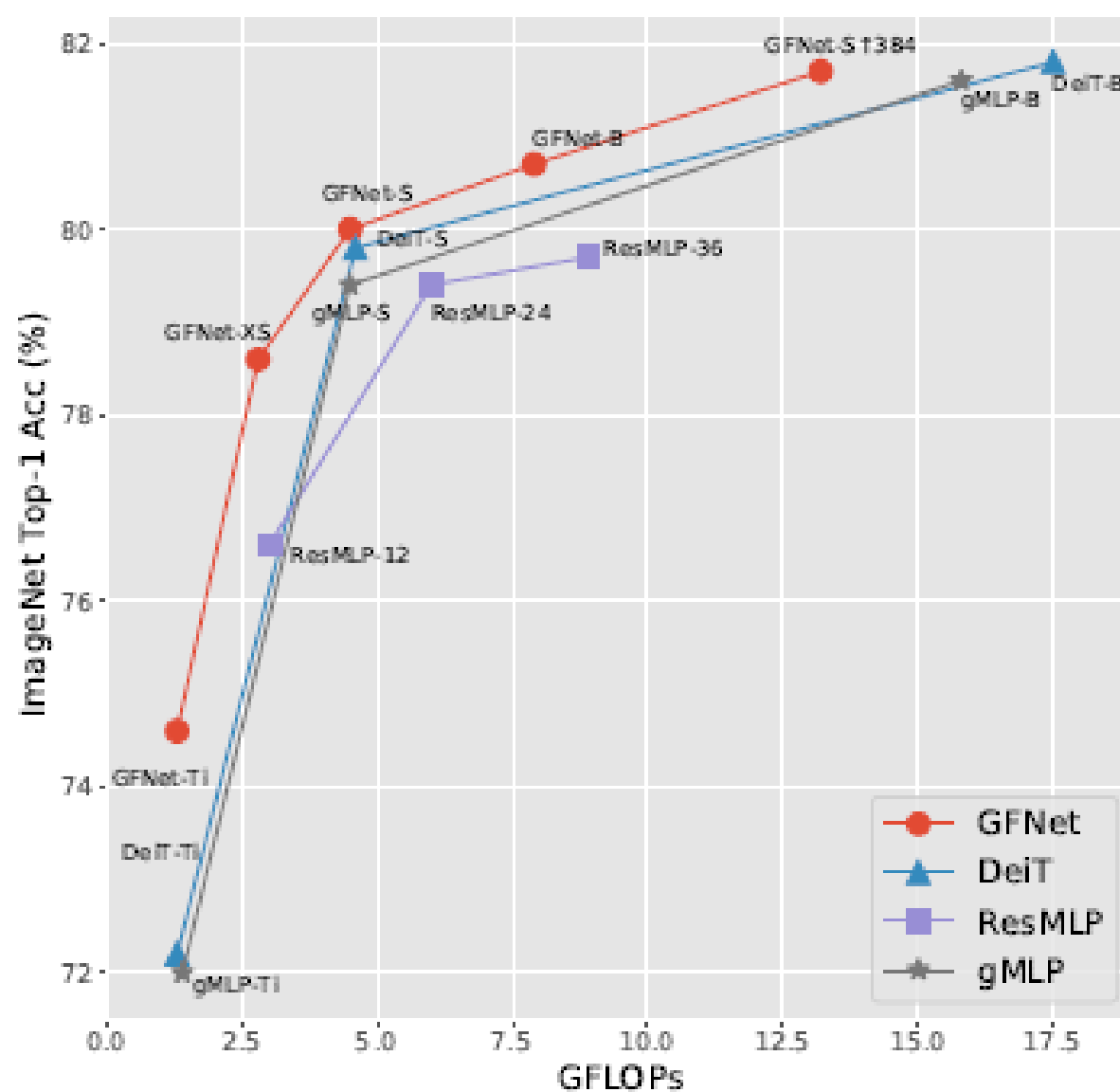
$$\varepsilon = y - \hat{g}(x)$$
$$\varepsilon \perp\!\!\!\perp x$$





**Figure 2: Comparisons among GFNet, ViT [10] and ResMLP in (a) FLOPs (b) latency and (c) GPU memory with respect to the number of tokens (feature resolution). The dotted lines indicate the estimated values when the GPU memory has run out. The latency and GPU memory is measured using a single NVIDIA RTX 3090 GPU with batch size 32 and feature dimension 384.**

**FLOPs: floating point operations per second**



Top-N Accuracy takes the N model predictions with higher probability. If one of them is a true label, it classifies the prediction as correct.

Figure 3: ImageNet acc. vs model complexity

## 13. 5. FourierFormer: Transformer Meets Generalized Fourier Integral Theorem

- 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.
- the generalized Fourier integral kernels can automatically capture dependency and remove the need to tune the covariance matrix.
- Our PyTorch code with documentation can be found at [https://github.com/minhtannguyen/FourierFormer\\_NeurIPS](https://github.com/minhtannguyen/FourierFormer_NeurIPS)

Tan M. Nguyen←  
Department of Mathematics  
University of California, Los Angeles  
tanmnguyen89@ucla.edu

Khai Nguyen  
Department of Statistics and Data Sciences  
University of Texas at Austin  
khainb@utexas.edu

Minh Pham←  
Department of Mathematics  
University of California, Los Angeles  
minhrose@ucla.edu

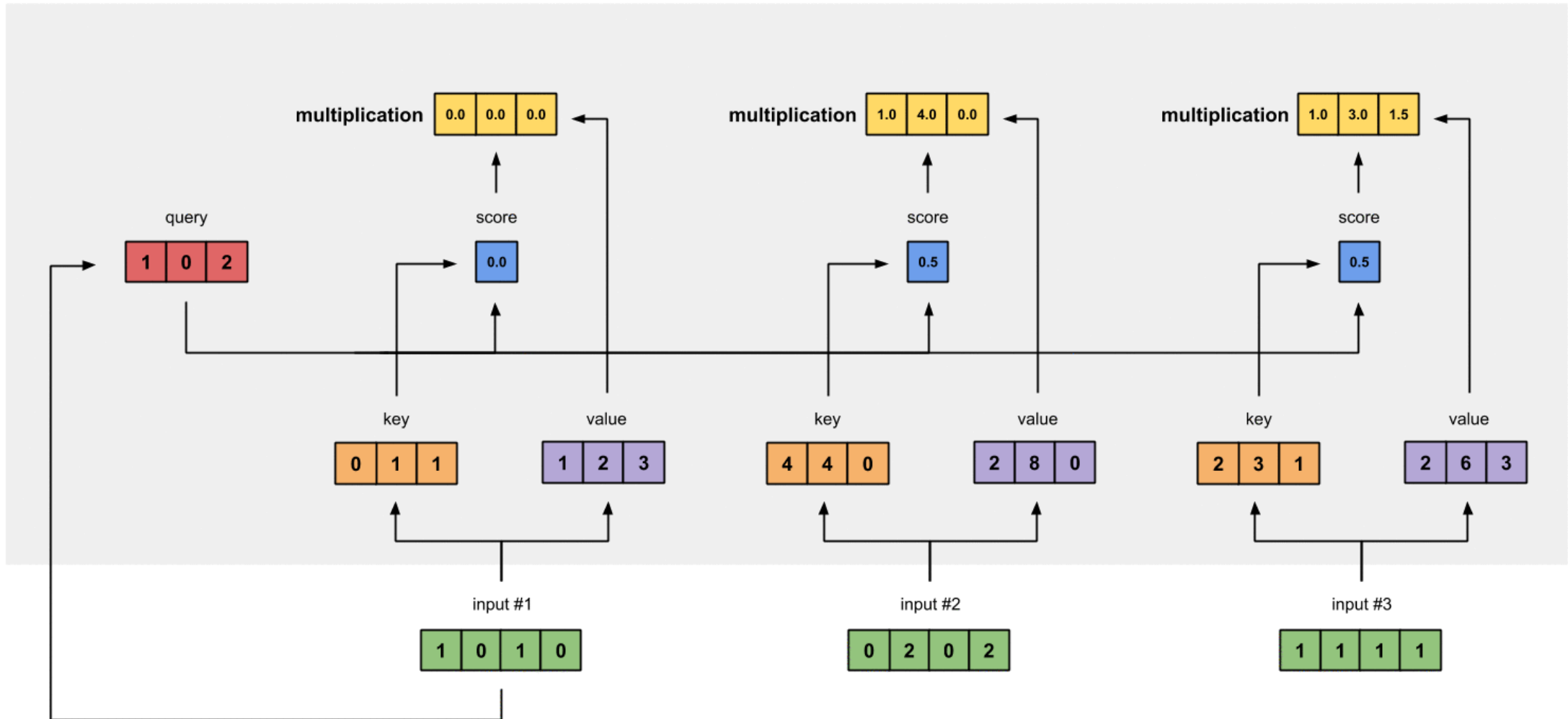
Stanley J. Osher←←  
Department of Mathematics  
University of California, Los Angeles  
sjo@math.ucla.edu

Tam Nguyen  
Department of ECE  
Rice University  
nguyenminhtam9520@gmail.com

Nhat Ho←←  
Department of Statistics and  
Data Sciences  
University of Texas at Austin  
minhnhat@utexas.edu

# Scheme of Self Attention

Self-attention



## 13.5.1. A Nonparametric Regression Interpretation of Self-attention

- The key vectors  $k_j$  and value vectors  $v_j$  are training inputs and training targets
- The query vectors  $q_i$  and the output vectors  $h_i$  form a set of new inputs and their corresponding targets that need to be estimated.

- **nonparametric regression model:**

$$v_j = f(k_j) + \varepsilon_j, j = 1, \dots, N, \quad k_j \sim p, (v_j, k_j) \sim p$$

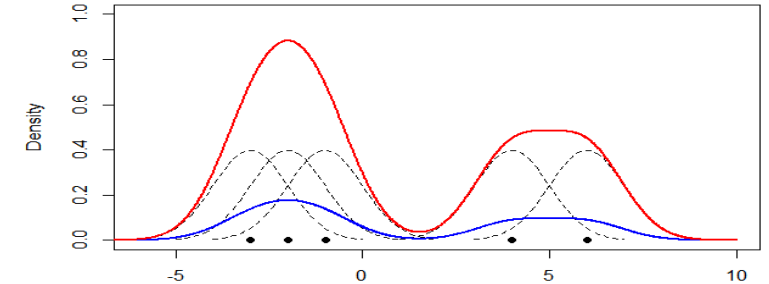
- **Nadaraya–Watson’s nonparametric kernel regression**

$$f(k) = E[v|k]$$
$$E[v|k] = \int_{R^D} v p(v|k) dv = \int_{R^D} \frac{v p(v, k)}{p(k)} dv = \begin{bmatrix} \int_R \frac{v_1 p(v, k)}{p(k)} dv_1 \\ \vdots \\ \int_R \frac{v_N p(v, k)}{p(k)} dv_N \end{bmatrix}$$

- Kernel density estimator**

Using the isotropic Gaussian kernel with bandwidth  $\sigma$ , we have the following estimators:

$$\hat{p}_{\sigma}(v, k) = \frac{1}{N} \sum_{j=1}^N \varphi_{\sigma}(v - v_j) \varphi_{\sigma}(k - k_j)$$



$$\hat{p}_{\sigma}(k) = \frac{1}{N} \sum_{j=1}^N \varphi_{\sigma}(k - k_j)$$

the isotropic multivariate Gaussian density function with diagonal covariance matrix

- Given the kernel density estimators, we obtain the following estimation of the function  $f$ :

$$\begin{aligned} \hat{f}_{\sigma}(k) &= \int_{R^D} \frac{v \hat{p}_{\sigma}(v, k)}{\hat{p}_{\sigma}(k)} dv = \int_{R^D} \frac{v \sum_{j=1}^N \varphi_{\sigma}(v - v_j) \varphi_{\sigma}(k - k_j)}{\sum_{j'=1}^N \varphi_{\sigma}(k - k_{j'})} dv \\ &= \sum_{j=1}^N \frac{\varphi_{\sigma}(k - k_j)}{\sum_{j'=1}^N \varphi_{\sigma}(k - k_{j'})} \int_{R^D} v \varphi_{\sigma}(v - v_j) dv = \sum_{j=1}^N \frac{\varphi_{\sigma}(k - k_j)}{\sum_{j'=1}^N \varphi_{\sigma}(k - k_{j'})} v_j \quad (1) \end{aligned}$$

### 13.5.2. Connection between Self-Attention and nonparametric regression

The query vectors  $q_i$  and the output vectors  $h_i$  form a set of new inputs and their corresponding targets.

By plugging the query vectors  $q_i$  into the function  $\hat{f}_\sigma$  in equation (1), we obtain that

$$h_i = \hat{f}_\sigma(q_i) = \sum_{j=1}^N \frac{\varphi_\sigma(q_i - k_j)}{\sum_{j'=1}^N \varphi_\sigma(q_i - k_{j'})} v_j \quad (2)$$

Note that

$$\frac{\varphi_\sigma(q_i - k_j)}{\sum_{j'=1}^N \varphi_\sigma(q_i - k_{j'})} = \frac{\exp\left\{-\frac{\|q_i - k_j\|^2}{2\sigma^2}\right\}}{\sum_{j'=1}^N \exp\left\{-\frac{\|q_i - k_{j'}\|^2}{2\sigma^2}\right\}} = \frac{\exp\left\{-\frac{\|q_i\|^2 + \|k_j\|^2}{2\sigma^2}\right\} \exp\left\{\frac{q_i^T k_j}{\sigma^2}\right\}}{\sum_{j'=1}^N \exp\left\{-\frac{\|q_i\|^2 + \|k_{j'}\|^2}{2\sigma^2}\right\} \exp\left\{\frac{q_i^T k_{j'}}{\sigma^2}\right\}} \quad (3)$$

Assume that  $k_j$  is normalized, then equation (3) is reduced to

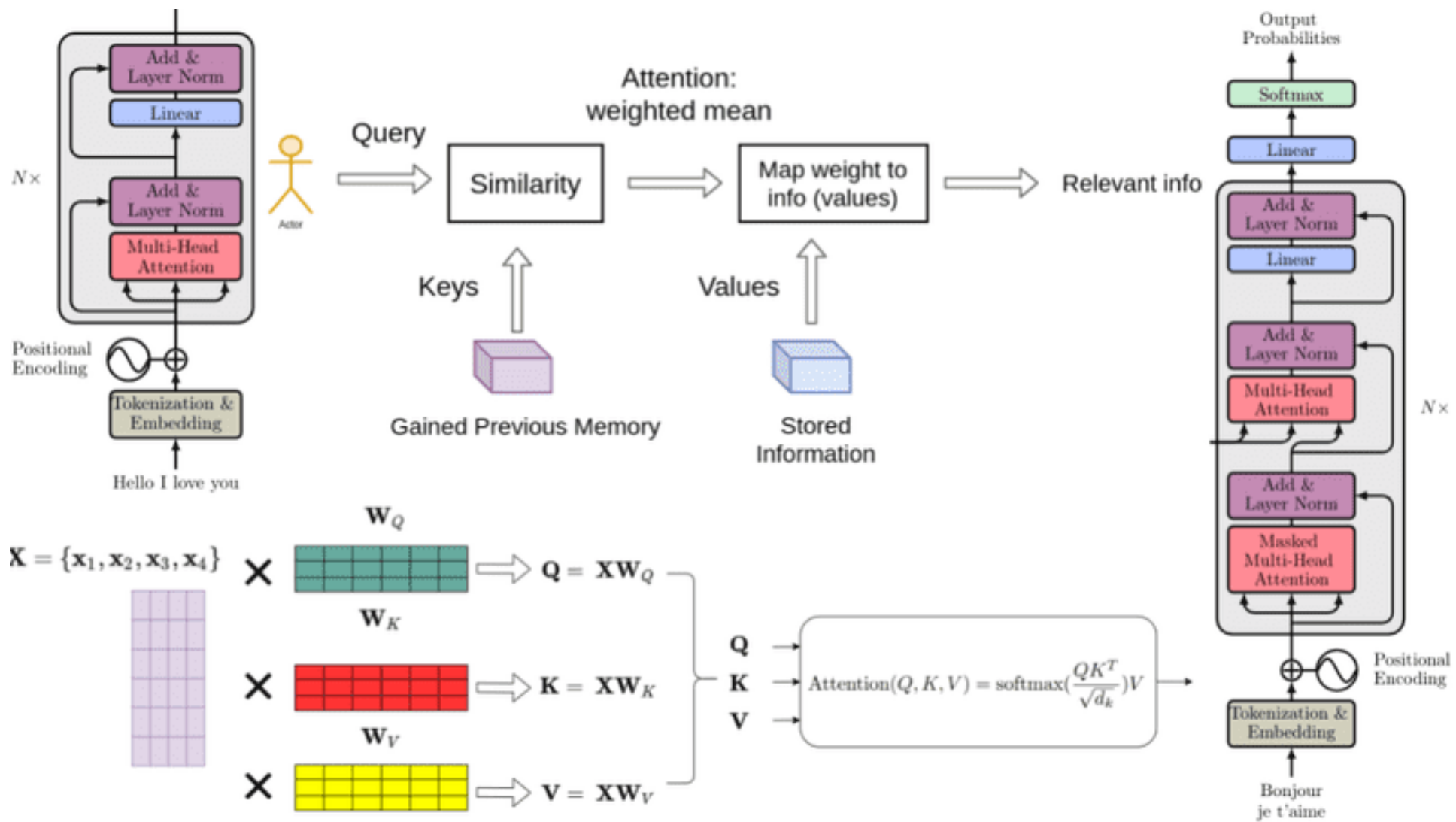
$$\frac{\varphi_{\sigma}(q_i - k_j)}{\sum_{j'=1}^N \varphi_{\sigma}(q_i - k_{j'})} = \frac{\exp\left\{\frac{q_i^T k_j}{\sigma^2}\right\}}{\sum_{j'=1}^N \exp\left\{\frac{q_i^T k_{j'}}{\sigma^2}\right\}} \quad (4)$$

Substituting equation (4) into equation (2), we obtain

$$\begin{aligned} h_i = \hat{f}_{\sigma}(q_i) &= \sum_{j=1}^N \frac{\exp\left\{\frac{q_i^T k_j}{\sigma^2}\right\}}{\sum_{j'=1}^N \exp\left\{\frac{q_i^T k_{j'}}{\sigma^2}\right\}} v_j \\ &= \sum_{j=1}^N \text{softmax}\left(\frac{q_i^T k_j}{\sigma^2}\right) v_j \end{aligned} \quad (5)$$

Choose  $\sigma^2 = \sqrt{D}$  where  $D$  is the dimension of  $q_i$  and  $k_j$ , equation (5) matches equation of self-attention.





- **Limitation**

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  $q_i$  and  $k_j$ .
- Using multivariate Gaussian kernels with dense covariance matrices can help capture such correlations; however, choosing good covariance matrices is challenging and inefficient.
- 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.

## 13.5.3. FourierFormer: Transformer via Generalized Fourier Integral Theorem

### 13.5.3.1. Generalized Fourier Integral Theorems and Their Applications

Fourier integral theorem is a combination of Fourier transform and Fourier inverse transform. Let  $\mathbf{p} \in L_1(\mathbb{R}^D)$ , using equations (F1) and (F3), we obtain the Fourier integral theorem:

$$\begin{aligned} \mathbf{p}(\mathbf{k}) &= \int_{-\mathbb{R}^D}^{\mathbb{R}^D} \hat{\mathbf{p}}(\mathbf{s}) e^{i2\pi \mathbf{s}^T \mathbf{k}} d\mathbf{s} = \int_{-\mathbb{R}^D}^{\mathbb{R}^D} \int_{-\mathbb{R}^D}^{\mathbb{R}^D} \mathbf{p}(\mathbf{y}) e^{-i2\pi \mathbf{s}^T \mathbf{y}} d\mathbf{y} e^{i2\pi \mathbf{s}^T \mathbf{k}} d\mathbf{s} \\ &= \int_{-\mathbb{R}^D}^{\mathbb{R}^D} \int_{-\mathbb{R}^D}^{\mathbb{R}^D} \mathbf{p}(\mathbf{y}) e^{i2\pi \mathbf{s}^T (\mathbf{k} - \mathbf{y})} d\mathbf{y} d\mathbf{s} = \frac{1}{(2\pi)^D} \int_{-\mathbb{R}^D}^{\mathbb{R}^D} \int_{-\mathbb{R}^D}^{\mathbb{R}^D} \mathbf{p}(\mathbf{y}) e^{i\mathbf{s}^T (\mathbf{k} - \mathbf{y})} d\mathbf{y} d\mathbf{s} \end{aligned} \quad (6)$$

Note that

$$\begin{aligned}\int_{-R^D}^{R^D} e^{is^T(k-y)} d\mathbf{s} &= \int_{-R^D}^{R^D} e^{i[\sum_{j=1}^D s_j(k_j-y_j)]} d\mathbf{s} \\ &= \prod_{j=1}^D \int_{-R}^R e^{is_j(k_j-y_j)} ds_j \\ &= \prod_{j=1}^D \int_{-R}^R \frac{1}{i(k_j-y_j)} d[e^{is_j(k_j-y_j)}] \\ &= 2^D \prod_{j=1}^D \frac{\text{Sin}[R(k_j-y_j)]}{(k_j-y_j)}\end{aligned}\tag{8}$$

Substituting equation (8) into equation (6) yields

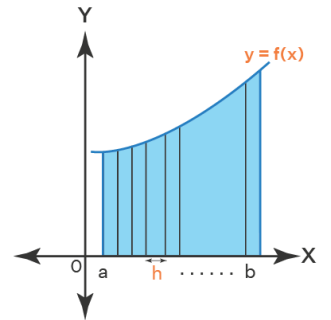
$$\mathbf{p}(\mathbf{k}) = \frac{1}{\pi^D} \int_{R^D} \prod_{j=1}^D \frac{\text{Sin}[R(k_j-y_j)]}{(k_j-y_j)} \mathbf{p}(\mathbf{y}) d\mathbf{y}\tag{9}$$

## 13.5.3.2. Generalized Fourier integral estimator

Definite Integral

- Generalized Fourier integral theorem:

$$p(k) = \lim_{n \rightarrow \infty} p_R^\emptyset = \lim_{n \rightarrow \infty} \frac{R^D}{A^D} \int_{R^D} \prod_{j=1}^D \emptyset \left( \frac{\sin [R(k_j - y_j)]}{R(k_j - y_j)} \right) p(y) dy \quad (10)$$



$$A = \int_R \emptyset \left( \frac{\sin(z)}{z} \right) dz, \emptyset: R \rightarrow R \text{ is a given function.}$$

- Generalized Fourier density estimator

Let  $p(y) = \frac{1}{N}$ , then we have

$$\begin{aligned} & \int_R g(y_j) p(y_j) dy_j \\ &= \sum_{i=1}^N g(y_{ij}) \frac{1}{N} \end{aligned}$$

$$\int_{R^D} \prod_{j=1}^D \emptyset \left( \frac{\sin [R(k_j - y_j)]}{R(k_j - y_j)} \right) p(y) dy = \frac{1}{N} \sum_{i=1}^N \prod_{j=1}^D \emptyset \left( \frac{\sin [R(k_j - y_{ij})]}{R(k_j - y_{ij})} \right) \quad (11)$$

Assume that  $\mathbf{k}_1, \dots, \mathbf{k}_i, \dots, \mathbf{k}_D \in R^D$  and be i.i.d, where  $\mathbf{k}_i = (k_{i1}, \dots, k_{ij}, \dots, k_{iD})$

Substituting equation (11) into equation (10) yields the Generalized Fourier density estimator:

$$p_{N,R}^{\emptyset}(\mathbf{k}) = \frac{R^D}{NA^D} \sum_{i=1}^N \prod_{j=1}^D \emptyset \left( \frac{\text{Sin}[R(k_j - k_{ij})]}{R(k_j - k_{ij})} \right) \quad (12)$$

where  $\mathbf{k} = (k_1, \dots, k_j, \dots, k_D)$

# 13.5.3.3. FourierFormer: Transformers with Fourier Attentions

Recall the nonparametric regression model

$$v = f(k) + \varepsilon$$

- The Nadaraya–Watson estimator of the function  $f$

$$f_{N,R}(k) = \frac{\sum_{i=1}^N v_i \prod_{j=1}^D \phi\left(\frac{\sin[R(k_j - k_{ij})]}{R(k_j - k_{ij})}\right)}{\sum_{i=1}^N \prod_{j=1}^D \phi\left(\frac{\sin[R(k_j - k_{ij})]}{R(k_j - k_{ij})}\right)} \quad (13)$$

$$k = (k_1, \dots, k_j, \dots, k_D)$$

- **FourierFormer**

Given the generalized Fourier nonparametric regression estimator  $f_{N,R}$  in equation (13), by plugging the query values  $q_1, \dots, q_N$  into that function, we obtain the following definition of the Fourier attention:

$$h_i = f_{N,R}(q_i) = \frac{\sum_{i=1}^N v_i \prod_{j=1}^D \phi \left( \frac{\sin [R(q_{ij} - k_{ij})]}{R(q_{ij} - k_{ij})} \right)}{\sum_{i=1}^N \prod_{j=1}^D \phi \left( \frac{\sin [R(q_{ij} - k_{ij})]}{R(q_{ij} - k_{ij})} \right)} \quad (14)$$

**Definition 2 (FourierFormer)**

**Define a** FourierFormer as a transformer that uses Fourier attention to capture dependency between tokens in the input sequence and the correlation between features in each token.



## 13.5.3.4. Results

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)		
Baseline (Dot-product) (Small)	33.15			34.29		
Fourier Former (small)	31.86			32.85		
Baseline (Medium)	27.90			29.60		
Fourier Former (medium)	26.51			28.01		

## Image Classification on ImageNet

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
Baseline Deit	72.23	91.13
Fourier Former	73.25	91.66