# PhaseDNN - A Parallel Phase Shift Deep Neural Network for Adaptive Wideband Learning

Wei Cai[a], Xiaoguang Li[b,a], Lizuo Liu[a]

[a]*Department of Mathematics, Southern Methodist University, Dallas, TX 75275, USA.*
[b]*LCSM, Ministry of Education, School of Mathematics and Statistics, Hunan Normal University, Changsha, Hunan 410081, P. R. China*

## Abstract

In this paper, we propose a phase shift deep neural network (PhaseDNN) which provides a wideband convergence in approximating a high dimensional function during its training of the network. The PhaseDNN utilizes the fact that many DNN achieves convergence in the low frequency range first, thus, a series of moderately-sized of DNNs are constructed and trained in parallel for ranges of higher frequencies. With the help of phase shifts in the frequency domain, implemented through a simple phase factor multiplication on the training data, each DNN in the series will be trained to approximate the target function's higher frequency content over a specific range. Due to the phase shift, each DNN achieves the speed of convergence as in the low frequency range. As a result, the proposed PhaseDNN system is able to convert wideband frequency learning to low frequency learning, thus allowing a uniform learning to wideband high dimensional functions with frequency adaptive training. Numerical results have demonstrated the capability of PhaseDNN in learning information of a target function from low to high frequency uniformly.

*Keywords:* Neural network, phase shift, wideband data.

## 1. Introduction

Deep neural network (DNNs) have shown great potential in approximating high dimensional functions without suffering the curse of dimensionality of traditional approximations based on Lagrangian interpolation or spectral methods. Recently, it has been shown in [1, 2] that some common NNs, including fully connected and CNN with tanh and ReLU activation functions,

demonstrated a frequency dependent convergence behavior, namely, the DNN during the training are able to approximate the low frequency components of the targeted functions first before the higher frequency components, named as F-Principle of DNNs [1]. The stalling of DNN convergence in the later stage of training could be mostly coming from learning the high frequency component. This F-principle behavior of DNN is just the opposite to the convergence behavior of the traditional multigrid method (MGM) [3] for approximating the solutions of PDEs where the convergence occurs first in the higher frequency end of the spectrum, due to the smoothing operator employed in the MGM. The MGM takes advantage of this fast high frequency error reduction in the smoothing iteration cycles and restrict the original solution of a fine grid to a coarse grid, then continuing the smoothing iteration on the coarse grid to reduce the higher end frequency spectrum in the context of the coarse grid. This downward restriction can be continued until errors over all frequency are reduced by a small number of iteration on each level of the coarse grids.

In this work, we propose a wideband DNN in error reductions in the approximation for all frequencies of the targeted function by making use of the faster convergence in the low frequencies of the DNN during its training. We first will train a neural network $T_\star(\boldsymbol{x}, \theta^{(n_0)})$ until it will achieve a sufficient error reduction for the frequency range $|k| < K_0$ for a prescribed $K_0$ within $n_0$ training epoches. In order to further learn the high frequency component of the target function, we decompose the residual, i.e. the difference between the target function and $T_\star(\boldsymbol{x}, \theta^{(n_0)})$, into a series of functions with frequency spectrum through a partition of the unit (POU) in the $k$-space. To learn each of these function of specific frequency range, we will employ a phase shift in the $k$-space to translate its frequency to the range $|k| < K_0$, the phase shifted function now only has low frequency content and can be learned by common DNN with a small number of epoches of training. This series of DNNs with phase shifts will be named Phase-shift deep Neural Network (PhaseDNN).

It should be noted that during the construction of PhaseDNN, only the original set of training data will be needed, but it will be modified to account for the frequency shifts, implemented simply by multiplying the training data by an phase factor $e^{\pm \omega_j x}$ on the residuals from the previous neural network approximations. We also allude the way PhaseDNN achieve approximation over increasing range of frequency with each additional network to that of wavelet approximations, which generate better and better resolution of approximation by adding wavelet subspaces produced by dialation of mother

2

wavelet to account for higher frequency [4][5]. In the case of PhaseDNN though phase shifts are used to cover higher frequencies. The substraction between the target function provided by the training data and the lower frequency network approximation is similar to the strategy employed in construction of multi-resolution interpolation for PDE solutions[5].

The rest of the paper will be organized as follows. In section 2, we will review the fast low frequency convergence property of neural network with tanh activation function. Section 3 will present the new phase shift deep neural network. Finally, a conclusion and discussions will be given in Section 4.

## 2. Deep Neural Network and frequency dependent training of DNN

A deep neuron network (DNN) is a sequential alternative composition of linear functions and nonlinear activation functions. Given $m, n \geqslant 1$, let $\Theta(\boldsymbol{x}) : \mathbb{R}^n \to \mathbb{R}^m$ is a linear function with the form $\Theta(\boldsymbol{x}) = \boldsymbol{W}\boldsymbol{x} + \boldsymbol{b}$, where $\boldsymbol{W} = (w_{ij}) \in \mathbb{R}^{n \times m}$, $\boldsymbol{b} \in \mathbb{R}^m$. The nonlinear activation function $\sigma(u) : \mathbb{R} \to \mathbb{R}$. By applying this function componentwisely, we can extend activation function to $\sigma(u) : \mathbb{R}^n \to \mathbb{R}^n$ naturally. A DNN with $L + 1$ layers can be expressed as a compact form

$$
\begin{aligned}
T(\boldsymbol{x}) &= T^L(\boldsymbol{x}), \\
T^l(\boldsymbol{x}) &= [\Theta^l \circ \sigma](T^{l-1}(\boldsymbol{x})), \quad l = 1, 2, \dots L,
\end{aligned}
\tag{1}
$$

with $T^0(\boldsymbol{x}) = \Theta^0(\boldsymbol{x})$. Equivalently, we can also express it explicitly:

$$
T(\boldsymbol{x}) = \Theta^L \circ \sigma \circ \Theta^{L-1} \circ \sigma \cdots \circ \Theta^1 \circ \sigma \circ \Theta^0(\boldsymbol{x}).
\tag{2}
$$

Here, $\Theta^l(\boldsymbol{x}) = \boldsymbol{W}^l\boldsymbol{x} + \boldsymbol{b}^l : \mathbb{R}^{n^l} \to \mathbb{R}^{n^{l+1}}$ are linear functions. This DNN is also said to have $L$ hidden layers. The $l$-th layer has $n^l$ neurons.

In the following text, we only consider DNN with one dimensional input layer and one dimensional output layer. Namely, $\Theta^0(x) : \mathbb{R} \to \mathbb{R}^{n^1}$, $\Theta^L(\boldsymbol{x}) : \mathbb{R}^{n^{L-1}} \to \mathbb{R}$. That is, $T(\boldsymbol{x}) : \mathbb{R} \to \mathbb{R}$. The activation function is chosen to be $\sigma(u) = \tanh(u)$.

In this paper, we consider the problem of approximating a function by DNN through training. Given a function $f(x) \in \mathbb{L}^1(\mathbb{R}) \cap \mathbb{L}^2(\mathbb{R})$ (This choice is only to ensure the existence of Fourier transform and loss function for

analysis reason.), we are going to minimize

$$L(\boldsymbol{W}^0, \boldsymbol{b}^1, \boldsymbol{W}^1, \boldsymbol{b}^1, \ldots, \boldsymbol{W}^L, \boldsymbol{b}^L) = \|f(\boldsymbol{x}) - T(\boldsymbol{x})\|_2^2 = \int_{-\infty}^{+\infty} |f(\boldsymbol{x}) - T(\boldsymbol{x})|^2 \, x.$$
(3)

We consider this problem in frequency space. We define Fourier transform and its inverse of a function $f(\boldsymbol{x})$ by

$$\mathcal{F}[f](k) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} f(x) e^{-ikx} \, \mathrm{d}x, \qquad \mathcal{F}^{-1}[\widehat{f}](x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} \hat{f}(k) e^{ikx} \, \mathrm{d}k.$$
(4)

Let $D(k) = \mathcal{F}[T - f](k) = A(k) e^{i\varphi(k)}$, $L(k) = |D(k)|^2$. By Paseval's equality,

$$L(\boldsymbol{W}^0, \boldsymbol{b}^1, \boldsymbol{W}^1, \boldsymbol{b}^1, \ldots, \boldsymbol{W}^L, \boldsymbol{b}^L) = \int_{-\infty}^{+\infty} L(k) k.$$
(5)

Let $\theta$ denotes all the parameters in DNN. Namely,

$$\theta = (\boldsymbol{W}_1^0, \boldsymbol{W}_2^0 \ldots, \boldsymbol{W}_{n^1}^0, \boldsymbol{b}^0, \boldsymbol{W}_{11}^1, \ldots, \boldsymbol{W}_{n^1 n^2}^1, \boldsymbol{b}_1^1 \ldots \boldsymbol{b}_{n^2}^1 \ldots) \in \mathbb{R}^p.$$

Here, $p = n^1 + 1 + (n^1 + 1) \times n^2 + (n^2 + 1) \times n^3 + \ldots (n^L + 1)$ is the number of the parameters. The F-principle states the relative changing rate of $\frac{\partial L(k)}{\partial \theta_j}$ for different frequency $k$. We cite the following result from [2]

**Theorem 1.** *Given a function $f(\boldsymbol{x}) \in \mathbb{L}^1(\mathbb{R}) \cap \mathbb{L}^2(\mathbb{R})$ and a DNN function $T(\boldsymbol{x})$ defined by (1), we scale $T(\boldsymbol{x})$ by $\Theta^0(\boldsymbol{x}) = \varepsilon \boldsymbol{W}^0 \boldsymbol{x} + \boldsymbol{b}$, where $\varepsilon > 0$ is a small parameter. Suppose $\boldsymbol{W}_{rs}^l \neq 0$ for all $0 \leqslant l \leqslant L - 1$, $1 \leqslant r \leqslant n^l$ and $1 \leqslant s \leqslant n^{l+1}$. For any frequency $k_1$ and $k_2$ such that $|\mathcal{F}[f](k_1)| > 0$, $|\mathcal{F}[f](k_2)| > 0$ and $|k_2| > |k_1| > 0$, there exists a $\delta > 0$ such that when $\varepsilon < \delta$,*

$$\frac{\mu\left(\{\theta| \left|\frac{\partial L(k_1)}{\partial \theta_j}\right| > \left|\frac{\partial L(k_2)}{\partial \theta_j}\right|\} \cap B_\delta\right)}{\mu(B_\delta)} > 1 - C \exp(-c/\delta)$$

*holds for all $j$, where $B_\delta = \{\boldsymbol{x} \in \mathbb{R}^p | |\boldsymbol{x}| < \delta\}$, $\mu(\cdot)$ is the Lebesgue measure and $c > 0$, $C > 0$ are two constants.*

## 3. Parallel Phase shift DNN (PhaseDNN)

The F-principle and estimation of $L(k)$ in Theorem 1 suggest that when we apply a gradient-based training method to loss function, the low frequency

4

part of loss function converges faster than the high frequency part. So when we train a DNN $T_\star(x)$ to approximate $f(x)$, there exists positive frequency number $K_0$ and an integer $n_0$ such that after $n_0$ steps of training, the Fourier transform of training function $\mathcal{F}[T_\star](\boldsymbol{k}; \theta^{(n_0)})$ should be a good approximation of $\mathcal{F}[f](\boldsymbol{k})$ for $|k| < K_0$, where $\theta^{(n_0)}$ are parameters obtained after $n_0$ steps of training.

- Frequency selection kernel $\varphi_j^\vee(x)$

In order to speed up the learning of the higher frequency content of the target function $f(x)$, we will employ a phase shift technique to translate higher frequency spectrum $\widehat{f}(k)$ to the frequency range of $[-K_0, K_0]$. Such a shift in frequency is a simple phase factor multiplication on the training data in the physical space, which can be easily implemented.

For a given $\Delta k$, let us assume that

$$\operatorname{supp} \widehat{f}(k) \subset [-m\Delta k, m\Delta k].$$

Before we can carry out this scheme, we construct a mesh for the interval $[-m\Delta k, m\Delta k]$ by

$$\omega_j = j\Delta k, j = -m, \cdots, m, \tag{6}$$

Then, we introduce a POU for the domain $[-M, M]$ associated with the mesh as

$$1 = \sum_{j=-m}^{m-1} \varphi_j(k), \ \ k \in [-M, M]. \tag{7}$$

The simplest choice of $\varphi_j(k)$ is

$$\varphi_j(k) = \varphi(\frac{k - \omega_j}{\Delta k}), \tag{8}$$

where $\varphi(k)$ is just the characteristic function of $[-1, 1]$, i.e.,

$$\varphi(k) = \chi_{[0,1]}(k).$$

The inverse Fourier transform of $\varphi(k)$, indicated by $\vee$, is

$$\varphi^\vee(x) = \frac{1}{\sqrt{2\pi}} e^{i\frac{x}{2}} \frac{\sin x}{x}. \tag{9}$$

5

For a smoother function $\varphi(k)$ whose inverse Fourier transform, which will act as a frequency selection kernel in $x$-space, has a faster decay condition, we could use B-Splines $B_m(k)$, defined by the following recursive convolutions. Namely,

$$B_1(k) = \chi_{[0,1]}(k),$$

$$B_m(k) = B_{m-1}(k) * \chi_{[0,1]}(k), m = 2, \cdots . \tag{10}$$

It can be easily shown that

$$\text{supp } B_m(k) = [0, m],$$

and the inverse Fourier transform of $B_m(k)$ is

$$B_m^\vee(x) = (B_1^\vee(x))^m = \left( \frac{1}{\sqrt{2\pi}} e^{i\frac{x}{2}} \frac{\sin\frac{x}{2}}{\frac{x}{2}} \right)^m$$

$$= \frac{e^{i\frac{mx}{2}}}{(2\pi)^{m/2}} \left( \frac{\sin\frac{x}{2}}{\frac{x}{2}} \right)^m.$$

In most cases, we will just use the cubic B-Spine. We choose

$$\varphi(k) = B_4(k + 2), \tag{11}$$

and $\varphi_j(k) = \varphi(k/\Delta k - j)$. Therefore,

$$\varphi^\vee(x) = \frac{1}{(2\pi)} \left( \frac{\sin\frac{x}{2}}{\frac{x}{2}} \right)^4 = O(\frac{1}{|x|^4}) \text{ as x} \to \infty,$$

which defines the *frequency selection kernel*

$$\varphi_j^\vee(x) = e^{-ij\Delta kx}\varphi^\vee(\Delta kx) = e^{-ijK_0x}\varphi^\vee(K_0x), \Delta k = K_0,$$

and

$$\text{supp } \mathcal{F}[\varphi_j^\vee(x)] \subset [\omega_{j-2}, \omega_{j+2}].$$

From the basic property of B-spline functions, we know that

$$\sum_{j\in\mathbb{Z}} \varphi_j(k) = \sum_{j\in\mathbb{Z}} B_4(\frac{k}{\Delta k} + 2 - j) \equiv 1, \quad \forall k \in \mathbb{R}.$$

6

Thus,
$$\sum_{j=-m-1}^{m+2} \varphi_j(k) = 1, \quad \forall k \in [-m\Delta k, m\Delta k]. \tag{12}$$

In Eqn. (12), we define $\varphi_{-m-1} = B_4(k/\Delta k + m + 3)$, $\varphi_m = B_4(k/\Delta k - m + 2)$ and $\varphi_{m+1} = B_4(k/\Delta k - m + 1)$.

Now, with the POU in (7), we can decompose the target function $f(x)$ in the Fourier space as follows,

$$\widehat{f}(k) = \sum_{j=-m}^{m-1} \varphi_j(k)\widehat{f}(k) \triangleq \sum_{j=-m}^{m-1} \widehat{f}_j(k). \tag{13}$$

Equation (13) will give a corresponding decomposition in $x$-space

$$f(x) = \sum_{j=0}^{m} f_j(x), \tag{14}$$

where
$$f_j(x) = \mathcal{F}^{-1}[\widehat{f}_j](x).$$

The decomposition (14) involves $m + 1$ function $f_j(x)$ whose frequency spectrum is limited to $[\omega_{j-1}, \omega_{j+1}]$, therefore, a simple phase shift could translate its spectrum to $[-K_0, K_0]$, then it could be learned quickly by a DNN with $n_0$ training epoches.

- A parallel phase shift DNN (PhaseDNN) algorithm

We define the first approximation to $f(x)$ by $T_\star(x, \theta^{(n_0)})$, $x \in R$,

$$T_\star(x, \theta^{(n_0)}) = f_{\text{DNN}}(x) \approx f(x), \tag{15}$$

based on given $N$-training data

$$\{x_i, y_i = f(x_i)\}_{i=1}^{N}. \tag{16}$$

Let us consider the residual between the function $f(x)$ and its first DNN approximation $f_{\text{DNN}}$,
$$r(x) = f(x) - f_{\text{DNN}}(x). \tag{17}$$

We should have
$$|\mathcal{F}[r](k)| \ll 1 \text{ for } |k| < K_0, \tag{18}$$
namely,
$$\mathcal{F}[f](k) \approx \mathcal{F}[f_{\text{DNN}}(x)](k), \quad k \in (-K_0, K_0). \tag{19}$$

Our next step is to learn this residual function $r(x)$ based on its training data
$$\{x_i, r_i = r(x_i)\}_{i=1}^N. \tag{20}$$

In order to learn $r(x)$ for its frequency information for $|k| > K_0$, we will apply the decomposition (14) to $r(x)$
$$r(x) = \sum_{j=0}^m r_j(x), \tag{21}$$
where
$$r_j(x) = \mathcal{F}^{-1}[\widehat{r_j}(k)] = \mathcal{F}^{-1}[\varphi_j(k)\widehat{r}(k)](x). \tag{22}$$
The training data for $\{r_j(x_i)\}_{i=1}^N$ will be computed by (22), which can be implemented in the data $x$-space through the following convolution

$$r_j(x_i) = \varphi_j^\vee * r(x_i) = \int_{-\infty}^\infty \varphi_j^\vee(x_i - s)r(s)ds$$
$$\approx \frac{2\delta}{N_s} \sum_{x_s \in (x_i - \delta, x_i + \delta)} \varphi_j^\vee(x_i - x_s)r(x_s), \tag{23}$$

where $\delta$ is chosen such that the kernel function $|\varphi_j^\vee(\delta)| = |\varphi^\vee(K_0\delta)|$ is smaller than some given error tolerance, $N_s$ is the number of data points inside $(x_i - \delta, x_i + \delta)$.

Now as the support of $\widehat{r_j}(k)$ is $[\omega_{j-1}, \omega_{j+1}]$, then $\widehat{r_j}(k - \omega_j)$ will have its frequency spectrum in $[-\Delta k, \Delta k] = [-K_0, K_0]$, therefore its inverse Fourier transform $\mathcal{F}^{-1}[\widehat{r_j}(k - \omega_j)]$, denoted as
$$r_j^{\text{shift}}(x) = \mathcal{F}^{-1}[\widehat{r_j}(k - \omega_j)](x), \tag{24}$$
can be learned quickly by a DNN
$$T_j(x, \theta^{(n_0)}) \approx r_j^{\text{shift}}(x), 0 \leq j \leq m, \tag{25}$$

8

in $n_0$-epoches of training. Moreover, we know that

$$r_j^{\text{shift}}(x_i) = e^{i\omega_j x_i} r_j(x_i), 1 \le i \le N, \tag{26}$$

which provides the training data for $r_j^{\text{shift}}(x)$. Equation (26) shows that once $r_j^{\text{shift}}(x)$ is learned, $r_j(x)$ is also learned by removing the phase factor.

$$r_j(x) \approx e^{-i\omega_j x} T_j(x, \theta^{(n_0)}). \tag{27}$$

Now all $r_j(x)$ $0 \le j \le m$ have been learned, we will have an approximation to $r(x)$ as

$$r(x) \approx \sum_{j=0}^{m} e^{-i\omega_j x} T_j(x, \theta^{(n_0)}). \tag{28}$$

Finally, we will obtain our Phase shift DNN (PhaseDNN) approximation of $f(x)$, as illustrated in the flowchart in Fig. 1 over all frequency range $[-M, M]$ as follows

$$f_{\text{DNN}}(x) \leftarrow f_{\text{DNN}}(x) + \sum_{j=0}^{m} e^{-i\omega_j x} T_j(x, \theta^{(n_0)}). \tag{29}$$

**Remark 1.** (*Parallelism of PhaseDNN*) *It is clear the learning for each part of residual function $r_j(x)$ can be done in parallel, the PhaseDNN is a version of domain decomposition method in k-space.*

**Remark 2.** (*Recursive version of PhaseDNN*) *After the updating step (29), in principle, we can repeat the process by starting at (17) with the new $f_{DNN}(x)$ to further reduce any remaining error in all frequency range (see Fig. 1).*

**Remark 3.** (*Frequency Adaptivity of PhaseDNN*) *The frequency-selected residual $r_j(x_i)$ requires the evaluation of the convolution (23), which implicitly assumes sufficient data information is available in the neighborhood of $x_i$. In case such an assumption is not valid, we basically do not have enough data to extract the frequency information for the target function. In all practical sense of learning, we should assume the target function do not have the information of this specific frequency, and there is no need to learn this function $r_j(x)$ near $x_i$ and we could simply set $r_j(x) = 0$ for $x \in (x_i - \delta, x_i + \delta)$, which*
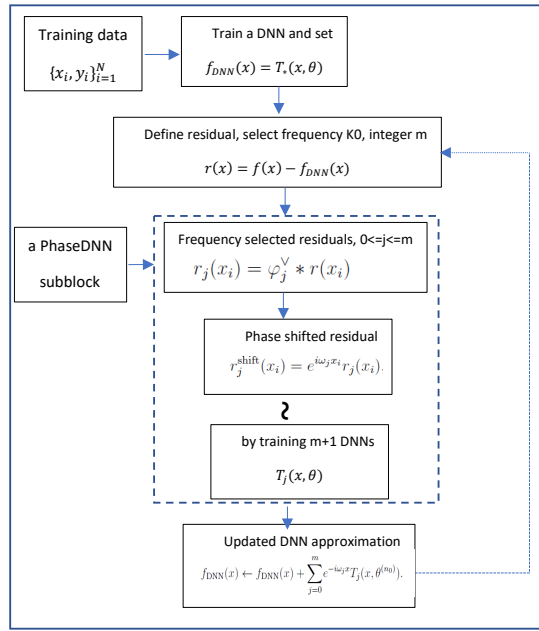
9

Training data
$\{x_i, y_i\}_{i=1}^{N}$

Train a DNN and set
$f_{DNN}(x) = T_*(x, \theta)$

Define residual, select frequency K0, integer m
$r(x) = f(x) - f_{DNN}(x)$

a PhaseDNN
subblock

Frequency selected residuals, 0<=j<=m
$r_j(x_i) = \varphi_j^{\vee} * r(x_i)$

Phase shifted residual
$r_j^{\text{shift}}(x_i) = e^{i\omega_j x_i} r_j(x_i).$

by training m+1 DNNs
$T_j(x, \theta)$

Updated DNN approximation
$f_{\text{DNN}}(x) \leftarrow f_{\text{DNN}}(x) + \sum_{j=0}^{m} e^{-i\omega_j x} T_j(x, \theta^{(n_0)}).$

Figure 1: Flowchart of Parallel PhaseDNN
with shifts for $\boldsymbol{k} \in R^d$ by an alternating coordinate direction sweep.

*is what could happen for a constant (or low) frequency function $r(x)$. As this process is done locally in space and frequency, the PhaseDNN will be able to achieve frequency adaptivity commonly associated with wavelet approximation of signals and images.*

*• **Adaptivity in frequency and local radial Basis fitting** A pre-analysis of the training data could be done to investigate the frequency content of the data, so only those frequencies will be considered when phase shifts are used for error reduction. Firstly, we group the training data $x_i$ into $N_c$ clusters (disjoint or overlapping). Secondly, in a sphere centered at the center of a cluster, a local radial basis function fitting of training data $\{x_i, y_i\}$ falling in the sphere will be carried out [6]. Finally, we apply an 1-D discrete Fourier transform over the fitting function along each coordinate direction within the cluster to get the frequency content for the data in the cluster. Once this analysis is done for each cluster, we will have a collection of frequencies and those with significant magnitude will be used as candidates for the phase shifts in the PhaseDNN.*

**Remark 4.** *(**Selection of $K_0$ and $n_0$** and pre-processing training data) In principle, any nonzero frequency threshold $K_0$ can be used in the implementation of PhaseDNN, which will dictate the size of $n_0$ to achieve sufficient convergence of the underlying DNN over the frequency range $[-K_0, K_0]$. Therefore, a careful analysis and numerical experiment of the specific low frequency convergence of the DNN will be needed.*

*As $n_0$ is selected to train the initial network in PhaseDNN to learn low frequent content of the data, a low-pass filtering on the training data could be used to select appropriate size of $n_0$ needed for achieve convergence within $K_0$ frequency.*

**Remark 5.** *(**Avoiding curse of dimensionality**) It is clear the proposed phase shift DNN can be extended to learn function $y = f(\overrightarrow{x}), \overrightarrow{x} \in R^d$ by making frequency shift along various $\overrightarrow{k}$ unit vector directions. However, for high dimension problems, the shifts employed in the frequency $\boldsymbol{k} \in R^d$-space to learn $r(x)$ in (17) should be done in an alternating direction sweep over the $\boldsymbol{k}$-coordinate directions to avoid the curse of dimensionality. Namely, we will carry out $m+1$ shifts along $k_1$ direction, learn $m+1$ DNNs for the phase shifts only along $k_1$ direction, which will be used to update $f_{DNN}(x)$ in (29). Then, we learn another $m+1$ DNNs for the shifts along $k_2$ direction, etc, until all directions in $\boldsymbol{k}$ are touched. In all, $d(m+1)$ DNNs will have been trained*

*for one sweep of all **k** coordinate directions at a cost $O(d)$. This alternating direction sweep in **k** space can be repeated to improve the accuracy.*

*Another way to improve the efficiency of the **k** alternating direction sweep is to introduce matching random rotation of the coordinate systems in both **x** and **k** spaces, so the rotated **k** coordinate directions will be able to detect the high frequency in non-uniform data during the sweep more effectively.*

## 4. Numerical Result

In this section, we will present some preliminary numerical results to demonstrate the capability of PhaseDNN to learn high frequency content of target functions. In practice, we could sweep over all frequency ranges. For the test function we have some rough idea of the the range of frequencies in the data, only a few frequency intervals are selected for the phase shift.

We choose a target function $f(x)$ defined in $[-\pi, \pi]$ by

$$f(x) = \begin{cases} 10(\sin x + \sin 3x), & \text{if } x \in [-\pi, 0] \\ 10(\sin 23x + \sin 137x + \sin 203x), & \text{if } x \in [0, \pi]. \end{cases} \tag{30}$$

Because the frequencies of this function is well separated, we need not to sweep all the frequencies in $[-\infty, +\infty]$. Instead, we choose $\Delta k = 5$, and use the following functions

$$\varphi_1(k) = \chi_{[-205,-200]}(k) \quad \varphi_2(k) = \chi_{[-140,-135]}(k)$$
$$\varphi_3(k) = \chi_{[-25,-20]}(k) \qquad \varphi_4(k) = \chi_{[-5,0]}(k)$$
$$\varphi_5(k) = \chi_{[0,5]}(k) \qquad \varphi_6(k) = \chi_{[20,25]}(k)$$
$$\varphi_7(k) = \chi_{[135,140]}(k) \qquad \varphi_8(k) = \chi_{[200,205]}(k)$$

to collect the frequency information in the corresponding frequency intervals and shift the center of the interval to the origin by a phase factor. For each $f_j(x) = \mathcal{F}^{-1}[\widehat{f}\varphi_j](x)$, we construct two DNNs to learn its real part and imaginary part, separately. Every DNN have 4 hidden layers and each layer has 40 neurons. Namely, the DNN has the structure 1-40-40-40-40-1. The training data is obtained by 10,000 samples from the uniform distribution on $[-\pi, \pi]$ and the testing data is 500 uniformly distributed samples. We train these DNNs by Adam optimizer with training rate 0.0002 with 10 epochs of training for each DNN. The result is shown in Fig. 2.

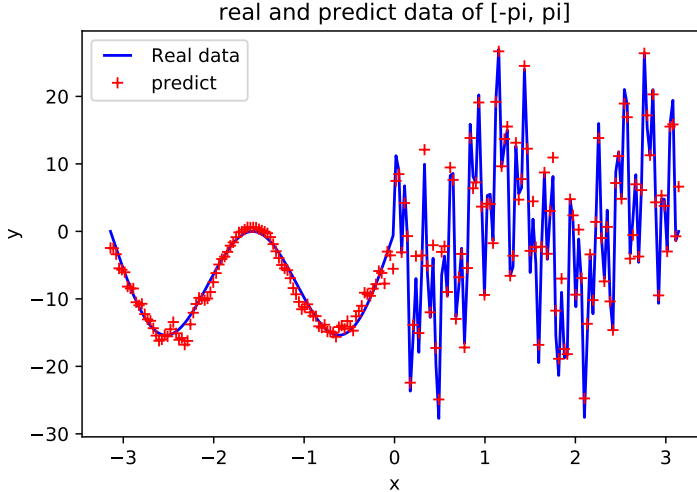The detail of this training is shown in Fig. 3.

12

Figure 2: The testing result of $f_{DNN}(x)$ trained by phase DNN. The blue solid line is $f(x)$ and the data marked by $+$ in red is the value of $f_{DNN}(x)$ at testing data set.

These figures clearly shows that phase DNN can capture the various high frequencies, from low frequency $\pm 1$, $\pm 3$ to high frequency $\pm 203$ quite well. The training error and time of phase DNN are collected in table 1

It is shown that taking the advantage of vectorization, the convolution step costs about 20% time of training. Because in different interval, $f_j(x)$ can be trained in parallel, PhaseDNN is ideal to take advantage of parallel computing architectures.

In comparison with one single DNN, we have used a 24 hidden layers with 80 neurons per hidden layer with the same amount of training data, the loss is still around 100 after more than 5000 epoch of training taking over 22 hours nonstop running on the same workstation.

## 5. Discussion and Conclusion

In this paper, we proposed a parallel PhaseDNN, consisting of a series of commonly-used DNNs in parallel, to achieve uniform frequency convergence of deep neural network in approximating a function. In contrast to the wavelet multi-resolution approximation with its higher resolution wavelet subspaces generated by dilation of mother wavelet to provide approximation to higher frequency components, the PhaseDNN uses a phase shift strategy
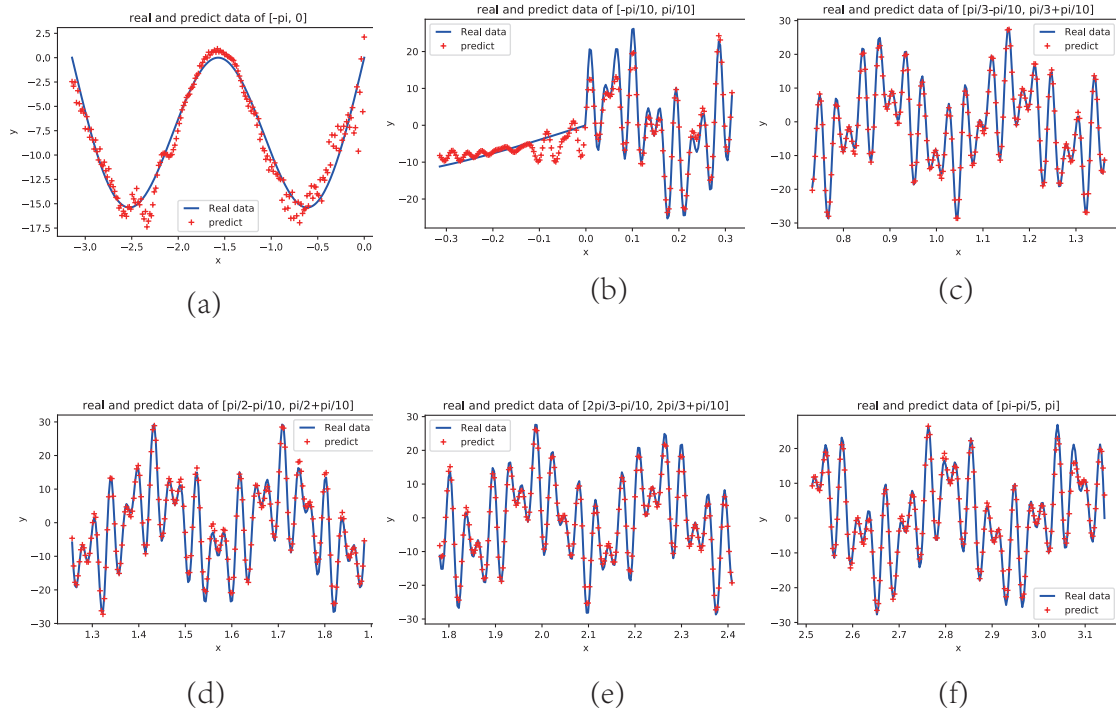
13

Figure 3: The detail result of training in different intervals. The subfigure (a)-(f) shows the results in interval $[-\pi, 0]$, $[-\pi/10, \pi/10]$, $[\pi/3-\pi/10, \pi/3+\pi/10]$, $[\pi/2-\pi/10, \pi/2+\pi/10]$, $[2\pi/3 - \pi/10, 2\pi/3 + \pi/10]$ and $[\pi - \pi/10, \pi]$ correspondingly. The blue solid line is $f(x)$ and the data marked by $+$ in red is the value of $f_{DNN}(x)$ at testing data set.

14

|  | Convolution time(s) | Training Time(s) | Training Error | Test Error |
|---|---|---|---|---|
| $[-205, -200]$ | 34.43 | 169.82 | - | - |
| $[-140, -135]$ | 33.74 | 185.40 | - | - |
| $[-25, -20]$ | 33.22 | 199.48 | - | - |
| $[-5, 0]$ | 9.70 | 214.39 | - | - |
| $[0, 5]$ | 9.78 | 230.46 | - | - |
| $[20, 25]$ | 33.21 | 247.34 | - | - |
| $[135, 140]$ | 33.38 | 263.89 | - | - |
| $[200, 205]$ | 32.68 | 279.39 | - | - |
| Total | 220.14 | 1790.17 | 0.01508 | 0.06831 |

Table 1: The training time and error statistics. For each $j$, the training time is the sum of training time of real and imaginary part. For each DNN, epoch is chosen to be 10.

to handle higher frequency, which allows us to make use of the fast convergence of common DNN in the low frequency range. The PhaseDNN employs original training data, which are easily modified (i.e., a multiplication by a phase factor) to carry out the required phase shift operation.

Two additional comments on the PhaseDNN are given below.

- **Training data and PhaseDNN**. The convolution computation (23) required in obtaining the training data for the selected frequency component of the residual $r_j(x)$ implicitly implies that the training data location $x$ are dense enough to extract the information of the given frequency range. Similarly, The frequency shift in the PhaseDNN in (26) also implicitly requires that the training data location $x$ are dense enough to be uni-solvent to represent the Fouirer mode $e^{\pm i\omega_j x}$. Otherwise due to the alias phenomenon of the Fourier modes on a discrete set, the simple multiplication by $e^{\pm i\omega_j x}$ on the training data will not achieve the desired frequency shift.

- **Spatial and Frequency adaptive approximation**. The frequency shifts can be done in any order, therefore, for problems for prior knowledge of frequency information or need of knowledge for specific frequency such as segmentation, the PhaseDNN can be used to obtain convergence in the desired frequency ranges. As mentioned in Remark

3, each DNN in the PhaseDNN is linked to the local spatial and frequency, therefore PhaseDNN can achieve also spatial adaptivity in the construction of the its DNNs.

In future work, more numerical tests will be carried out to understand the potential and behavior of the proposed PhaseDNN and calibrate the efficiency and accuracy of the PhaseDNN over a stand-alone DNN, especially for high dimensional problems.

**Acknowledgment**

**References**

[1] Xu, Zhi-Qin John, Understanding training and generalization in deep learning by Fourier analysis, arXiv:1808.04295, November, 2018.

[2] Xu, Zhi-Qin John, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency Principle: Fourier Analysis Sheds Light on Deep Neural Networks. arXiv preprint arXiv:1901.06523 (2019).

[3] Brandt, Achi. Multi-level adaptive solutions to boundary-value problems. Mathematics of computation 31.138 (1977): 333-390.

[4] Daubechies, Ingrid. Ten lectures on wavelets. Vol. 61. Siam, 1992.

[5] Cai, Wei, and Jianzhong Wang. Adaptive multiresolution collocation methods for initial-boundary value problems of nonlinear PDEs. SIAM Journal on Numerical Analysis 33, no. 3 (1996): 937-970.

[6] Wendland, Holger. Scattered data approximation. Vol. 17. Cambridge university press, 2004.