# Parareal Neural Networks Emulating a Parallel-in-time Algorithm[*]

Chang-Ock Lee[†], Youngkyu Lee[†], and Jongho Park [‡]

**Abstract.** As deep neural networks (DNNs) become deeper, the training time increases. In this perspective, multi-GPU parallel computing has become a key tool in accelerating the training of DNNs. In this paper, we introduce a novel methodology to construct a parallel neural network that can utilize multiple GPUs simultaneously from a given DNN. We observe that layers of DNN can be interpreted as the time steps of a time-dependent problem and can be parallelized by emulating a parallel-in-time algorithm called parareal. The parareal algorithm consists of fine structures which can be implemented in parallel and a coarse structure which gives suitable approximations to the fine structures. By emulating it, the layers of DNN are torn to form a parallel structure, which is connected using a suitable coarse network. We report accelerated and accuracy-preserved results of the proposed methodology applied to VGG-16 and ResNet-1001 on several datasets.

**Key words.** deep neural network, parallel computing, time-dependent problem, parareal

**AMS subject classifications.** 68T01, 68U10, 68W10

**1. Introduction.** Deep neural networks (DNNs) have demonstrated success for many classification and regression tasks such as image recognition [15, 22] and natural language processing [5, 19]. A principal reason for why DNN performs well is the depth of DNN, i.e., the number of sequential layers of DNN. Each layer of DNN is composed of a combination of an affine transformation and a nonlinear activation function, e.g., a rectified linear unit (ReLU). A broad range of functions can be generated by stacking a number of layers with nonlinear activation functions so that DNN can be a model that fits the given data well [6, 17]. However, there are undesirable side effects of using many layers for DNN. Due to the large number of layers in DNNs, DNN training is time-consuming and there are demands to reduce training time these days. Recently, multi-GPU parallel computing has become an important topic for accelerating DNN training [2, 3, 13].

Data parallelism [2] is a commonly used parallelization technique. In data parallelism, the training dataset is distributed across multiple GPUs and then processed separately. For instance, suppose that we have 2 GPUs and want to apply data parallelism to the mini-batch gradient descent with the batch size 128. In this case, each GPU computes 64-batch and the computed gradients are averaged. In data parallelism, each GPU must possess a whole copy of the DNN model so that inter-GPU communication is required at every step of the training process in order to update all parameters of the model. Therefore, the training time is seriously deteriorated when the number of layers in the model is large. In order to resolve such a drawback, several asynchronous methodologies of data parallelism were proposed [4, 25, 38].

[†]Department of Mathematical Sciences, KAIST, Daejeon 34141, Korea (colee@kaist.edu, lyk92@kaist.ac.kr).
 [‡]Natural Science Research Institute, KAIST, Daejeon 34141, Korea (jongho.park@kaist.ac.kr).

In asynchronous data parallelism, a parameter server in charge of parameter update is used; it collects computed gradients from other GPUs in order to update parameters, and then distributes the updated parameters to other GPUs.

On the other hand, model parallelism [18, 30] is usually utilized when the capacity of a DNN exceeds the available memory of a single GPU. In model parallelism, layers of DNN and their corresponding parameters are partitioned into multiple GPUs. Since each GPU owns part of the model's parameters, the cost of inter-GPU communication in model parallelism is much less than the cost of data parallelism. However, only one GPU is active at a time in the naive application of model parallelism. To resolve the inefficiency, a pipelining technique called PipeDream [30] which uses multiple mini-batches concurrently was proposed. PipeDream has a consistency issue in parameter update that a mini-batch may start the training process before its prior mini-batch updates parameters. To avoid this issue, another pipelining technique called Gpipe [18] was proposed; it divides each mini-bath into micro-batches and utilizes micro-batches for the simultaneous update of parameters. However, experiments [3] have shown that the possible efficiency of Gpipe can not exceed 29% of that of Pipedream. Recently, further improvements of PipeDream and Gpipe were considered; see SpecTrain [3] and PipeMare [35].

There are several notable approaches of parallelism based on layerwise decomposition of the model [9, 13]. Unlike the aforementioned ones, these approaches modify data propagation in the training process of the model. Günther et al. [13] replaced the sequential data propagation of layers in DNN by a nonlinear in-time multigrid method [8]. It showed strong scalability in a simple ResNet [15] when it was implemented on a computer cluster with multiple CPUs. Fok et al. [9] introduced WarpNet which was based on ResNet. They replaced residual units (RUs) in ResNet by the first-order Taylor approximations, which enabled parallel implementation. In WarpNet, $(N-1)$ RUs are replaced by a single warp operator which can be treated in parallel using $N$ GPUs. However, this approach requires data exchange at every warp operation so that it may suffer from a communication bottleneck as the DNN becomes deeper.

In this paper, we propose a novel paradigm of multi-GPU parallel computing for DNNs, called *parareal neural network*. In general, DNN has a feed-forward architecture. That is, the output of DNN is obtained from the input by sequential compositions of functions representing layers. We observe that sequential computations can be interpreted as time steps of a time-dependent problem. In the field of numerical analysis, after a pioneering work of Lions et al. [26], there have been numerous researches on parallel-in-time algorithms to solve time-dependent problems in parallel; see, e.g., [12, 27, 29]. Motivated by these works, we present a methodology to transform a given feed-forward neural network to another neural network called parareal neural network which naturally adopts parallel computing. The parareal neural network consists of fine structures which can be processed in parallel and a coarse structure which approximates the fine structures by emulating one of the parallel-in-time algorithms called parareal [26]. Unlike the existing methods mentioned above, the parareal neural network can significantly reduce the time for inter-GPU communication because the fine structures do not communicate with each other but communicate only with the coarse structure. Therefore, the proposed methodology is effective in reducing the elapsed time for dealing with very deep neural networks. Numerical results confirm that the parareal neural network provides similar or better performance to the original network even with less training time.

The rest of this paper is organized as follows. In Section 2, we briefly summarize the parareal

algorithm for time-dependent differential equations. An abstract framework for the construction of the parareal neural network is introduced in Section 3. In Section 4, we present how to apply the proposed methodology to two popular neural networks VGG-16 [32] and ResNet-1001 [16] with details. Also, accelerated and accuracy-preserved results of parareal neural networks for VGG-16 and ResNet-1001 with datasets CIFAR-10, CIFAR-100 [20], MNIST [23], SVHN [31], and ImageNet [7] are given. We conclude this paper with remarks in Section 5.

**2. The parareal algorithm.** The parareal algorithm proposed by Lions et al. [26] is a parallel-in-time algorithm to solve time-dependent differential equations. For the purpose of description, the following system of ordinary differential equations is considered:

$$\dot{\mathbf{u}}(t) = A\mathbf{u}(t) \ \ \text{in} \ [0,T], \ \ \mathbf{u}(0) = \mathbf{u}_0, \tag{2.1}$$

where $A \colon \mathbb{R}^m \to \mathbb{R}^m$ is an operator, $T > 0$, and $\mathbf{u}_0 \in \mathbb{R}^m$. The time interval $[0, T]$ is decomposed into $N$ subintervals $0 = T_0 < T_1 < \cdots < T_N = T$. First, an approximated solution $\{\mathbf{U}_j^1\}_{j=0}^N$ of (2.1) on the *coarse* grid $\{T_j\}_{j=0}^N$ is obtained by the backward Euler method with the step size $\Delta T_j = T_{j+1} - T_j$:

$$\frac{\mathbf{U}_{j+1}^1 - \mathbf{U}_j^1}{\Delta T_j} = A\mathbf{U}_{j+1}^1, \ \ \mathbf{U}_0^1 = \mathbf{u}_0 \ \ \text{for} \ j = 0, \ldots, N-1.$$

Then in each time subinterval $[T_j, T_{j+1}]$, we construct a local solution $\mathbf{u}_j^1$ by solving the following initial value problem:

$$\dot{\mathbf{u}}_j^1(t) = A\mathbf{u}_j^1(t) \ \ \text{in} \ [T_j, T_{j+1}], \ \ \mathbf{u}_j^1(T_j) = \mathbf{U}_j^1. \tag{2.2}$$

The computed solution $\mathbf{u}_j^1$ does not agree with the exact solution $\mathbf{u}$ in general since $\mathbf{U}_j^1$ differs from $\mathbf{u}(T_j)$. For $k \geq 1$, a better coarse approximation $\{\mathbf{U}_j^{k+1}\}_{j=0}^N$ than $\{\mathbf{U}_j^k\}_{j=0}^N$ is obtained by the *coarse grid correction*: Let $\mathbf{U}_0^{k+1} = \mathbf{U}_0^k, \mathbf{S}_0^k = 0$; for $j = 0, \ldots, N-1$, we repeat the followings:

1. Compute the difference at the coarse node: $\mathbf{S}_{j+1}^k = \mathbf{u}_j^k(T_{j+1}) - \mathbf{U}_{j+1}^k$.
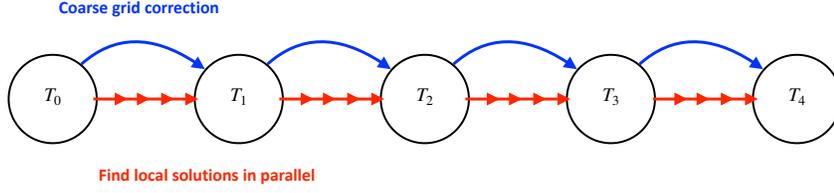2. Propagate the difference to the next coarse node by the backward Euler method:
   $\frac{\boldsymbol{\delta}_{j+1}^k - \boldsymbol{\delta}_j^k}{\Delta T_j} = A\boldsymbol{\delta}_{j+1}^k + \mathbf{S}_j^k, \ \boldsymbol{\delta}_0^k = 0.$
3. Set $\mathbf{U}_{j+1}^{k+1} = \mathbf{U}_{j+1}^k + \boldsymbol{\delta}_{j+1}^k$.

That is, $\{\mathbf{U}_j^{k+1}\}_{j=0}^N$ is made by the correction with the propagated residual $\{\boldsymbol{\delta}_j^k\}_{j=0}^N$. Using the updated coarse approximation $\{\mathbf{U}_j^{k+1}\}_{j=0}^N$, one obtains a new local solution $\mathbf{u}_j^{k+1}$ in the same manner as (2.2):
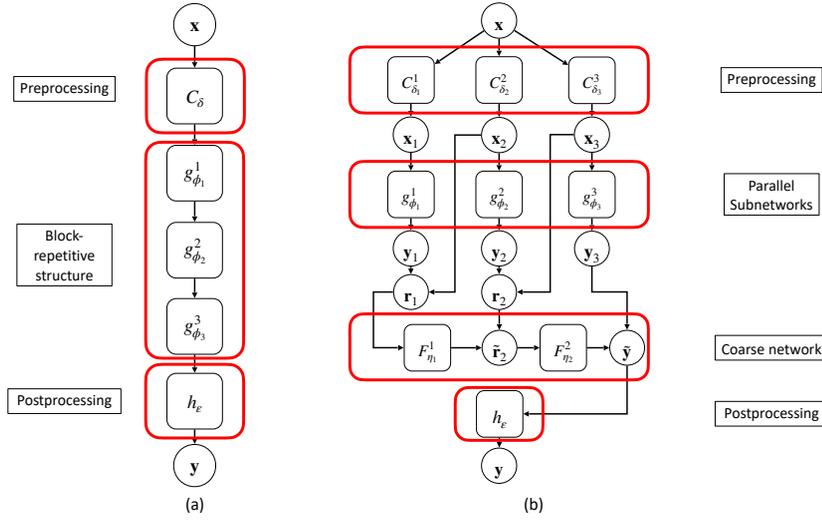
$$\dot{\mathbf{u}}_j^{k+1}(t) = A\mathbf{u}_j^{k+1}(t) \ \ \text{in} \ [T_j, T_{j+1}], \ \ \mathbf{u}_j^{k+1}(T_j) = \mathbf{U}_j^{k+1}. \tag{2.3}$$

It is well-known that $\mathbf{u}_j^k$ converges to the exact solution $\mathbf{u}$ uniformly as $k$ increases [1, 11].

Since (2.3) can be solved independently in each time subinterval, we may assign the problem in each $[T_j, T_{j+1}]$ to the processor one by one and compute $\mathbf{u}_j^{k+1}$ in parallel. In this sense, the parareal algorithm is suitable for parallel computation on distributed memory architecture. A diagram illustrating the parareal algorithm is presented in Figure 1.

**Figure 1.** *Fine and coarse propagations in the parareal algorithm: the red lines which propagate from $T_j$ to $T_{j+1}$ represent* (2.2) *and the blue curves which directly connect $T_j$ and $T_{j+1}$ represent* (2.3).



**Figure 2.** *A feed-forward neural network and its corresponding parareal neural network:* **(a)** *Feed-forward neural network $f_\theta$,* **(b)** *Parareal neural network $\bar{f}_{\bar{\theta}}$ with $N$ parallel subnetworks ($N = 3$).*

**3. Parareal neural networks.** In this section, we propose a methodology to design a *parareal neural network* by emulating the parareal algorithm introduced in Section 2 from a given feed-forward neural network. The resulting parareal neural network has an intrinsic parallel structure and is suitable for parallel computation using multiple GPUs with distributed memory simultaneously.

**3.1. Parallelized forward propagation.** Let $f_\theta \colon X \to Y$ be a feed-forward neural network, where $X$ and $Y$ are the spaces of inputs and outputs, respectively, and $\theta$ is a vector consisting of parameters. Since many modern neural networks such as [16, 32, 36] have block-repetitive substructures, we may assume that $f_\theta$ can be written as the composition of three functions $C_\delta \colon X \to W_0$, $g_\phi \colon W_0 \to W_1$, and $h_\varepsilon \colon W_1 \to Y$, i.e.,

$$f_\theta = h_\varepsilon \circ g_\phi \circ C_\delta, \quad \theta = \delta \oplus \phi \oplus \varepsilon,$$

where $W_0$ and $W_1$ are vector spaces, $g_\phi$ is a block-repetitive substructure of $f_\theta$ with parameters $\phi$, $C_\delta$ is a *preprocessing operator* with parameters $\delta$, and $h_\varepsilon$ is a *postprocessing operator* with parameters $\varepsilon$. Note that $\oplus$ represents a concatenation. Examples of VGG-16 [32] and ResNet-1001 [16] will be given in Section 4.

For appropriate vector spaces $X_0, X_1, \ldots, X_N$, we further assume that $g_\phi$ can be partitioned into $N$ subnetworks $\{g_{\phi_j}^j \colon X_{j-1} \to X_j\}_{j=1}^N$ which satisfy the followings:

- $X_0 = W_0$ and $X_N = W_1$,
- $\phi = \bigoplus_{j=1}^N \phi_j$,
- $g_\phi = g_{\phi_N}^N \circ g_{\phi_{N-1}}^{N-1} \circ \cdots \circ g_{\phi_1}^1$.

See Figure 2(a) for a graphical description for the case $N = 3$. In the computation of $g_\phi$, the subnetworks $\{g_{\phi_j}^j\}_{j=1}^N$ are computed in the sequential manner. Regarding the subnetworks as subintervals of a time-dependent problem and adopting the idea of the parareal algorithm introduced in Section 2, we construct a new neural network $\bar{f}_{\bar\theta} \colon X \to Y$ which contains $\{g_{\phi_j}^j\}_{j=1}^N$ as parallel subnetworks; the precise definition for parameters $\bar\theta$ will be given in (3.3).

Since the dimensions of the spaces $\{X_j\}_{j=0}^{N-1}$ are different for each $j$ in general, we introduce preprocessing operators $C_{\delta_j}^j \colon X \to X_{j-1}$ such that $C_{\delta_1}^1 = C_\delta$ and $C_{\delta_j}^j$ for $j = 2, \ldots, N$ play similar roles to $C_\delta$; particular examples will be given in Section 4. We write $\mathbf{x}_j \in X_{j-1}$ and $\mathbf{y}_j \in X_j$ as follows:

$$(3.1) \qquad \mathbf{x}_j = C_{\delta_j}^j(\mathbf{x}) \text{ for } \mathbf{x} \in X, \quad \mathbf{y}_j = g_{\phi_j}^j(\mathbf{x}_j).$$

Then, we consider neural networks $F_{\eta_j}^j \colon X_j \to X_{j+1}$ with parameters $\eta_j$ for $j \geq 1$ such that it approximates $g_{\phi_{j+1}}^{j+1}$ well while it has a cheaper computational cost than $g_{\phi_{j+1}}^{j+1}$, i.e., $F_{\eta_j}^j \approx g_{\phi_{j+1}}^{j+1}$ and $\dim(\eta_j) \ll \dim(\phi_{j+1})$. Emulating the coarse grid correction of the parareal algorithm, we assemble a network called *coarse network* with building blocks $F_{\eta_j}^j$. With inputs $\mathbf{x}_{j+1}, \mathbf{y}_j$, and an output $\mathbf{y} \in Y$, the coarse network is described as follows:

$$(3.2a) \qquad \mathbf{r}_N = \mathbf{0}, \quad \mathbf{r}_j = \mathbf{y}_j - \mathbf{x}_{j+1} \quad \text{for } j = 1, \ldots, N-1,$$

$$(3.2b) \qquad \tilde{\mathbf{r}}_1 = \mathbf{r}_1, \quad \tilde{\mathbf{r}}_{j+1} = \mathbf{r}_{j+1} + F_{\eta_j}^j(\tilde{\mathbf{r}}_j) \quad \text{for } j = 1, \ldots, N-1,$$

$$(3.2c) \qquad \tilde{\mathbf{y}} = \mathbf{y}_N + \tilde{\mathbf{r}}_N.$$

That is, in the coarse network, the residual $\mathbf{r}_j$ at the interface between layers $g_{\phi_j}^j$ and $g_{\phi_{j+1}}^{j+1}$ propagates through shallow neural networks $F_{\eta_1}^1, \ldots, F_{\eta_{N-1}}^{N-1}$. Then the propagated residual is added to the output.

Finally, the parareal neural network $\bar{f}_{\bar\theta}$ corresponding to the original network $f_\theta$ is defined as

$$(3.3) \qquad \bar{f}_{\bar\theta}(\mathbf{x}) = h_\varepsilon(\tilde{\mathbf{y}}), \quad \bar\theta = \left( \bigoplus_{j=1}^N (\delta_j \oplus \phi_j) \right) \oplus \left( \bigoplus_{j=1}^{N-1} \eta_j \right) \oplus \varepsilon.$$

That is, $\bar{f}_{\bar\theta}$ is composed of the preprocessing operators $\{C_{\delta_j}^j\}$, parallel subnetworks $\{g_{\phi_j}^j\}$, the coarse network $\{F_{\eta_j}^j\}$, and the postprocessing operator $h_\varepsilon$. Figure 2(b) illustrates $\bar{f}_{\bar\theta}$. Since

---

**Algorithm 3.1** Forward propagation of the parareal neural network $\bar{f}_{\bar{\theta}}$

---

Broadcast $\mathbf{x}$ to all processors.
**for** $j = 1, 2, \ldots, N$ **in parallel do**

$$\mathbf{x}_j = C_{\delta_j}^j(\mathbf{x}), \ \mathbf{y}_j = g_{\phi_j}^j(\mathbf{x}_j)$$

**end for**
Gather $\mathbf{x}_j, \mathbf{y}_j$ from all processors.
**for** $j = 1, 2, \ldots, N-1$ **do**

$$\mathbf{r}_j = \mathbf{y}_j - \mathbf{x}_{j+1}$$

**end for**
$\mathbf{r}_N = 0, \ \tilde{\mathbf{r}}_1 = \mathbf{r}_1.$
**for** $j = 1, 2, \ldots, N-1$ **do**

$$\tilde{\mathbf{r}}_{j+1} = \mathbf{r}_{j+1} + F_{\eta_j}^j(\tilde{\mathbf{r}}_j)$$

**end for**
$\tilde{\mathbf{y}} = \mathbf{y}_N + \tilde{\mathbf{r}}_N$
$\mathbf{y} = h_\varepsilon(\tilde{\mathbf{y}})$

---

**Algorithm 3.2** Gradient computation for the parareal neural network $\bar{f}_{\bar{\theta}}$

---

$\dfrac{\partial \bar{f}_{\bar{\theta}}}{\partial \varepsilon} = \dfrac{\partial h_\varepsilon}{\partial \varepsilon}, \ \mathcal{D}_N = \dfrac{\partial h_\varepsilon}{\partial \tilde{\mathbf{y}}}, \ \mathcal{D}_0 = 0$
**for** $j = N-1, N-2, \ldots 1$ **do**

$$\frac{\partial \bar{f}_{\bar{\theta}}}{\partial \eta_j} = \mathcal{D}_{j+1} \cdot \frac{\partial F_{\eta_j}^j}{\partial \eta_j}, \ \mathcal{D}_j = \mathcal{D}_{j+1} \cdot \frac{\partial F_{\eta_j}^j}{\partial \tilde{\mathbf{r}}_j}$$

**end for**
Send $\mathcal{D}_{j-1}$ and $\mathcal{D}_j$ to the $j$-th processor.
**for** $j = 1, 2, \ldots, N$ **in parallel do**

$$\frac{\partial \bar{f}_{\bar{\theta}}}{\partial \phi_j} = \mathcal{D}_j \cdot \frac{\partial g_{\phi_j}^j}{\partial \phi_j}, \ \frac{\partial \bar{f}_{\bar{\theta}}}{\partial \delta_j} = \left( \mathcal{D}_j \cdot \frac{\partial g_{\phi_j}^j}{\partial \mathbf{x}_j} - \mathcal{D}_{j-1} \right) \frac{\partial C_{\delta_j}^j}{\partial \delta_j}$$

**end for**

---

each $g_{\phi_j}^j \circ C_{\delta_j}^j$ lies in parallel, all computations related to $g_{\phi_j}^j \circ C_{\delta_j}^j$ can be done independently; parallel structures of forward and backward propagations for $\bar{f}_{\bar{\theta}}$ are described in Algorithm 3.1 and Algorithm 3.2, respectively; detailed derivation of Algorithm 3.2 will be provided in Subsection 3.2. Therefore, multiple GPUs can be utilized to process $\{g_{\phi_j}^j \circ C_{\delta_j}^j\}$ simultaneously for each $j$. In this case, one may expect significant decrease of the elapsed time for training

$\bar{f}_{\bar{\theta}}$ compared to the original network $f_\theta$. On the other hand, the coarse network cannot be parallelized since $\{F_{\eta_j}^j\}$ is computed in the sequential manner. One should choose $F_{\eta_j}^j$ whose computational cost is as cheap as possible in order to reduce the bottleneck effect of the coarse network.

In the following proposition, we show that the proposed parareal neural network $\bar{f}_{\bar{\theta}}$ is constructed consistently in the sense that it recovers the original neural network $f_\theta$ under a simplified setting.

**Proposition 3.1 (Consistency).** *Assume that the original network $f_\theta$ is linear and $F_{\eta_j}^j = g_{\phi_{j+1}}^{j+1}$ for $j = 1, \dots, N-1$. Then we have $\bar{f}_{\bar{\theta}}(\mathbf{x}) = f_\theta(\mathbf{x})$ for all $\mathbf{x} \in X$.*

*Proof.* We define a function $P_j \colon X \to X_{j+1}$ inductively as follows:

$$(3.4) \quad P_0(\mathbf{x}) = 0, \quad P_j(\mathbf{x}) = F_{\eta_j}^j \left( (g_{\phi_j}^j \circ C_{\delta_j}^j)(\mathbf{x}) - C_{\delta_{j+1}}^{j+1}(\mathbf{x}) + P_{j-1}(\mathbf{x}) \right), \quad 1 \le j \le N-1.$$

Then it follows that

$$(3.5) \qquad \bar{f}_{\bar{\theta}}(\mathbf{x}) = h_\varepsilon \left( (g_{\phi_N}^N \circ C_{\delta_N}^N)(\mathbf{x}) + P_{N-1}(\mathbf{x}) \right).$$

First, we show by mathematical induction that

$$(3.6) \qquad P_j(\mathbf{x}) = (g_{\phi_{j+1}}^{j+1} \circ g_{\phi_j}^j \circ \dots g_{\phi_1}^1 \circ C_{\delta_1}^1)(\mathbf{x}) - (g_{\phi_{j+1}}^{j+1} \circ C_{\delta_{j+1}}^{j+1})(\mathbf{x}), \quad 1 \le j \le N-1.$$

The case $j = 1$ is straightforward from (3.4). Suppose that (3.6) holds for $j = m-1$. Since the original network $f_\theta$ is linear and $F_{\eta_j}^j = g_{\phi_{j+1}}^{j+1}$, we get

$$\begin{aligned}
P_m(\mathbf{x}) &= F_{\eta_m}^m \left( (g_{\phi_m}^m \circ C_{\delta_m}^m)(\mathbf{x}) - C_{\delta_{m+1}}^{m+1}(\mathbf{x}) + P_{m-1}(\mathbf{x}) \right) \\
&= (g_{\phi_{m+1}}^{m+1} \circ g_{\phi_m}^m \circ C_{\delta_m}^m)(\mathbf{x}) - (g_{\phi_{m+1}}^{m+1} \circ C_{\delta_{m+1}}^{m+1})(\mathbf{x}) + (g_{\phi_{m+1}}^{m+1} \circ P_{m-1})(\mathbf{x}) \\
&= (g_{\phi_{m+1}}^{m+1} \circ g_{\phi_m}^m \circ \dots g_{\phi_1}^1 \circ C_{\delta_1}^1)(\mathbf{x}) - (g_{\phi_{m+1}}^{m+1} \circ C_{\delta_{m+1}}^{m+1})(\mathbf{x}),
\end{aligned}$$

where the last equality is due to the induction hypothesis. Hence, (3.6) also holds for $j = m$, which implies that it is true for all $j$.

Combining (3.5) and (3.6), we obtain

$$\begin{aligned}
\bar{f}_{\bar{\theta}}(\mathbf{x}) &= h_\varepsilon \left( (g_{\phi_N}^N \circ C_{\delta_N}^N)(\mathbf{x}) + P_{N-1}(\mathbf{x}) \right) \\
&= (h_\varepsilon \circ g_{\phi_N}^N \circ g_{\phi_{N-1}}^{N-1} \circ \dots g_{\phi_1}^1 \circ C_{\delta_1}^1)(\mathbf{x}) \\
&= (h_\varepsilon \circ g_{\phi_N}^N \circ g_{\phi_{N-1}}^{N-1} \circ \dots g_{\phi_1}^1 \circ C_\delta)(\mathbf{x}) \\
&= f_\theta(\mathbf{x}),
\end{aligned}$$

which completes the proof. ∎

Proposition 3.1 presents a guideline on how to design the coarse network of $\bar{f}_{\bar{\theta}}$. Under the assumption that $f_\theta$ is linear, a sufficient condition to ensure that $\bar{f}_{\bar{\theta}} = f_\theta$ is $F_{\eta_j}^j = g_{\phi_{j+1}}^{j+1}$ for all $j$. Therefore, we can say that it is essential to design the coarse network with $F_{\eta_j}^j \approx g_{\phi_{j+1}}^{j+1}$ to

ensure that the performance of $\bar{f}_{\bar{\theta}}$ is as good as that of $f_{\theta}$. Detailed examples will be given in Section 4.

On the other hands, the propagation of the coarse network in the parareal neural network is similar to gradient boosting [10, 28], one of the ensemble techniques. In (3.2), the coarse network satisfies $F_{\eta_j}^j(\tilde{\mathbf{r}}_j) = \tilde{\mathbf{r}}_{j+1} - \mathbf{r}_{j+1}$. Let $\mathcal{L}(\mathbf{x}, \mathbf{y}) = \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|^2$, then we have

$$F_{\eta_j}^j(\tilde{\mathbf{r}}_j) = \tilde{\mathbf{r}}_{j+1} - \mathbf{r}_{j+1} = -\frac{\partial \mathcal{L}(\tilde{\mathbf{r}}_{j+1}, \mathbf{r}_{j+1})}{\partial \mathbf{r}_{j+1}}.$$

From the viewpoint of gradient boosting, the propagation of coarse network can be expressed by the following gradient descent method:

$$\tilde{\mathbf{r}}_{j+1} = \mathbf{r}_{j+1} - \frac{\partial \mathcal{L}(\tilde{\mathbf{r}}_{j+1}, \mathbf{r}_{j+1})}{\partial \mathbf{r}_{j+1}}.$$

Thus, it can be understood that the coarse network constructed by emulating the coarse grid correction of the parareal algorithm corrects the residuals at the interface by the gradient descent method, i.e., it reduces the difference between $\mathbf{x}_{j+1}$ and $\mathbf{y}_j$ at each interface.

Furthermore, we can think of the preprocessing $C_{\delta_j}^j$ and subnetwork $g_{\phi_j}^j$ in parareal neural network as a single shallow network $g_{\phi_j}^j \circ C_{\delta_j}^j$. Then, the parareal neural network can be thought of as a network in which several shallow neural networks are stacked, such as the Stacked generalization [34] of the ensemble technique. Thanks to the coarse network, the parareal neural network does not simply stack the shallow networks, but behaves like a deep neural network that sequentially computes the parallel subnetworks as mentioned in Proposition 3.1.

**3.2. Details of backward propagation.** We present a detailed description on the backward propagation for the parareal neural network $\mathbf{y} = \bar{f}_{\bar{\theta}}(\mathbf{x})$. Partial derivatives $\frac{\partial \bar{f}_{\bar{\theta}}}{\partial \varepsilon}$ and $\frac{\partial \mathbf{y}}{\partial \tilde{\mathbf{y}}}$ regarding to the postprocessing operator $h_\varepsilon$ are computed directly from (3.3):

$$\frac{\partial \bar{f}_{\bar{\theta}}}{\partial \varepsilon} = \frac{\partial h_\varepsilon}{\partial \varepsilon}, \quad \frac{\partial \mathbf{y}}{\partial \tilde{\mathbf{y}}} = \frac{\partial h_\varepsilon}{\partial \tilde{\mathbf{y}}}.$$

It is clear from (3.2) that

$$(3.7) \qquad \frac{\partial \mathbf{r}_j}{\partial \mathbf{y}_j} = 1, \quad \frac{\partial \mathbf{r}_j}{\partial \mathbf{x}_{j+1}} = -1, \quad \frac{\partial \tilde{\mathbf{r}}_j}{\partial \mathbf{r}_j} = 1, \quad \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{y}_N} = 1, \quad \frac{\partial \tilde{\mathbf{y}}}{\partial \tilde{\mathbf{r}}_N} = 1.$$

Moreover, by (3.2b), we get

$$(3.8) \qquad \frac{\partial \tilde{\mathbf{r}}_{j+1}}{\partial \eta_j} = \frac{\partial F_{\eta_j}^j}{\partial \eta_j}, \quad \frac{\partial \tilde{\mathbf{r}}_{j+1}}{\partial \tilde{\mathbf{r}}_j} = \frac{\partial F_{\eta_j}^j}{\partial \tilde{\mathbf{r}}_j}.$$

Invoking the chain rule with (3.7) and (3.8), $\frac{\partial \bar{f}_{\bar{\theta}}}{\partial \eta_j}$ is described as

$$(3.9) \qquad \frac{\partial \bar{f}_{\bar{\theta}}}{\partial \eta_j} = \frac{\partial \mathbf{y}}{\partial \tilde{\mathbf{y}}} \frac{\partial \tilde{\mathbf{y}}}{\partial \tilde{\mathbf{r}}_N} \left( \prod_{l=j+1}^{N-1} \frac{\partial \tilde{\mathbf{r}}_{l+1}}{\partial \tilde{\mathbf{r}}_l} \right) \frac{\partial \tilde{\mathbf{r}}_{j+1}}{\partial \eta_j} = \frac{\partial h_\varepsilon}{\partial \tilde{\mathbf{y}}} \left( \prod_{l=j+1}^{N-1} \frac{\partial F_{\eta_l}^l}{\partial \tilde{\mathbf{r}}_l} \right) \frac{\partial F_{\eta_j}^j}{\partial \eta_j}.$$

On the other hand, partial derivatives $\frac{\partial \mathbf{x}_j}{\partial \delta_j}$, $\frac{\partial \mathbf{y}_j}{\partial \phi_j}$, and $\frac{\partial \mathbf{y}_j}{\partial \mathbf{x}_j}$ can be computed in parallel by (3.1):

$$(3.10) \qquad \frac{\partial \mathbf{x}_j}{\partial \delta_j} = \frac{\partial C_{\delta_j}^j}{\partial \delta_j}, \quad \frac{\partial \mathbf{y}_j}{\partial \phi_j} = \frac{\partial g_{\phi_j}^j}{\partial \phi_j}, \quad \frac{\partial \mathbf{y}_j}{\partial \mathbf{x}_j} = \frac{\partial g_{\phi_j}^j}{\partial \mathbf{x}_j}.$$

Using (3.7), (3.8), and (3.10), it follows that

$$
\frac{\partial \bar{f}_{\bar{\theta}}}{\partial \phi_j} = 
\begin{cases}
\dfrac{\partial \mathbf{y}}{\partial \tilde{\mathbf{y}}} \dfrac{\partial \tilde{\mathbf{y}}}{\partial \tilde{\mathbf{r}}_N} \left( \displaystyle\prod_{l=j}^{N-1} \dfrac{\partial \tilde{\mathbf{r}}_{l+1}}{\partial \tilde{\mathbf{r}}_l} \right) \dfrac{\partial \tilde{\mathbf{r}}_j}{\partial \mathbf{r}_j} \dfrac{\partial \mathbf{r}_j}{\partial \mathbf{y}_j} \dfrac{\partial \mathbf{y}_j}{\partial \phi_j} & \text{if } j < N, \\[2ex]
\dfrac{\partial \mathbf{y}}{\partial \tilde{\mathbf{y}}} \dfrac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{y}_N} \dfrac{\partial \mathbf{y}_N}{\partial \phi_N} & \text{if } j = N
\end{cases}
$$

$$(3.11) \qquad = \frac{\partial h_\varepsilon}{\partial \tilde{\mathbf{y}}} \left( \prod_{l=j}^{N-1} \frac{\partial F_{\eta_l}^l}{\partial \tilde{\mathbf{r}}_l} \right) \frac{\partial g_{\phi_j}^j}{\partial \phi_j}.$$

Similarly, with the convention $\frac{\partial F_{\eta_0}^0}{\partial \tilde{\mathbf{r}}_0} = 0$, we have

$$
\frac{\partial \bar{f}_{\bar{\theta}}}{\partial \delta_j} = 
\begin{cases}
\dfrac{\partial \mathbf{y}}{\partial \tilde{\mathbf{y}}} \dfrac{\partial \tilde{\mathbf{y}}}{\partial \tilde{\mathbf{r}}_N} \left( \displaystyle\prod_{l=j}^{N-1} \dfrac{\partial \tilde{\mathbf{r}}_{l+1}}{\partial \tilde{\mathbf{r}}_l} \right) \left( \dfrac{\partial \tilde{\mathbf{r}}_j}{\partial \mathbf{r}_j} \dfrac{\partial \mathbf{r}_j}{\partial \mathbf{y}_j} \dfrac{\partial \mathbf{y}_j}{\partial \mathbf{x}_j} + \dfrac{\partial \tilde{\mathbf{r}}_j}{\partial \tilde{\mathbf{r}}_{j-1}} \dfrac{\partial \tilde{\mathbf{r}}_{j-1}}{\partial \mathbf{r}_{j-1}} \dfrac{\partial \mathbf{r}_{j-1}}{\partial \mathbf{x}_j} \right) \dfrac{\partial \mathbf{x}_j}{\partial \delta_j} & \text{if } j < N, \\[2ex]
\dfrac{\partial \mathbf{y}}{\partial \tilde{\mathbf{y}}} \left( \dfrac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{y}_N} \dfrac{\partial \mathbf{y}_N}{\partial \mathbf{x}_N} + \dfrac{\partial \tilde{\mathbf{y}}}{\partial \tilde{\mathbf{r}}_N} \dfrac{\partial \tilde{\mathbf{r}}_N}{\partial \tilde{\mathbf{r}}_{N-1}} \dfrac{\partial \tilde{\mathbf{r}}_{N-1}}{\partial \mathbf{r}_{N-1}} \dfrac{\partial \mathbf{r}_{N-1}}{\partial \mathbf{x}_N} \right) \dfrac{\partial \mathbf{x}_N}{\partial \delta_N} & \text{if } j = N
\end{cases}
$$

$$(3.12) \qquad = \frac{\partial h_\varepsilon}{\partial \tilde{\mathbf{y}}} \left( \prod_{l=j}^{N-1} \frac{\partial F_{\eta_l}^l}{\partial \tilde{\mathbf{r}}_l} \right) \left( \frac{\partial g_{\phi_j}^j}{\partial \mathbf{x}_j} - \frac{\partial F_{\eta_{j-1}}^{j-1}}{\partial \tilde{\mathbf{r}}_{j-1}} \right) \frac{\partial C_{\delta_j}^j}{\partial \delta_j}.$$

For efficient computation, the value of $\frac{\partial h_\varepsilon}{\partial \tilde{\mathbf{y}}} \left( \prod_{l=j}^{N-1} \frac{\partial F_{\eta_l}^l}{\partial \tilde{\mathbf{r}}_l} \right)$ can be stored during the evaluation of (3.9) and then used in (3.11) and (3.12). Such a technique is described in Algorithm 3.2.

**4. Applications.** In this section, we present applications of the proposed methodology to two existing convolutional neural networks (CNNs): VGG-16 [32] and ResNet-1001 [16]. For each network, we deal with details on the construction of parallel subnetworks and a coarse network. Also, numerical results are presented showing that the proposed parareal neural network gives comparable or better results than given feed-forward neural network and other variants in terms of both training time and classification performance.

First, we present details on the datasets we used. The CIFAR-$m$ ($m = 10, 100$) dataset consists of $32 \times 32$ colored natural images and includes 50,000 training and 10,000 test samples with $m$ classes. The SVHN dataset is composed of $32 \times 32$ colored digit images; there are 73,257 and 26,032 samples for training and test, respectively, with additional 531,131 training samples. However, we did not use the additional ones for training. MNIST is a classic dataset which contains handwritten digits encoded in $28 \times 28$ grayscale images. It includes 55,000 training, 5,000 validation, and 10,000 test samples. In our experiments, the training and validation

| Layer | Output size | VGG-16 |
|---|---|---|
| Preprocessing | $224 \times 224$ | $[3 \times 3, 64]$ |
| | $112 \times 112$ | $[3 \times 3, 64]$ + maxpool |
| Block-repetitive substructure | $56 \times 56$ | $[3 \times 3, 128] \times 2$ + maxpool |
| | $28 \times 28$ | $[3 \times 3, 256] \times 3$ + maxpool |
| | $14 \times 14$ | $[3 \times 3, 512] \times 3$ + maxpool |
| | $7 \times 7$ | $[3 \times 3, 512] \times 3$ + maxpool |
| Postprocessing | $1 \times 1$ | [4096-d FC]$\times 2$ + 1000-d FC |

samples are used as training data and the test samples as test data. ImageNet is a dataset which contains 1000 classes of $224 \times 224$ colored natural images. It includes 1,280,000 training and 50,000 test samples.

We adopted a data augmentation technique in [24] for CIFAR datasets; four pixels are padded on each side of images, and $32 \times 32$ crops are randomly sampled from the padded images and their horizontal flips. All neural networks in this section were trained using the stochastic gradient descent with the batch size 128, weight decay 0.0005, Nesterov momentum 0.9, and weights initialized as in [14]. The initial learning rate was set to 0.1, and was reduced by a factor of 10 in the 80th and 120th epochs. For ImageNet datasets, the input image is $224 \times 224$ randomly cropped from a resized image using the scale and aspect ratio augmentation [33]. Hyperparameter settings are the same as other cases except the followings; the weight decay 0.0001, total epoch 90, and the learning rate was reduced by a factor of 10 in the 30th and 60th epochs. All networks were implemented in Python with PyTorch and all computations were performed on a cluster equipped with Intel Xeon Gold 5515 (2.4GHz, 20C) CPUs, NVIDIA Titan RTX GPUs, and the operating system Ubuntu 18.04 64bit.

**4.1. VGG-16.** In general, CNNs without skip connections (see, e.g., [21, 32]) can be represented as

$$\mathbf{x}_l = H_l(\mathbf{x}_{l-1}),$$

where $\mathbf{x}_l$ is an output of the $l$th layer of the network and $H_l$ is a nonlinear transformation consisting of convolutions, batch normalization and ReLU activation. Each layer of VGG-16 [32], one of the most popular CNNs without skip connections, consists of multiple $3 \times 3$ convolutions. The network consists of 5 stages of convolutional blocks and 3 fully connected layers. Each convolutional block is a composition of double or triple convolutions and a max pooling operation.

**4.1.1. Parareal transformation.** First, we describe the structure of VGG-16 which was designed for the classification problem of ImageNet [7] with the terminology introduced in Section 3. Inputs for VGG-16 are 3-channel images with $224 \times 224$ pixels, i.e., $X = \mathbb{R}^{3 \times 224 \times 224}$.

**Table 2**
*Error rates (%) on the ImageNet dataset of VGG-16 and Parareal VGG-16-4.*

| Network | Subnetwork Parameters | Coarse network Parameters | Total Parameters | Error rate (%) |
|---------|----------------------|--------------------------|------------------|----------------|
| VGG-16 | - | - | 138.4M | 32.30 |
| Parareal VGG-16-4 | 3.7M | 9.1M | 147.5M | 30.97 |

The output space $Y$ is given by $Y = \mathbb{R}^{1000}$, where $1,000$ is the number of classes of ImageNet. We set the preprocessing operator $C_\delta \colon X \to W_0 = \mathbb{R}^{64 \times 224 \times 224}$ by the first $3 \times 3$ convolution layer in VGG-16. We refer to the remaining parts of VGG-16 as the block-repetitive structure $g_\phi \colon W_0 \to W_1$ with $W_1 = \mathbb{R}^{512 \times 7 \times 7}$ except for the last three fully connected layers. Finally, the postprocessing operator $h_\varepsilon \colon W_1 \to Y$ is the composition of the three fully connected layers. Table 1 shows the detailed architecture of VGG-16.

In order to construct a parareal neural network with $N$ parallel subnetworks for VGG-16, we have to specify its components $g_{\phi_j}^j$, $C_{\delta_j}^j$, and $F_{\eta_j}^j$. For simplicity, we assume that $N = 4$. We decompose $g_\phi$ into 4 parts such that the output size of each part is $56 \times 56$, $28 \times 28$, $14 \times 14$, and $7 \times 7$, respectively. Then, the block-repetitive structure $g_\phi$ can be decomposed as

$$g_\phi = g_{\phi_4}^4 \circ g_{\phi_3}^3 \circ g_{\phi_2}^2 \circ g_{\phi_1}^1,$$

where each of $g_{\phi_j}^j \colon X_{j-1} \to X_j$ with

$$X_j = \begin{cases} \mathbb{R}^{128 \times 56 \times 56} & \text{for } j = 1 \\ \mathbb{R}^{256 \times 28 \times 28} & \text{for } j = 2, \\ \mathbb{R}^{512 \times 14 \times 14} & \text{for } j = 3, \\ \mathbb{R}^{512 \times 7 \times 7} & \text{for } j = 4, \end{cases} \qquad \phi = \bigoplus_{j=1}^{4} \phi_j.$$

Recall that the main role of the preprocessing operator $C_{\delta_j}^j \colon X \to X_{j-1}$ is to transform an input $\mathbf{x} \in X$ to fit in the space $X_{j-1}$. In this perspective, we simply set $C_{\delta_1}^1 = C_\delta$ and $C_{\delta_j}^j$ for $j > 1$ by a $1 \times 1$ convolution to match the number of channels after appropriate number of $3 \times 3$ max pooling layers with stride 2 to match the image size.

According to Proposition 3.1, it is essential to design the coarse network such that $F_{\eta_j}^j \approx g_{\phi_{j+1}}^{j+1}$ in order to ensure the performance of the parareal neural network. We simply define $F_{\eta_j}^j \colon X_j \to X_{j+1}$ by the composition of two $3 \times 3$ convolutions and a max pooling with kernel size 2 and stride 2, which has a simplified structure of $g_{\phi_{j+1}}^{j+1}$ with fewer parameters.

**4.1.2. Numerical results.** We present the comparison results with Parareal VGG-16-4 and VGG-16 on ImageNet dataset. Note that Parareal VGG-16-4 denotes the parareal neural network version of VGG-16 with $N = 4$. Table 2 shows that the error rate of Parareal VGG-16-4 is smaller than that of VGG-16. From this result, it can be seen that the accuracy is guaranteed even when the parareal neural network is applied to a network with small number of layers.

**Table 3**

*Forward/backward computation time for VGG-16 and Parareal VGG-16-4. It is a measure of the time taken in one iteration for ImageNet dataset input* $\mathbf{x} \in \mathbb{R}^{3 \times 224 \times 224}$ *with batch size* 128. *One interation means one step of updating all parameters with SGD.*

| | Virtual wall-clock time (ms) | | | | |
|---|---|---|---|---|---|
| Network | Preprocessing | Parallel subnetworks | Coarse network | Postprocessing | Total |
| VGG-16 | 18.73/321.22 | 187.21/3302.11 | - | 3.33/5.82 | 209.27/3629.15 |
| Parareal VGG-16-4 | 18.78/334.10 | 46.59/453.25 | 62.84/1413.49 | 2.96/5.90 | 131.17/2206.74 |

**Table 4**

*Error rates (%) and wall-clock times on the ImageNet dataset. The wall-clock time is the total time taken to train a given network by* 200 *epochs. Relative speed-up is measured according to* (4.1).

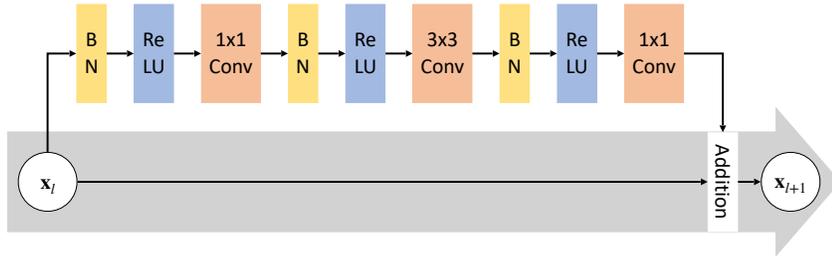| Network | Parameters | Error rate (%) | Wall-clock time (h:m:s) | RS (%) |
|---|---|---|---|---|
| VGG-16 | 138.4M | 32.30 | 191:17:00 | 0.0 |
| Data Parallel VGG-16-4 | 553.6M | 27.73 | 142:56:18 | 25.3 |
| Parareal VGG-16-4 | 147.5M | 30.97 | 188:22:18 | 1.5 |

Next, we investigate the elapsed time for forward and backward propagations of parareal neural networks, which are the most time-consuming part of network training. Table 3 shows the virtual wall-clock time for forward and backward computation of VGG-16 and Parareal VGG-16-4 for the input $\mathbf{x} \in \mathbb{R}^{3 \times 224 \times 224}$. Note that the virtual wall-clock time is measured under the assumption that all parallelizable procedures indicated in Algorithm 3.2 are executed simultaneously. It means that, it excludes the communication time among GPUs. In Table 3, even if the VGG-16 has the small number of layers, it can be seen that the computation time of Parareal-VGG-16-4 is reduced in terms of virtual wall clock time.

Now, we compare the performance of the proposed methodology to data parallelism. In data parallelism, the batch is split into subsets at each epoch. Each subset is assigned one for each GPU and the gradient corresponding to the subset is computed in parallel. Then the parameters of the neural network are updated by the averaged gradient over all subsets. In what follows, Data Parallel VGG-16-4 denotes the data parallelized VGG-16 with 4 GPUs. We compare the error rate and the wall-clock time of each parallelized network with the ImageNet dataset. To investigate the speed-up provided by each parallelism, we use the relative speed-up (RS) introduced in [9], which is defined by

$$(4.1) \qquad \mathrm{RS} = \frac{t_r - t_p}{t_r},$$

where $t_p$ is the total elapsed time taken to complete training of the given parallelism and $t_r$ is that of the given feed-forward network.

Table 4 shows that the wall-clock time of Data Parallel VGG-16-4 is the shortest. This is because Data Parallel VGG-16-4 has very small number of layers, which takes a short time for communication to average the gradients for updating parameters. On the other hand, in

**Figure 3.** *Bottleneck structure of an RU used in ResNet-1001. The first $1 \times 1$ convolution squeezes the number of channels of an input and the last $1 \times 1$ convolution increases the number of channels of an intermediate result. If $\mathbf{x}_l$ and $\mathbf{x}_{l+1}$ have different numbers of channels, then a $1 \times 1$ convolution is added to the skip connection in order to match the number of channels.*

**Table 5**

*Architecture of ResNet-1001 for the dataset CIFAR-100. Downsampling with stride 2 is performed in the second and third stages of the block-repetitive structure. A layer consisting of an $n \times n$ convolution with $k$-channel ouput, an average pooling with output size $1 \times 1$, and a $k$-way fully connected layer are denoted by $[n \times n, k]$, avgpool, and $k$-d FC, respectively.*

| Layer | Output size | ResNet-1001 |
|---|---|---|
| Preprocessing | $32 \times 32$ | $[3 \times 3, \ 16]$ |
| | $32 \times 32$ | $\begin{bmatrix} 1 \times 1, 16 \\ 3 \times 3, 16 \\ 1 \times 1, 64 \end{bmatrix} \times 111$ |
| Block-repetitive substructure | $16 \times 16$ | $\begin{bmatrix} 1 \times 1, 32 \\ 3 \times 3, 32 \\ 1 \times 1, 128 \end{bmatrix} \times 111$ |
| | $8 \times 8$ | $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 111$ |
| Postprocessing | $1 \times 1$ | avgpool + 100-d FC |

the case of Parareal VGG-16-4, as shown in Table 3, the computation time of the parallel subnetwork is very small, and the majority of the total computation time is the computation time of the coarse network. Therefore, Parareal VGG-16-4 has a parallel structure, but the effect of parallelization is not significant. In fact, the parareal algorithm works more effectively in networks with a large number of layers, but maintains accuracy and does not slow down even in networks with a small number of layers.

**4.2. ResNet-1001.** Next, we will apply the proposed parareal neural network to ResNet [15, 16], which is typically one of the very deep neural networks. In ResNet, an RU is constructed by adding a skip connection, i.e., it is written as

$$\mathbf{x}_l = H_l(\mathbf{x}_{l-1}) + \mathbf{x}_{l-1}.$$

Figure 3 shows the structure of an RU, called bottleneck used in ResNet. In particular, we deal with ResNet-1001 for the CIFAR dataset [20], where $1,001$ is the number of layers. ResNet-1001 consists of a single $3 \times 3$ convolution layer, a sequence of 333 RUs with varying feature map dimensions, and a global average pooling followed by a fully connected layer. Table 5 shows the details of the architecture of ResNet-1001 for the CIFAR-100 dataset.

As presented in Section 3, ResNet-1001 can be decomposed as follows: a preprocessing operator $C_\delta \colon X \to W_0$ as the $3 \times 3$ convolution layer, a block-repetitive substructure $g_\phi \colon W_0 \to W_1$ as 333 RUs, and a postprocessing operator $h_\varepsilon \colon W_1 \to Y$ consisting of the global average pooling and the fully connected layer. More specifically, we have $X = \mathbb{R}^{3 \times 32 \times 32}$, $W_0 = \mathbb{R}^{16 \times 32 \times 32}$, $W_1 = \mathbb{R}^{256 \times 8 \times 8}$, and $Y = \mathbb{R}^m$ where $m$ is the number of classes of images.

**4.2.1. Parareal transformation.** The design of a parareal neural network with $N$ parallel subnetworks for ResNet-1001, denoted as *Parareal ResNet-N*, can be completed by specifying the structures $g_{\phi_j}^j$, $C_{\delta_j}^j$, and $F_{\eta_j}^j$. For convenience, the original neural network ResNet-1001 is called Parareal ResNet-1. We assume that $N = 3N_0$ for some positive integer $N_0$. We observe that $g_\phi$ can be decomposed as

$$g_\phi = g_{\phi_N}^N \circ \cdots \circ g_{\phi_{2N_0+1}}^{2N_0+1} \circ g_{\phi_{2N_0}}^{2N_0} \circ \cdots \circ g_{\phi_{N_0+1}}^{N_0+1} \circ g_{\phi_{N_0}}^{N_0} \circ \cdots \circ g_{\phi_1}^1,$$

where each of $g_{\phi_j}^j \colon X_{j-1} \to X_j$ consists of $\lceil 333/N \rceil$ RUs with

$$X_j = \begin{cases} \mathbb{R}^{64 \times 32 \times 32} & \text{for } j = 1, \ldots, N_0, \\ \mathbb{R}^{128 \times 16 \times 16} & \text{for } j = N_0 + 1, \ldots, 2N_0, \\ \mathbb{R}^{256 \times 8 \times 8} & \text{for } j = 2N_0 + 1, \ldots, N, \end{cases} \quad \phi = \bigoplus_{j=1}^{N} \phi_j.$$

Similarly to the case of VGG-16, the preprocessing operators are defined as follows: $C_{\delta_1}^1 = C_\delta$ and $C_{\delta_j}^j$ for $j > 1$ consists of a $1 \times 1$ convolution to match the number of channels after appropriate number of $3 \times 3$ max pooling layers with stride 2 to match the image size. For the coarse network, we first define a coarse RU consisting of two $3 \times 3$ convolutions and skip-connection. If the downsampling is needed, then the stride of first convolution in coarse RU is set to 2. We want to define $F_{\eta_j}^j \colon X_j \to X_{j+1}$ having smaller number of (coarse) RUs than $g_{\phi_{j+1}}^{j+1}$ but a similar coverage to $g_{\phi_{j+1}}^{j+1}$. Note that the receptive field of $g_{\phi_j}^j$ covers the input size $32 \times 32$. In terms of receptive field, even if we construct $F_{\eta_j}^j$ with fewer coarse RUs than $\lceil 333/N \rceil$, it can have similar coverage to the parallel subnetwork $g_{\phi_j}^j$.

Let $N_c$ be the number of coarse RUs in $F_{\eta_j}^j$ of the coarse network. Here, we present how $N_c$ affects the performance of the parareal neural network. Specifically, we experimented with Parareal ResNet-3 with $g_{\phi_j}^j$ consisting of 111 RUs. Table 6 shows the error rates of Parareal ResNet-3 with respect to various $N_c$. The performance of Parareal ResNet-3 with $N_c = 4$ is similar to ResNet-1001. In fact, four coarse RUs consist of one $3 \times 3$ convolution with stride 2 and seven $3 \times 3$ convolutions so that the receptive field covers $31 \times 31$ pixels, which is almost all pixels of the input $\mathbf{x} \in X$. Parareal ResNet-3 shows better performance than ResNet-1001 when $N_c > 4$. Since Parareal ResNet-3 with $N_c = 4$ shows similar performance to ResNet-1001, we may say that 1 RUs in $F_{\eta_j}^j$ can approximate $111/4 \approx 28$ RUs in $g_{\phi_j}^j$. Generally, if we use $N$

*Error rates (%) on the CIFAR-100 dataset of Parareal ResNet-3, where $N_c$ is the number of coarse RUs in each component $F_{\eta_j}^j$ of the coarse network.*

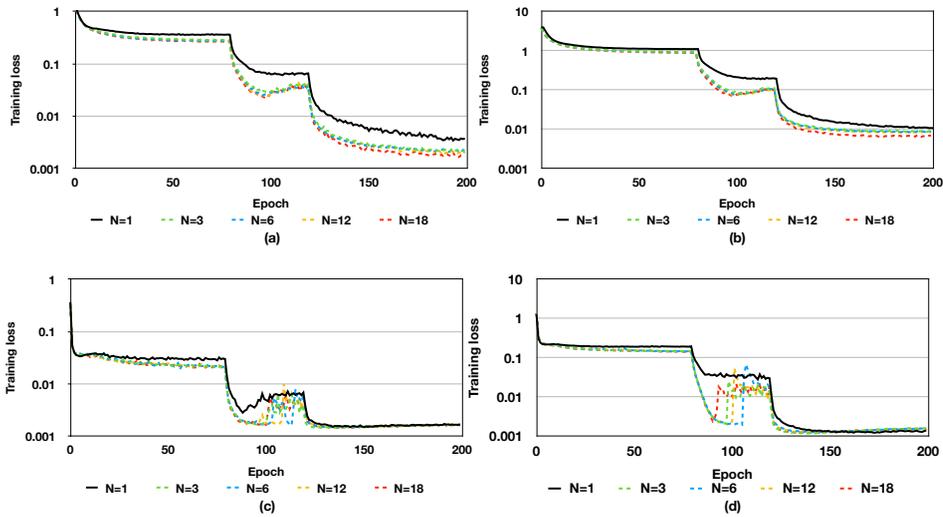| $N_c$ | Error rate (%) |
|:---:|:---:|
| 1 | 23.47 |
| 2 | 22.20 |
| 4 | 21.14 |
| 8 | 20.85 |
| 16 | 20.83 |

Reference (ResNet-1001): 21.13



**Figure 4.** *Comparison of the training loss for Parareal ResNet-N  ($N = 1, 3, 6, 12, 18$) on various datasets:* *(a) CIFAR-10, (b) CIFAR-100, (c) MNIST, and (d) SVHN results.*

parallel subnetworks ($N \geq 3$), each $333/N$ RUs in $g_{\phi_j}^j$ can be approximated by the $N_c$ RUs in $F_{\eta_j}^j$ whenever we select $N_c = \lceil 12/N \rceil$.

**4.2.2. Numerical results.** With fixed $N_c = \lceil 12/N \rceil$, we report the classification results of Parareal ResNet with respect to various $N$ on datasets CIFAR-10, CIFAR-100, SVHN, and MNIST. Decay of the training loss of Parareal ResNet-$N$ ($N = 1, 3, 6, 12, 18$) for various datasets is depicted in Figure 4. As shown in Figure 4, the training loss converges to a smaller value for larger $N$. It seems that such a phenomenon is due to the increase of the number of parameters in Parareal ResNet when $N$ increases. In the cases of MNIST and SVHN, oscillations of the training loss are observed. It is well-known that such oscillations are caused by excessive weight decay and can be removed by dropout [36].

We next study how the classification performance is affected by $N$, the number of sub-networks. Table 7 shows that the error rates of Parareal ResNet-$N$ are usually smaller than

*Error rates (%) on the CIFAR-10, CIFAR-100, MNIST, and SVHN datasets of Parareal ResNet-N $(N = 1, 3, 6, 12, 18)$ with $N_c = \lceil 12/N \rceil$.*

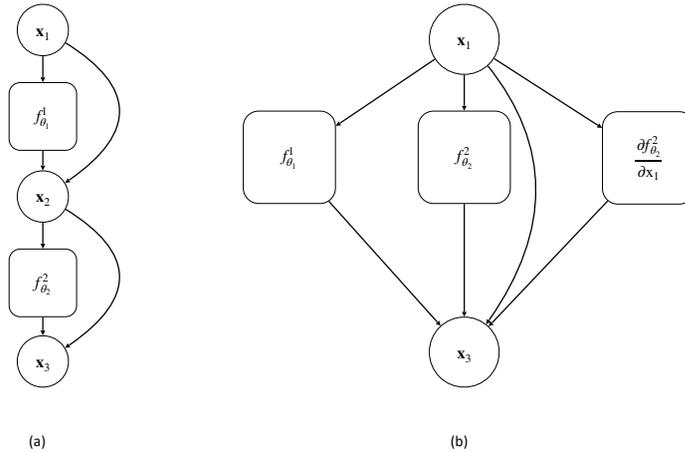| $N$ | Subnetwork Parameters | Coarse network Parameters | Total Parameters | CIFAR-10 | CIFAR-100 | MNIST | SVHN |
|---|---|---|---|---|---|---|---|
| 1 | - | - | 10.3M | 4.96 | 21.13 | 0.34 | 3.17 |
| 3 | 3.4M | 5.6M | 15.9M | 4.61 | 21.14 | 0.31 | 3.11 |
| 6 | 1.7M | 5.7M | 16.1M | 4.20 | 20.87 | 0.31 | 3.21 |
| 12 | 0.9M | 5.8M | 16.2M | 4.37 | 20.42 | 0.28 | 3.25 |
| 18 | 0.6M | 8.9M | 19.4M | 4.02 | 20.40 | 0.33 | 3.29 |

*Forward/backward computation time for Parareal ResNet-N $(N = 1, 3, 6, 12, 18)$. The time is measured in one iteration for CIFAR-100 dataset input $\mathbf{x} \in \mathbb{R}^{3 \times 32 \times 32}$ with batch size $128$.*

| | Virtual wall-clock time (ms) | | | | |
|---|---|---|---|---|---|
| $N$ | Preprocessing | Parallel subnetworks | Coarse network | Postprocessing | Total |
| 1 | 0.25/6.46 | 443.81/1387.62 | - | 0.06/3.18 | 444.12/1397.26 |
| 3 | 0.25/6.45 | 131.92/458.87 | 10.01/97.60 | 0.06/3.71 | 142.24/566.63 |
| 6 | 0.27/6.42 | 67.59/219.72 | 14.68/137.08 | 0.06/3.33 | 82.60/366.55 |
| 12 | 0.28/6.59 | 48.47/113.33 | 17.97/149.52 | 0.06/3.63 | 66.78/273.07 |
| 18 | 0.29/6.17 | 30.40/77.84 | 27.87/163.25 | 0.06/3.64 | 58.62/250.90 |
| 24 | 0.29/6.58 | 22.71/58.04 | 41.03/242.87 | 0.06/3.54 | 64.09/311.03 |

ResNet-1001. There are some exceptional cases that the error rate of Parareal ResNet-$N$ exceeds that of ResNet-1001: $N = 6, 12$ and 18 for SVHN. It is known that these cases occur when there are oscillations in the decay of the training loss [37]; see Figure 4. As we mentioned above, such oscillations can be handled by dropout.

Like the case of VGG-16, we investigate the elapsed time for forward and backward propagations of parareal neural networks. Table 8 shows the virtual wall-clock time for forward and backward computation of Parareal ResNet-$N$ with various $N$ for the input $\mathbf{x} \in \mathbb{R}^{3 \times 32 \times 32}$. As shown in Table 8, the larger $N$, the shorter the computing time of the parallel subnetworks $g_{\phi_j}^j$, while the longer the computing time of the coarse network. This is because as $N$ increases, the depth of each parallel subnetwork $g_{\phi_j}^j$ becomes shallower while the number of $F_{\eta_j}^j$ in the coarse network increases. On the other hand, each preprocessing operator $C_{\delta_j}^j$ is designed to be the same as or similar to the preprocessing operator $C_\delta$ of the original neural network and the postprocessing operator $h_\varepsilon$ is the same as the original one. Therefore, the computation time for the pre- and postprocessing operators does not increase even as $N$ increases. As a trade-off between the decrease in time in parallel subnetworks and the increase in time in the coarse network, we observe in Table 8 that the elapsed time of forward and backward computation decreases as $N$ increases when $N \leq 18$.

(a)                                          (b)

**Figure 5.** *Structure of the $K$-warp operator, $K = 2$: **(a)** Two RUs $f^1_{\theta_1}$ and $f^2_{\theta_2}$ in ResNet. **(b)** The corresponding $2$-warp operator in WarpNet.*

*Remark* 4.1. In our experiments, the decomposition of the block-repetitive structure of ResNet-1001 was done so that each parallel subnetwork of Parareal ResNet had the same number of RUs. However, such a decomposition is not optimal because each RU has the different computational cost. It is expected that the uniform decomposition of the Parareal ResNet in terms of computational cost will further improve the parallel efficiency reducing virtual wall-clock time even above $N = 18$.

Finally, we compare the performance of the Parareal ResNet to two existing approaches of parallel computing for neural networks: WarpNet [9] and data parallelism. WarpNet replaces every $K$ RUs in ResNet by a $K$-*warp operator*, which can be evaluated using RUs arranged in a parallel manner and their derivatives. For example, we display the parallel structure of the 2-warp operator in Figure 5; see [9] for further details. In what follows, Data Parallel ResNet-$N$ and WarpNet-$N$ denote the data parallelized ResNet-1001 with $N$ GPUs and WarpNet with $(N-1)$-warp operators using $N$ GPUs, respectively. We compare the error rate and the wall-clock time of each parallelized network with the CIFAR-100 dataset. Table 9 shows that only RS of Parareal ResNet is greater than 0. This is because when the number of layers of WarpNet and Data Parallel ResNet is very large, data communication between GPUs is required, which can frequently cause communication bottlenecks at each layer. On the other hand, in the forward propagation of Parareal ResNet, only a single communication process among GPUs is needed: communication between each parallel subnetwork and the coarse network. Therefore, Parareal ResNet is relatively free from communication bottlenecks compared to the other methods. In conclusion, Parareal ResNet outperforms the other two methods in parallelization of deep neural networks in terms of training time reduction.

*Remark* 4.2. In terms of memory efficiency, Data Parallel VGG-16 and Data Parallel ResNet allocate all parameters to each GPU, which is a waste of memory of GPUs. WarpNet-$N$ requires duplicates of parameters since it computes RUs and their derivatives simultaneously at forward propagation so that the memory efficiency is deteriorated. On the other hand, in Parareal

**Table 9**
*Error rates (%) and wall-clock times on the CIFAR-100 dataset. The wall-clock time is the total time taken to train a given network by 200 epochs.*

| Network | Parameters | Error rate (%) | Wall-clock time (h:m:s) | RS (%) |
|---|---|---|---|---|
| ResNet-1001 | 10.3M | 21.13 | 22:44:53 | 0.0 |
| Data Parallel ResNet-3 | 30.9M | 21.09 | 46:50:20 | -105.9 |
| Parareal ResNet-3 | 15.9M | 21.14 | 16:28:38 | 27.6 |
| WarpNet-3 | 15.3M | 19.95 | 90:56:46 | -299.8 |
| Data Parallel ResNet-6 | 61.8M | 20.42 | 68:10:06 | -199.7 |
| Parareal ResNet-6 | 16.1M | 20.87 | 11:48:13 | 48.1 |

VGG-16 and Parareal ResNet, the entire model can be equally distributed to each GPU as a parallel subnetwork; only a single GPU requires additional memory to store the parameters of the coarse network which is much smaller than the entire model. Therefore, the proposed Parareal neural network outperforms the other two methods in the terms of memory efficiency as well.

**5. Conclusion.** In this paper, we proposed a novel methodology to construct a parallel neural network called the parareal neural network, which is suitable for parallel computation using multiple GPUs from a given feed-forward neural network. Motivated by the parareal algorithm for time-dependent differential equations, the block-repetitive part of the original neural network was partitioned into small pieces to form parallel subnetworks of the parareal neural network. The coarse network that corrects differences at the interfaces among subnetworks was introduced so that the performance of the resulting parareal network agrees with the original network. As a concrete example, we presented how to design the parareal neural network corresponding to VGG-16 and ResNet-1001. Numerical results were provided to highlight the robustness and parallel efficiency of Parareal VGG-16 and Parareal ResNet. To the best of our knowledge, the proposed methodology is a new kind of multi-GPU parallelism in the field of deep learning.

**REFERENCES**

[1] G. Bal, *On the convergence and the stability of the parareal algorithm to solve partial differential equations*, in Domain Decomposition Methods in Science and Engineering, Springer, Berlin, 2005, pp. 425–432.
[2] T. Ben-Nun and T. Hoefler, *Demystifying parallel and distributed deep learning: An in-depth concurrency analysis*, ACM Computing Surveys (CSUR), 52 (2019), pp. 1–43.
[3] C.-C. Chen, C.-L. Yang, and H.-Y. Cheng, *Efficient and robust parallel DNN training through model parallelism on multi-GPU platform*, Sept. 2014, https://arxiv.org/abs/1809.02839.
[4] J. Chen, R. Monga, S. Bengio, and R. Jozefowicz, *Revisiting distributed synchronous SGD*, in International Conference on Learning Representations Workshop Track, 2016, https://arxiv.org/abs/1604.00981.

[5] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa, *Natural language processing (almost) from scratch*, Journal of Machine Learning Research, 12 (2011), pp. 2493–2537, http://jmlr.org/papers/v12/collobert11a.html.

[6] G. Cybenko, *Approximation by superpositions of a sigmoidal function*, Mathematics of Control, Signals and Systems, 2 (1989), pp. 303–314.

[7] 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, Ieee, 2009, pp. 248–255.

[8] R. D. Falgout, S. Friedhoff, T. V. Kolev, S. P. MacLachlan, and J. B. Schroder, *Parallel time integration with multigrid*, SIAM Journal on Scientific Computing, 36 (2014), pp. C635–C661.

[9] R. Fok, A. An, Z. Rashidi, and X. Wang, *Decoupling the layers in residual networks*, in International Conference on Learning Representations, 2018.

[10] J. H. Friedman, *Greedy function approximation: a gradient boosting machine*, Annals of Statistics, (2001), pp. 1189–1232.

[11] M. J. Gander and E. Hairer, *Nonlinear convergence analysis for the parareal algorithm*, in Domain Decomposition Methods in Science and Engineering XVII, Springer, Berlin, 2008, pp. 45–56.

[12] M. J. Gander and S. Vandewalle, *Analysis of the parareal time-parallel time-integration method*, SIAM Journal on Scientific Computing, 29 (2007), pp. 556–578.

[13] S. Günther, L. Ruthotto, J. B. Schroder, E. C. Cyr, and N. R. Gauger, *Layer-parallel training of deep residual neural networks*, SIAM Journal on Mathematics of Data Science, 2 (2020), pp. 1–23.

[14] K. He, X. Zhang, S. Ren, and J. Sun, *Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification*, in Proceedings of the IEEE International Conference on Computer Vision, 2015, pp. 1026–1034.

[15] K. He, X. Zhang, S. Ren, and J. Sun, *Deep residual learning for image recognition*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 770–778.

[16] K. He, X. Zhang, S. Ren, and J. Sun, *Identity mappings in deep residual networks*, in European Conference on Computer Vision, Springer, Cham, 2016, pp. 630–645.

[17] K. Hornik, M. Stinchcombe, and H. White, *Multilayer feedforward networks are universal approximators*, Neural Networks, 2 (1989), pp. 359–366.

[18] Y. Huang, Y. Cheng, A. Bapna, O. Firat, M. X. Chen, D. Chen, H. Lee, J. Ngiam, Q. V. Le, Y. Wu, and Z. Chen, *GPipe: Efficient training of giant neural networks using pipeline parallelism*, in Advances in Neural Information Processing Systems, 2019, pp. 103–112.

[19] S. Jean, K. Cho, R. Memisevic, and Y. Bengio, *On using very large target vocabulary for neural machine translation*, in Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), 2015, pp. 1–10.

[20] A. Krizhevsky, V. Nair, and G. Hinton, *CIFAR-10 and CIFAR-100 datasets*, https://www.cs.toronto.edu/kriz/cifar.html, 6 (2009).

[21] A. Krizhevsky, I. Sutskever, and G. E. Hinton, *ImageNet classification with deep convolutional neural networks*, in Advances in Neural Information Processing Systems, 2012, pp. 1097–1105.

[22] Y. LeCun, Y. Bengio, and G. Hinton, *Deep learning*, Nature, 521 (2015), pp. 436–444.

[23] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, *Gradient-based learning applied to document recognition*, Proceedings of the IEEE, 86 (1998), pp. 2278–2324.

[24] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu, *Deeply-supervised nets*, in Artificial Intelligence and Statistics, 2015, pp. 562–570.

[25] X. Lian, Y. Huang, Y. Li, and J. Liu, *Asynchronous parallel stochastic gradient for nonconvex optimization*, in Advances in Neural Information Processing Systems, 2015, pp. 2737–2745.

[26] J.-L. Lions, Y. Maday, and G. Turinici, *Résolution d'EDP par un schéma en temps «pararéel»*, Comptes Rendus de l'Académie des Sciences-Series I-Mathematics, 332 (2001), pp. 661–668.

[27] Y. Maday and G. Turinici, *A parareal in time procedure for the control of partial differential equations*, Comptes Rendus Mathematique, 335 (2002), pp. 387–392.

[28] L. Mason, J. Baxter, P. L. Bartlett, and M. R. Frean, *Boosting algorithms as gradient descent*, in Advances in Neural Information Processing Systems, 2000, pp. 512–518.

[29] M. Minion, *A hybrid parareal spectral deferred corrections method*, Communications in Applied Mathe-

matics and Computational Science, 5 (2011), pp. 265–301.

[30] D. Narayanan, A. Harlap, A. Phanishayee, V. Seshadri, N. R. Devanur, G. R. Ganger, P. B. Gibbons, and M. Zaharia, *Pipedream: generalized pipeline parallelism for DNN training*, in Proceedings of the 27th ACM Symposium on Operating Systems Principles, 2019, pp. 1–15.

[31] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, *Reading digits in natural images with unsupervised feature learning*, NIPS Workshop on Deep Learning and Unsupervised Feature Learning, (2011).

[32] K. Simonyan and A. Zisserman, *Very deep convolutional networks for large-scale image recognition*, in International Conference on Learning Representations, 2015.

[33] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, *Going deeper with convolutions*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 1–9.

[34] D. H. Wolpert, *Stacked generalization*, Neural Networks, 5 (1992), pp. 241–259.

[35] B. Yang, J. Zhang, J. Li, C. Ré, C. R. Aberger, and C. De Sa, *Pipemare: Asynchronous pipeline parallel DNN training*, Oct. 2019, https://arxiv.org/abs/1910.05124.

[36] S. Zagoruyko and N. Komodakis, *Wide residual networks*, in Proceedings of the British Machine Vision Conference (BMVC), BMVA Press, September 2016, pp. 87.1–87.12.

[37] G. Zhang, C. Wang, B. Xu, and R. Grosse, *Three mechanisms of weight decay regularization*, in International Conference on Learning Representations, 2018.

[38] W. Zhang, S. Gupta, X. Lian, and J. Liu, *Staleness-aware async-SGD for distributed deep learning*, in Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, 2016, pp. 2350–2356.