

Optimizing Convolutional Neural Networks on Sunway TaihuLight Supercomputer

WENLAI ZHAO, HAOHUA FU, JIARUI FANG, WEIJIE ZHENG, LIN GAN, and GUANG-WEN YANG, Tsinghua University

Sunway TaihuLight supercomputer is powered by SW26010, a new 260-core processor designed with on-chip fusion of heterogeneous cores. In this paper, we present our work on optimizing the training process of convolutional neural networks (CNNs) on Sunway TaihuLight supercomputer. Specifically, a highly-efficient library (swDNN) and a customized Caffe framework (swCaffe) are proposed. Architecture-oriented optimization methods targeting the many-core architecture of SW26010 are introduced and are able to achieve 48 times speedup for the convolution routine in swDNN and 4 times speedup for the complete training process of VGG-16 network using swCaffe, compared with the unoptimized algorithm and framework. Compared with cuDNN library and the Caffe framework based on NVIDIA K40m GPU, the proposed swDNN library and swCaffe framework on SW26010 have nearly half the performance of K40m in single-precision, and have 3.6 times and 1.8 times speedup over K40m in double-precision, respectively.

CCS Concepts: • **Computing methodologies** → **Neural networks**; • **Computer systems organization** → **Multicore architectures**;

Additional Key Words and Phrases: Convolutional Neural Network, Deep Learning, Heterogeneous Many-core Architecture, Sunway TaihuLight Supercomputer

This manuscript is an *extension of conference paper*: "swDNN: A Library for Accelerating Deep Learning Applications on Sunway TaihuLight" published in IPDPS 2017 [10]. We consider this manuscript an improved edition of the conference paper with new contributions listed as follow:

- We present a more systemic algorithm design and optimization process with methods related to LDM usage, register communication, and instruction pipeline, including a modified performance model and a new register blocking strategy based on the previous work.
- We propose algorithm design and optimization methods to support single-precision on SW26010.
- We present the swDNN library with 4-CG parallelization to support different CNN layers.
- We propose swCaffe, an optimized Caffe framework that can support highly-efficient CNN training process on Sunway TaihuLight supercomputer.
- We present algorithm and framework evaluation with both float and double-precision for training practical CNN models, so that to provide more comprehensive performance results.

This work is supported in part by the National Key R&D Program of China (Grant No. 2016YFA0602200), by the National Natural Science Foundation of China (Grant No. 4137411, 91530323, 61702297, 61672312) and by the China Postdoctoral Science Foundation (No. 2016M601031).

Author address: W. Zhao, J. Fang, W. Zheng, Lin Gan, G. Yang, Department of Computer Science and Technology, Tsinghua University, Beijing 100084; H. Fu, Department of Earth System Science, Tsinghua University, Beijing 100084, China. All authors are concurrently with the National Supercomputing Center in Wuxi, Wuxi, 214000, Jiangsu Province, China.

Corresponding author: Haohuan Fu (Email address: haohuan@tsinghua.edu.cn).

Authors' address: Wenlai Zhao; Haohuan Fu; Jiarui Fang; Weijie Zheng; Lin Gan; Guangwen Yang, Tsinghua University, Computer Science and Technology.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Association for Computing Machinery.

XXXX-XXXX/2018/1-ART1 \$15.00

<https://doi.org/10.1145/3177885>

ACM Reference Format:

Wenlai Zhao, Haohuan Fu, Jiarui Fang, Weijie Zheng, Lin Gan, and Guangwen Yang. 2018. Optimizing Convolutional Neural Networks on Sunway TaihuLight Supercomputer. *ACM Transactions on Architecture and Code Optimization* 1, 1, Article 1 (January 2018), 25 pages. <https://doi.org/10.1145/3177885>

1 INTRODUCTION

Convolutional neural network (CNN [14]) is one of the most successful deep learning models in modern artificial intelligence applications [8, 12, 20, 22, 24]. The training process of CNN involves a large amount of computation, and has become a popular research topic in the field of high performance computing (HPC). GPUs have currently been considered as the most efficient hardware choice for deep learning tasks, and can support high level deep learning frameworks [1, 4, 7, 13].

Sunway TaihuLight [11], a supercomputer that ranks No.1 in the latest release (Nov. 2017) of the TOP 500 list with over 100 PFlops computing capacity, is powered by SW26010 many-core processor, which is designed with on-chip fusion of heterogeneous cores and is able to provide a peak double-precision performance of 3.06 TFlops. SW26010 introduces a number of unique features that could potentially accelerate the training process of CNNs, such as the user-controlled local directive memory (LDM), the hardware-supported register-level data sharing, and a unified memory space shared by all processing elements.

Our previous publication [10] has introduced the optimization of convolutional algorithm targeting the many-core architecture of SW26010. As an extension of the previous work, in this paper, we present a more systemic optimization process to accelerate CNN training tasks on Sunway TaihuLight supercomputer. Specifically, a highly-efficient library and a customized Caffe framework for SW26010 many-core processor are proposed.

The major contributions of this paper include:

- We propose algorithm design and optimization methods related to the LDM usage, register communication, and instruction pipeline, guided by a performance model. The optimized convolution routine can achieve 48 times speedup over the basic implementation on SW26010.
- A customized deep learning library for SW26010 many-core processor is developed, called swDNN, to provide the support for various computation and data processing layers in CNN models.
- An optimized Caffe framework for SW26010 many-core processor is proposed, called swCaffe, which is integrated with swDNN library, and supports 4-CG (core-group) parallelization on a SW26010 processor. The swCaffe framework can achieve about 4 times speedup over a BLAS-based Caffe framework on SW26010.

Evaluation results also show that the proposed convolution implementation and the swCaffe framework have nearly half the performance of NVIDIA K40m GPU in single-precision, while achieving 3.6 times and 1.8 times speedup over K40m in double-precision, respectively.

The paper is organized as follows. Section 2 introduces the background of the work, including the CNN algorithms, the detailed architecture of SW26010 and the related work on the optimization of CNN algorithms. Section 3 presents the performance model and architecture-oriented optimization methods targeting the convolution algorithm, including the evaluation of the implementation. Section 4 presents the swDNN library and the swCaffe framework, as well as the evaluation of the complete training process with swCaffe on a SW26010 many-core processor. Section 5 is the conclusion.

2 BACKGROUND

2.1 Convolutional Neural Networks

Table 1. Configurations of a convolutional layer

N_i	Number of input feature maps
R_i	Height of an input feature map
C_i	Width of an input feature map
N_o	Number of output feature maps
R_o	Height of an output feature map
C_o	Width of an output feature map
K	Size of convolution kernel

ALGORITHM 1: Original algorithm of a convolutional layer

```

1: //  $IN[B_s][N_i][R_i][C_i]$ ,  $OUT[B_s][N_o][R_o][C_o]$ ,  $CONVW[N_o][N_i][K_r][K_c]$  and  $b[N_o]$  are input/output
   feature maps, convolutional kernels and bias
2: //  $K_r = K_c = K$  represent the number of rows and columns of a 2-dimensional convolutional kernel
3: // The output images  $OUT$  are initialized with the bias  $b$ 
4: for  $cB := 0 : 1 : B_s$  do
5:   for  $cN_o := 0 : 1 : N_o$  do
6:     for  $cR_o := 0 : 1 : R_o$  do
7:       for  $cC_o := 0 : 1 : C_o$  do
8:         for  $cN_i := 0 : 1 : N_i$  do
9:           for  $cK_r := 0 : 1 : K_r$  do
10:            for  $cK_c := 0 : 1 : K_c$  do
11:               $OUT[cB][cN_o][cR_o][cC_o] += CONVW[cN_o][cN_i][K_r - 1 - cK_r][K_c - 1 - cK_c]$ 
                 $* IN[cB][cN_i][cR_o + cK_r][cC_o + cK_c];$ 
12:            end for
13:          end for
14:        end for
15:      end for
16:    end for
17:  end for
18: end for

```

CNNs usually contain multiple computing layers, among which *convolutional layers* usually account for the majority of the computing time (over 90%). We first give the description of the convolutional layer configurations, listed in Tab. 1. The input data of a convolutional layer consists of N_i channels, each of which can be considered as a feature map with size of $R_i \times C_i$. Similarly, the output of a convolutional layer consists of N_o feature maps with size of $R_o \times C_o$. To calculate the values in an output feature map, N_i convolutional kernels with size of $K \times K$ and 1 bias value are required. Each kernel convolutes with an input feature map. The output value equals to the sum of N_i convolution results and the bias value. Therefore, there are $N_i \times N_o$ convolutional kernels and N_o bias in a convolutional layer.

The training process of a CNN model is based on the stochastic gradient descent (SGD) algorithm. In each training step, the network is trained with a batch of samples. We define the batch size as B_s ,

and the original algorithm of a convolutional layer in a training iteration can be described as Alg. 1. The input data, output data and convolution weights are organized in 4-dimension tensors and there are 7 nested loops in the algorithm, which provides possibilities for the parallel optimization on many-core processors like SW26010.

Besides convolutional layers, a CNN usually contains other kinds of layers such as *pooling layers*, *fully-connected layers*, *softmax layers* and other data processing layers such as *activation function layers* and *normalization layers*. Different CNN models have different network structures, which describe how different kinds of layers are stacked in the neural network.

The major algorithm of fully-connected layers is matrix multiplication, which can be supported by the high-performance Basic Linear Algebra Subprograms (BLAS). Other layers like pooling, activation functions and softmax can be considered as data processing layers, and they are not the critical points of performance optimization.

2.2 SW26010 Many-core Architecture

Fig. 1 shows the architecture of a SW26010 many-core processor. A SW26010 consists of four core-groups (CGs) and each CG includes 65 cores: one management processing element (MPE), and 64 computing processing elements (CPEs) organized as an 8 by 8 mesh. The MPE and CPE are both complete 64-bit RISC cores but serve as different roles in a computing task.

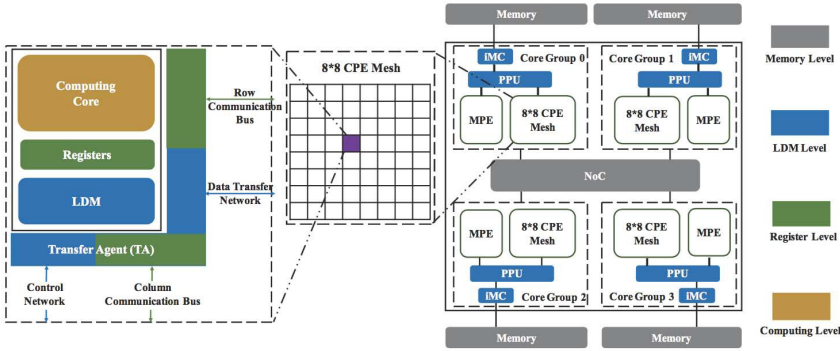


Fig. 1. SW26010 architecture

A MPE has a 32KB L1 instruction cache, a 32KB L1 data cache and a 256KB L2 cache, supporting the complete interrupt functions, memory management, superscalar, and out-of-order instruction issue/execution.

The CPE is designed for maximizing the aggregated computing throughput while minimizing the complexity of the micro-architecture. Each CPE has a 16KB L1 instruction cache and 64KB local directive memory (LDM). The LDM can be considered as a user-controlled fast buffer, which allows orchestrated memory usage strategies for different implementations, so that the LDM-level optimization is one of the important ways to improve the computation throughput.

A CPE has 32 vector registers (256 bits) and two execution pipelines (P0 and P1). P0 supports scalar and vectorized computing operations of both floating-point and integer, while P1 supports scalar and vectorized data load/store, compare, jump operations and scalar integer operations. The double pipelines provide opportunity for the overlapping of data accessing and computation operations. Therefore, register-level and instruction-level optimizations are also important to the performance.

Inside the 8×8 CPE mesh, there is a control network, a data transfer network (connecting the CPEs to the memory interface), 8 column communication buses, and 8 row communication buses. Each CPE has two 1024-bit send buffers and two 1024-bit receive buffers for column and row communication separately. The communication buses and buffers enable fast register-level data communication between CPEs of same column and same row, providing an important data sharing and cooperation capability within the CPE mesh.

In the instruction set, there are customized load/store instructions to support both vectorized data access and data sharing in a non-blocking mode. For example, a *vldr* instruction first loads 256-bit data into a vector register and then performs the **row broadcast**; and a *vlddec* instruction first loads 64-bit data into a scalar register, then extends (copies) the data to fill a vector register, and finally performs the **column broadcast**. Based on these instructions, highly efficient data access and register communication can be realized.

Each CG connects to a Memory Controller (MC), through which 8 GB memory space can be accessed and shared by the MPE and the CPE mesh. The maximum memory bandwidth of an MC is 36GB/s. An on-chip network (NoC) connects four CGs, so that the memory of a CG can also be shared to other CGs. Users can explicitly set the size of each CG's private memory space, and the size of the shared memory space. Through NoC, data sharing between four CGs can be implemented without memory data copy, which enables highly efficient CG-level parallelism for communication-intensive problems. Under the sharing mode, the maximum memory bandwidth of four CGs is up to 144GB/s.

2.3 Related Works

A straightforward implementation of the original convolution algorithm involves strong data dependency in the innermost accumulation computation. To improve the parallelism, several optimization methods are proposed, which can be summarized into the following three categories.

- **Time-domain transformation methods** are first introduced in the early phase of CNN optimization researches [2, 6, 15]. By expanding convolution operations into matrix multiplications, the performance can be improved with the help of the BLAS on different hardware platforms. However, additional data transformation is required, which either consumes more memory space and extra data copy operations, or involves complicated memory address remapping. Therefore, the memory consumption and bandwidth are major problems for time-domain transformation methods, and the overall performance is limited by the performance of BLAS.
- **Frequency-domain transformation methods** can reduce the arithmetic complexity of convolution operations. FFT-based [18, 23] and Winograd's filtering [16] based convolution algorithms are proposed and perform well in the cases with both large and small convolution kernel sizes. Similar to time-domain based methods, additional data transformation, as well as extra memory consumption, is required, and the overall performance is limited by the performance of transformation.
- **Direct convolution optimization methods** can reduce the data dependency by re-designing the convolution algorithm with loop reordering and data blocking, so as to improve the parallelism of the core computation. Instead of relying on existing BLAS or FFT libraries, direct convolution implementations require hardware-oriented optimization methods to take full advantage of the hardware architecture, and therefore, the overall performance can approach to the peak performance of the processor. Moreover, by carefully designing the data blocking strategies, additional data transformation and extra memory consumption can be avoided, which is more suitable for memory and bandwidth bounded architectures.

Besides the algorithm optimization, various hardware accelerators are employed to accelerate the convolution computation, such as GPU, FPGA and ASIC, focusing on both classification and training process of CNNs. FPGAs [19, 25–27] and ASICs [3, 5, 9, 17] are usually used for classification tasks due to the customizability of data precision, low latency and high energy efficiency. GPUs have currently dominated the competition of the HPC platforms for training tasks. Especially, NVIDIA has launched GPU like V100, which includes deep learning specific units like tensor cores. Correspondingly, the cuDNN [6] library is released to provide highly-efficient routines for deep learning algorithms on NVIDIA GPUs, and can be neatly integrated to widely used deep learning frameworks, such as Caffe [13], Tensorflow [1], etc.

To explore the potential of training CNNs on other off-the-shelf many-core processors, in this paper, we present the detailed architecture-oriented optimization methods for the convolution algorithm on the SW26010 processor. Then, we show the design of the deep learning library and the customized Caffe framework dedicated for the SW26010 processor, so as to support a highly-efficient training process of CNN on Sunway TaihuLight supercomputer.

3 CONVOLUTION ALGORITHM OPTIMIZATION

We first introduce a performance model that shows the features of the SW26010 architecture and indicates the key factors that could affect the performance of an implementation. Guided by the performance model, we re-design the convolution algorithm and propose LDM-related, register-related and instruction-related methods for further optimizations.

3.1 Performance model

We consider different factors that affect the performance of one CG and propose a performance model shown in Fig. 2. The frequency of a CPE is 1.45GHz and the vectorization size is 4. Assuming that each CPE executes one vector floating-point multiplication and addition (*vmad*) instruction, the peak performance of a CG can be derived as:

$$2 \times 4 \times 1.45 \times 64 = 742.4GFlops \quad (1)$$

For an implementation, we define the *execution efficiency* (EE) as the ratio of *vmad* instructions to the total execution cycles. Therefore, considering the loss from EE, the theoretical performance of an implementation is $742.4GFlops \cdot EE$.

Before a computing instruction can be executed, we need to make sure the data has been loaded into registers. For a *vmad* instruction, 12 double-precision numbers ($12 \times 64 = 768bits$) are needed. In Fig. 2, the *required bandwidth* (RBW) of an implementation is defined as the minimum data access bandwidth that could overlap the data access and computation.

A CPE supports two data access patterns to load the data into registers. One is the global memory access (*gload* instruction), which can load 64 bits data into a scalar register directly from main memory. In this case, to guarantee the overlapping of computation and data access, the data accessed by a *gload* instruction should be involved in at least 12 *vmad* instructions ($768bits : 64bits$). Here we define the *computation to data access ratio* (CDR), which represents the ratio of computation instructions (*vmad*) to data access instructions. In the global memory access pattern, in order to overlap the computation and data access, CDR should be greater than 12, which can hardly be met by most algorithms. Therefore, the global memory access pattern is relatively low efficient.

The performance model of the global memory access pattern is shown in Fig. 2. The maximum memory bandwidth of one CG is about 8GB/s. We denote the RBW by $RBW_{MEM \rightarrow REG}$. Here we assume that the computation and the data access are parallel processes and are independent, which can be realized through some optimization methods like double buffering (see Section 3.3.2).

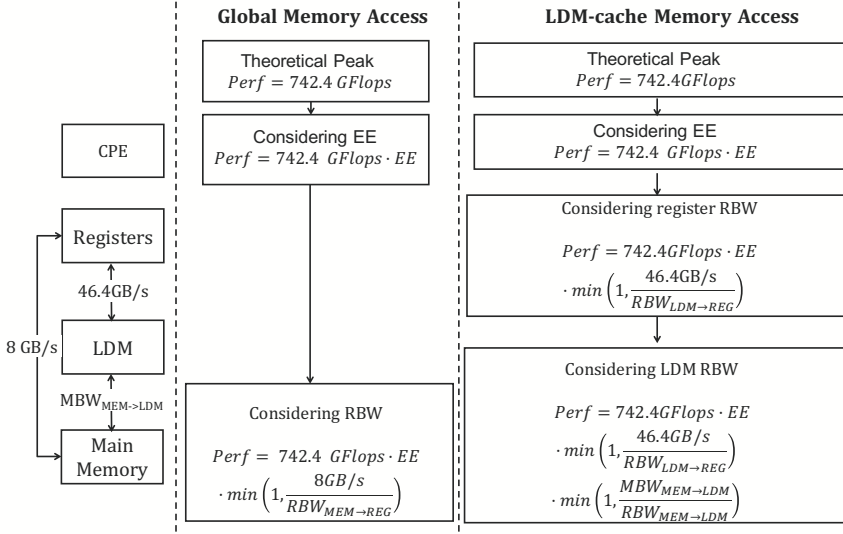


Fig. 2. Performance model for one CG (EE: Execution Efficiency, RBW: Required Bandwidth, MBW: Measured Bandwidth)

Therefore, if the $RBW_{MEM \rightarrow REG}$ is greater than 8GB/s, it will lower the performance by a rate of $\frac{8GB/s}{RBW_{MEM \rightarrow REG}}$.

The other memory access pattern is to use the LDM as a data cache, which means the data will be loaded first from the main memory into the LDM, and then from the LDM into registers. There are two stages of data accessing in this case. We denote the RBW of the two stages by $RBW_{MEM \rightarrow LDM}$ and $RBW_{LDM \rightarrow REG}$. When loading data from the LDM to registers, vectorized load instruction (*vload*) is supported. Each *vload* instruction can load 256-bit (32Bytes) data into a vector register. The execution of load/store instructions usually takes 3 or 4 CPU cycles and is a non-blocking process, so that we can issue an instruction every cycle and the bandwidth between the LDM and registers is $32Bytes \times 1.45GHz = 46.4GB/s$.

Data is transferred from main memory to LDM through the direct memory access interface (DMA), and the theoretical maximum bandwidth of DMA is 36GB/s. A DMA put/get operation will access one or more *memory blocks*, which has a size of 128 Bytes for SW26010. The latency of a DMA request from CPEs is usually more than 100 CPU cycles. Therefore, theoretically, successive DMA operations with large granularity can make full use of the DMA bandwidth. Practically, the actual bandwidth is not a constant value and is variant with the size of continuous memory access blocks of one CPE. We write a micro-benchmark on one CG to measure the actual DMA bandwidth and present the results in Tab. 2, where *Size* indicates the granularity of a DMA operation. We denote the measured DMA bandwidth (MBW) by $MBW_{MEM \rightarrow LDM}$. We can see that the bandwidth of DMA ranges from 4 GB/s to 36 GB/s. In general, a higher bandwidth is achieved when using a block size larger than 256 Bytes and aligned in 128 Bytes.

Fig. 2 also shows the performance model of the LDM-cache memory access pattern. Here the required CDR is 3 (768bits : 256bits), which is easier to be accomplished compared with the global memory access pattern. Our design is based on the LDM-cache memory access pattern. According to the performance model, we propose optimization methods to overlap the computation and data access, to increase the $MBW_{MEM \rightarrow LDM}$, EE , and to reduce the $RBW_{MEM \rightarrow LDM}$ and $RBW_{LDM \rightarrow REG}$.

Table 2. Measured DMA Bandwidth on one CG(GB/s)

Size(Byte)	Get	Put	Size(Byte)	Get	Put
32	4.31	2.56	512	27.42	30.34
64	9.00	9.20	576	25.96	28.91
128	17.25	18.83	640	29.05	32.00
192	17.94	19.82	1024	29.79	33.44
256	22.44	25.80	2048	31.32	35.19
384	22.88	24.67	4096	32.05	36.01

3.2 Algorithm Design

Considering the original algorithm of a convolutional layer (Alg. 1), the inner loops perform a $K \times K$ convolution. Usually, the value of K is relatively small and is odd, like 3, 5, 7, etc. Therefore, it is hard to map the inner loops onto the CPE mesh and is also inefficient for the vectorization of core computation.

To improve the parallelism, we re-schedule the 7 nested loops, making the inner computation to be a matrix multiplication with dimensions N_i , N_o and B_s , which are relatively large in most convolution layers and are suitable for mapping the inner computation onto the CPE mesh. Alg. 2 shows the optimized algorithm based on matrix multiplication. We call the inner matrix multiplication operation as the *core computation*. To complete the computation of an output matrix (D_o) of size $N_o \times B_s$, each CPE is responsible for a block of size $\frac{N_o}{8} \times \frac{B_s}{8}$. Correspondingly, the input data of a CPE includes a tile of the input matrix W (of size $\frac{N_o}{8} \times N_i$) and a tile of the input matrix D_i (of size $N_i \times \frac{B_s}{8}$), both of which can be shared between the CPEs either in the same row or in the same column. Therefore, for the core computation, the amount of data to be accessed by a CPE is $(N_i \times \frac{N_o}{8} + N_i \times \frac{B_s}{8} + \frac{N_o}{8} \times \frac{B_s}{8})$. The amount of *vmadd* instructions is $(N_i \times \frac{N_o}{8} \times \frac{B_s}{8})/4$. We use *vload* instruction for data access, so the theoretical CDR of the core computation is:

$$\frac{(N_i \times \frac{N_o}{8} \times \frac{B_s}{8})/4}{(N_i \times \frac{N_o}{8} + N_i \times \frac{B_s}{8} + \frac{N_o}{8} \times \frac{B_s}{8})/4} \quad (2)$$

Assuming that N_i , N_o , and B_s have the same value, the CDR can meet the requirement ($CDR \geq 3$) of LDM-cache pattern when the value is larger than 51, which can be realized in most of the convolution layers. For values which are not multiple of 8, zero padding can be adopted and will not cause too much decrease on the performance. Therefore, for brevity, we focus on the configurations that are multiple of 8 in the following discussion. The following subsections will show the detailed implementation and optimization methods based on Alg. 2.

3.3 LDM-Related Optimization

LDM-related optimization methods are focused on an effective implementation for outer loops of the algorithm. The targets are to realize the overlap of data access from main memory to LDM and the core computation of the CPE mesh, so as to increase $MBW_{MEM \rightarrow LDM}$ and reduce $RBW_{MEM \rightarrow LDM}$.

3.3.1 Optimized Data Layout. The input data of the core computation is a part of the input/output feature maps and the convolutional kernels. Based on the original data layout, data in W , D_i , D_o is not stored continuously in IN , OUT and $CONVW$, so that the $MBW_{MEM \rightarrow LDM}$ will be limited due to small data access block. To increase $MBW_{MEM \rightarrow LDM}$, we re-design the data layout of the input/output feature maps and the convolutional kernels as $IN[R_i][C_i][N_i][B_s]$,

ALGORITHM 2: Matrix-multiplication-based convolution algorithm

```

1: //  $IN[B_s][N_i][R_i][C_i]$ ,  $OUT[B_s][N_o][R_o][C_o]$ ,  $CONVW[N_o][N_i][K_r][K_c]$  and  $b[N_o]$  are input/output
   feature maps, convolutional kernels and bias
2: //  $K_r = K_c = K$  represent the number of rows and columns of a 2-dimensional convolutional kernel
3: // The output images  $OUT$  are initialized with the bias  $b$ 
4: for  $cR_o := 0 : 1 : R_o$  do
5:   for  $cC_o := 0 : 1 : C_o$  do
6:      $D_o[0 : N_o][0 : B_s] = (OUT[0 : B_s][0 : N_o][cR_o][cC_o])^T$ 
7:     for  $cK_r := 0 : 1 : K_r$  do
8:       for  $cK_c := 0 : 1 : K_c$  do
9:          $W[0 : N_o][0 : N_i] = CONVW[0 : N_o][0 : N_i][K - 1 - cK_r][K - 1 - cK_c]$ 
10:         $D_i[0 : N_i][0 : B_s] = (IN[0 : B_s][0 : N_i][cR_o + cK_r][cC_o + cK_c])^T$ 
11:        Core computation:  $D_o += W \times D_i$ 
12:      end for
13:    end for
14:     $OUT[0 : B_s][0 : N_o][cR_o][cC_o] = (D_o[0 : N_o][0 : B_s])^T$ 
15:  end for
16: end for

```

$OUT[R_o][C_o][N_o][B_s]$, and $CONVW[K_r][K_c][N_o][N_i]$. Besides, we rotate the convolutional kernels on K_r and K_c dimensions to eliminate the coordinate transform in Line 6 of Alg. 2. For IN and OUT , we put B_s as the lowest dimension, which can eliminate the data transposition in Line 3, 7 and 11 of Alg. 2, and can support vectorized operations on B_s dimension in the core computation.

3.3.2 Double Buffering. Double buffering is adopted to overlap the data access from main memory to LDM and the core computation. Because DMA is asynchronous, we design two LDM buffers of the same size. While the data in one buffer is used for core computation, the data to be used in next core computation can be loaded into another buffer. Note that the double buffering design halves the maximum available space of LDM for one computation iteration, which means for one CPE, only 32KB LDM is available for the core computation.

3.3.3 LDM Blocking. We consider the total LDM usage of 64 CPEs in the core computation with different convolutional layer configurations. It can be described as:

$$(N_i \times N_o + N_i \times B_s + N_o \times B_s) \times DataLen \quad (3)$$

where $DataLen$ is the number of bytes for the data type. Assuming N_i , N_o , and B_s are equal to 256, which are relatively large configurations for most convolutional layers, and the data type is double-precision, the total LDM usage of 64 CPEs is $3 \times 256 \times 256 \times 8Bytes = 1536KBytes$. By using register communication techniques, the data stored in one CPE's LDM can be shared to other CPEs (more details will be shown in Section 3.4.1), so that the exact LDM usage of each CPE is $1536KB/64 = 24KB$. Therefore, for most convolutional layers, 32KB LDM is enough for the core computation, and in other words, it is possible to take advantage of the remaining LDM spaces to improve the overall performance of the implementation.

In the convolution algorithm, the convolutional kernel is shared by the computation of values in the same output image. In the core computation of Alg. 2, the data of convolutional kernel (W) is only used for one core computation corresponding to the values in the output feature maps at coordinate (cR_o , cC_o). To improve the data reuse of W , and in the meantime to improve the CDR of the core computation, we propose an LDM blocking strategy shown in Alg. 3.

ALGORITHM 3: Optimized algorithm with LDM blocking

```

1: //IN[Ri][Ci][Ni][Bs], OUT[Ro][Co][No][Bs], CONVW[Kr][Kc][No][Ni] and b[No] are input/output
   feature maps, convolutional kernels and bias
2: //Kr = Kc = K represent the number of rows and columns of a 2-dimensional convolutional kernel
3: //W,  $\tilde{W}$  and Di,  $\tilde{D}_i$  represent the double buffering for weight and input feature maps
4: //The output images OUT are initialized with the bias b
5: for cRo := 0 : 1 : Ro do
6:   for cCo := 0 : bC : Co do
7:     DMA get Do[0 : bC][0 : No][0 : Bs] ← OUT[cRo][cCo : cCo + bC][0 : No][0 : Bs]
8:     for cKr := 0 : 1 : Kr do
9:       for cKc := 0 : 1 : Kc do
10:        DMA get:
11:         $\tilde{W}[0 : N_o][0 : N_i] \leftarrow CONVW[cK_r][cK_c][0 : N_o][0 : N_i]$ 
12:         $\tilde{D}_i[0 : b_C][0 : N_i][0 : B_s] \leftarrow IN[cR_o + cK_r][cC_o + cK_c : cC_o + cK_r + b_C][0 : N_i][0 : B_s]$ 
13:        //Core computation:
14:        for cbC := 0 : 1 : bC do
15:          Do[cbC] += W × Di[cbC]
16:        end for
17:        Check DMA get  $\tilde{W}, \tilde{D}_i$  finished.
18:        Exchange W, Di with  $\tilde{W}, \tilde{D}_i$ 
19:      end for
20:    end for
21:    DMA put Do[0 : bC][0 : No][0 : Bs] → OUT[cRo][cCo : cCo + bC][0 : No][0 : Bs]
22:  end for
23: end for

```

In the core computation of Alg. 3, we load b_C times more data of input/output feature maps and reuse the data of convolutional kernels to complete b_C matrix-multiplication computation. The $RBW_{MEM \rightarrow LDM}$ is reduced and the CDR of a CPE is:

$$\frac{b_C \times N_i \times \frac{N_o}{8} \times \frac{B_s}{8} / 4}{(N_i \times \frac{N_o}{8} + b_C \times N_i \times \frac{B_s}{8} + b_C \times \frac{N_o}{8} \times \frac{B_s}{8}) / 4} \quad (4)$$

which is greater than Equation (2). The larger b_C we choose, the greater CDR we can get. On the other hand, b_C is limited by the available size of LDM, and we can maximize the value to take fully advantage of the LDM.

3.4 Register-Related Optimization

Register-related optimization methods mainly focus on effectively mapping the core computation onto 8×8 CPE mesh. Two key problems are targeted in our work: (i) to realize the register-level data sharing between CPEs, so that to reduce the $RBW_{LDM \rightarrow REG}$ for each CPE; (ii) to take fully use of the vector register to implement the computation efficiently on a CPE.

3.4.1 Register Communication. In the core computation, a CPE is responsible for a $\frac{N_o}{8} \times \frac{B_s}{8}$ block of D_o , that requires an $\frac{N_o}{8} \times N_i$ tile of W and an $N_i \times \frac{B_s}{8}$ tile of N_i . CPEs in the same row of the mesh share the tile of W , and CPEs in the same row of the mesh share the tile of N_i , which perfectly matches the register communication feature of the CPE mesh. However, there are some limitations of the register communication feature: (i) the send and receive buffers designed for the register communication are simply FIFOs with limited size ($4 \times 256bits$); (ii) the data received

though the register communication buses has no information of the source CPE; (iii) if the send and receive buffer are both full, the source CPE will halt.

Considering the limitations, we carefully design a register communication strategy for the matrix multiplication computation. For simplicity, we take a 4×4 CPE mesh as an example to introduce the design, shown in Fig. 3. We label the CPEs with coordinates $(0, 0) - (3, 3)$ from top left to bottom right. D_i , W and D_o are divided into 4×4 parts, and are labeled as $D_i(0, 0) - D_i(3, 3)$, $W(0, 0) - W(3, 3)$ and $D_o(0, 0) - D_o(3, 3)$. For a given pair of (i, j) , the computation of $D_o(i, j)$ can be described as:

$$D_o(i, j) += \sum_{k=0}^3 W(i, k) \times D_i(k, j) \quad (5)$$

which can be done in 4 steps by CPE(i, j). $D_i(i, j)$, $W(i, j)$ and $D_o(i, j)$ are pre-loaded into the LDM of CPE(i, j) before executing the core computation. Without loss of generality, we take CPE(2, 1) as an example to show the process.

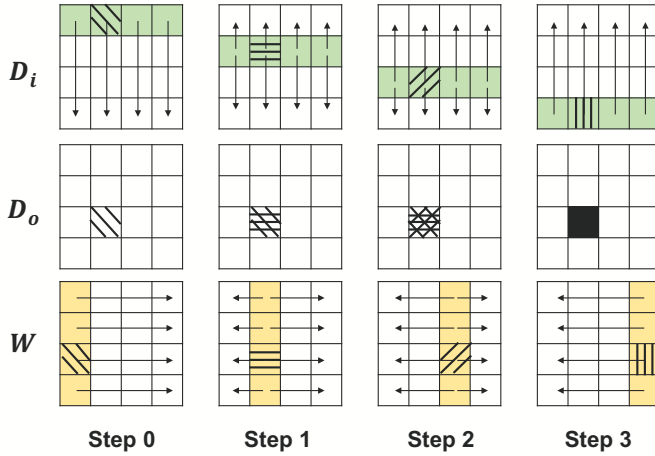


Fig. 3. Register communication example on 4×4 CPE mesh

- **Step 0** First, for all $j \in \{0, 1, 2, 3\}$, CPE($0, j$) loads data of $D_i(0, j)$ from LDM and sends the data to other CPEs in the same column by register communication. Thus, CPE(2, 1) can receive the data of $D_i(0, 1)$. Then, for all $i \in \{0, 1, 2, 3\}$, CPE($i, 0$) loads data of $W(i, 0)$ from LDM and sends the data to CPEs in the same row. CPE(2, 1) can receive the data of $W(2, 0)$. $D_o(2, 1)$ can be loaded from the LDM of CPE(2, 1), so that the computation of $D_o(2, 1) += W(2, 0) \times D_i(0, 1)$ can be done.
- **Step 1** First, CPEs with coordinates $(1, j)$ load data of $D_i(1, j)$ from LDM and send the data to CPEs in the same column. Then, CPEs with coordinates $(i, 1)$ load data of $W(i, 1)$ and send CPEs in the same row. Thus, CPE(2, 1) can receive the data of $D_i(1, 1)$ through column register communication, and can load $W(2, 1)$ and $D_o(2, 1)$ from LDM, so that to compute $D_o(2, 1) += W(2, 1) \times D_i(1, 1)$.
- **Step 2** CPEs with coordinates $(2, j)$ and $(i, 2)$ load the data of $D_i(2, j)$ and $W(i, 2)$, and send to the same column and same row respectively. Then, CPE(2, 1) can receive the data of $W(2, 2)$ through row register communication and load $W(2, 2)$ and $D_o(2, 1)$ from LDM. The computation of $D_o(2, 1) += W(2, 2) \times D_i(2, 1)$ can be done.

- **Step 3** Similarly, CPEs with coordinates $(3, j)$ and $(i, 3)$ load and send the data of $D_i(3, j)$ and $W(i, 3)$ respectively. Correspondingly, CPE(2, 1) can receive $W(2, 3)$ and $D_i(3, 1)$ through row and column register communication, and finally finish the computation of $D_o(2, 1) += W(2, 3) \times D_i(3, 1)$.

Based on the proposed register communication strategy, the core computation can be done on 8×8 CPE mesh following 8 steps with a highly efficient data sharing between CPEs.

3.4.2 Register Blocking. In each step of the register communication process, the computation task of a CPE is to calculate the matrix multiplication of $W(i, j)$ and $D_i(i, j)$. The size of the blocks are $(\frac{N_o}{8} \times \frac{N_i}{8})$ and $(\frac{N_i}{8} \times \frac{B_s}{8})$ respectively.

For each CPE, there are only 32 vector registers, including zero register and stack pointer (*sp*) register, so that the number of available registers is less than 30 for the implementation. We should consider to use vectorized computation, to improve the data reuse in registers, and to reduce the data dependency in order to achieve an efficient instruction flow. Therefore, we propose a register blocking strategy to implement the computation in each step. Fig. 4 shows the details.

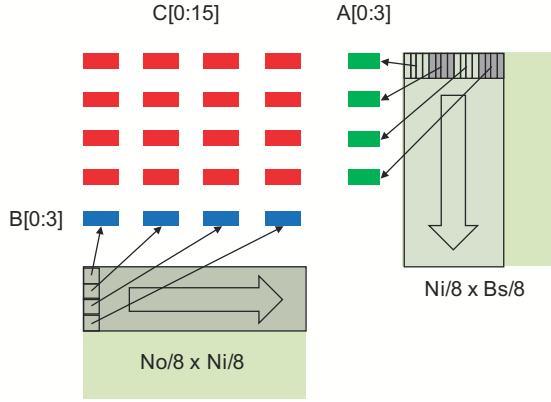


Fig. 4. Register blocking strategy on one CPE

We use 4 vector registers to load D_i , denoted by $A[0 : 3]$, and 4 vector registers to load W , denoted by $B[0 : 3]$. 16 vector registers are used for storing the data of D_o , denoted by $C[0 : 15]$. We define the following process as a **kernel task** of the register blocking design:

- First, we load 16 values in a row of $D(i, j)$ into $A[0:3]$, which can be done by 4 *vload* instructions. We load 4 values in a column of $W(i, j)$ and duplicate the values to fill $B[0:3]$, which can be done by 4 *vlde* instructions.
- Second, we load 4×16 values from $D_o(i, j)$ into $C[0:15]$ using 16 *vload* instructions.
- Third, for $i, j \in \{0, 1, 2, 3\}$, we calculate Equation (6) using 16 *vmad* instructions.

$$C[i + 4 * j] += A[i] \times B[j] \quad (6)$$

24 registers are used in the kernel task. As we can see from Fig. 4, to finish the calculation of 4×16 values of $D_o(i, j)$, $\frac{N_i}{8}$ kernel tasks are required. During this process, $A[0:3]$ and $B[0:3]$ are reloaded for $\frac{N_i}{8}$ times while $C[0:15]$ only need to be loaded once in the first kernel task, which improves the data reuse at register level, and thus, reduces the $RBW_{LDM \rightarrow REG}$. Because there is no data dependency between the *vmad* instructions in a kernel task, one instruction can be issued in each CPU cycle, which can increase the *EE* of the implementation.

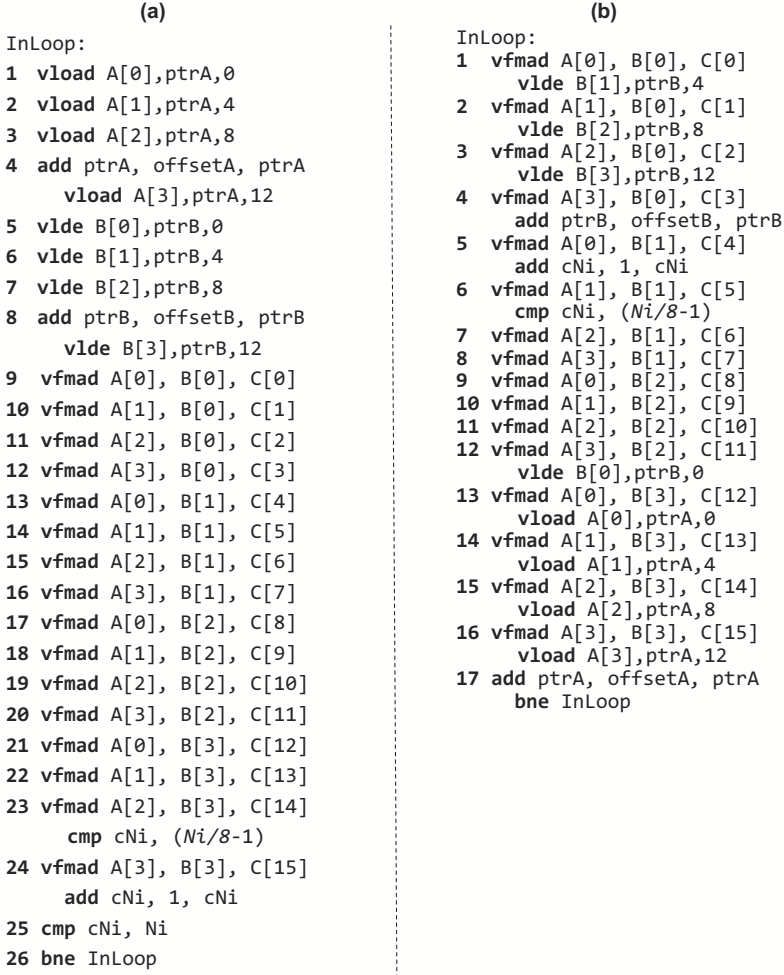


Fig. 5. Instruction-related optimization for the kernel task

3.5 Instruction-Related Optimization

We adopt instruction-related optimization methods to overlap the data loading and computation instructions and to further improve the *EE* in the kernel task. Fig. 5(a) shows the instruction flow based on a direct implementation of the kernel task. It takes 26 CPU cycles to issue the instructions, among which, there are 16 *vfmad* instructions. The *EE* is $16/26 = 61.5\%$. As we can see, in cycle 4, 8, 23 and 24, two instructions can be issued to pipeline P0 and P1 simultaneously, because there is no data dependency and the instructions can be executed on P0 and P1 separately. Only data loading instructions (*vldr* can load the data into a vector register and send out through row register communication) are issued in the first few cycles, which will lower the *EE* of the implementation.

Considering that $\frac{N_i}{8}$ kernel tasks are required to calculate a 4×16 block of $D_o(i, j)$, we unroll the $\frac{N_i}{8}$ kernel tasks and reorder the instructions to overlap the *vldr* instructions of a kernel task with the *vfmad* instructions at the end of the previous kernel task. The implementation after loop

unrolling and instructions reordering is shown in Fig. 5(b), where only 17 CPU cycles are required to finish a kernel task and the *EE* is improved to $16/17 = 94.1\%$.

3.6 Core-Group Level Parallel Scheme

Based on the above optimization methods, the convolution algorithms can be mapped onto a CG efficiently. Considering that there are 4 CGs in a SW26010 processor, we can further design the parallel scheme for 4 CGs. The simplest but most efficient way is to introduce parallelism on the outermost loop (R_o). As discussed in Section 2.2, data can be shared by 4 CGs without extra data copy. Therefore, we can set the data of input/output feature map, convolutional kernel and bias to the shared mode, and implement a four-CG convolution algorithm as shown in Alg. 4.

ALGORITHM 4: 4-CG implementation of convolution algorithm

```

1: //Assume that  $IN[R_i][C_i][N_i][B_s]$ ,  $OUT[R_o][C_o][N_o][B_s]$ ,  $CONVW[K_r][K_c][N_o][N_i]$  and  $b[N_o]$  are
   input/output feature maps, convolutional kernels and bias
2: //  $K_r = K_c = K$  represent the number of rows and columns of a 2-dimensional convolutional kernel
3: //  $W, \tilde{W}$  and  $D_i, \tilde{D}_i$  represent the double buffering for weight and input feature maps
4: //The output images  $OUT$  are initialized with the bias  $b$ 
5: //Parallel execution on 4 CGs
6: for  $cg := 0 : 1 : 4$  do
7:   for  $cR_o := 0 : 1 : \frac{R_o}{4}$  do
8:     for  $cC_o := 0 : b_C : C_o$  do
9:       DMA get  $D_o[0 : b_C][0 : N_o][0 : B_s] \leftarrow OUT[cg \times \frac{R_o}{4} + R_o][cC_o : cC_o + b_C][0 : N_o][0 : B_s]$ 
10:      for  $cK_r := 0 : 1 : K_r$  do
11:        for  $cK_c := 0 : 1 : K_c$  do
12:          DMA get:
13:           $\tilde{W}[0 : N_o][0 : N_i] \leftarrow CONVW[cK_r][cK_c][0 : N_o][0 : N_i]$ 
14:           $\tilde{D}_i[0 : b_C][0 : N_i][0 : B_s] \leftarrow IN[cg \times \frac{R_o}{4} + cR_o + cK_r][cC_o + cK_c : cC_o + cK_r + b_C][0 : N_o][0 : B_s]$ 
15:          for  $cb_C := 0 : 1 : b_C$  do
16:            Core computation:  $D_o[cb_C] += W \times D_i[cb_C]$ 
17:          end for
18:          Check DMA get  $\tilde{W}, \tilde{D}_i$  finished.
19:          Exchange  $W, D_i$  with  $\tilde{W}, \tilde{D}_i$ 
20:        end for
21:      end for
22:    end for
23:     $OUT[cg \times \frac{R_o}{4}][cC_o : cC_o + b_C][0 : N_o][0 : B_s] = D_o[0 : b_C][0 : N_o][0 : B_s]$ 
24:  end for
25: end for

```

3.7 single-precision Support

During the above design and optimization process, we consider double-precision (64-bit) as the data representation for both feature maps and weights. However, unlike in most scientific applications, single-precision (32-bit) is sufficient for training CNN models in deep learning applications. Therefore, to support practical CNN training tasks, we further improve our optimized algorithm implementation to support single-precision operations.

Originally designed for supporting major scientific applications that mostly rely on double-precision data types, the features of SW26010 architecture, such as vectorized instructions and

register communication operations, are generally more suitable to handling double-precision operations. There is no special optimization on hardware for single-precision operations on SW26010. Therefore, theoretically, the peak performance for single-precision computation is equal to that for double-precision. In practice, there will be performance loss due to the lack of support for single-precision operation in the instruction set, which will be discussed below.

A straightforward way to support single-precision is to re-design the kernel task instruction flow based on single-precision instructions. The major problem is that there is no instruction like *vldr* or *vlddec* for single-precision data in the instruction set of SW26010. Instead, we should first load 4 single-precision data into a vector register using *vlds* or *vldse* instruction, and then call register communication using *putr* or *putc* instruction. Therefore, for single-precision, the instruction flow of the kernel task has 8 more instructions than the double-precision implementation shown in Fig. 5, and, more importantly, these instructions can not be overlapped by computation instructions due to the register dependency. 8 more cycles in the kernel task will lower the *EE* to $16/(17 + 8) = 64\%$, which indicates the overall performance loss will be more than 30% (compared to 94.1%).

In the straightforward way, all data accessed in the kernel task requires extra cycles. Considering that there is data reused in the core computation, for example W will be reused for cb_C times, we propose another way to reduce the overall extra cycles for single-precision data access, called a *float2double* implementation. After we load the data from the main memory to LDM, we cast the single-precision data in the LDM to double-precision and then do the core computation in double-precision. Correspondingly, we cast the double-precision data to single-precision before storing the computation results back to the main memory. The data casting can be implemented using a flow of *vlds/vsts* (for single-precision) and *vldd/vstd* (for double-precision) instructions.

3.8 Evaluation

To show the performance improvement obtained from the proposed algorithm design and optimization methods, we first evaluate the performance of the implementation based on double-precision.

Different convolutional layer configurations listed in Table 1 will lead to different practical performance. Since the configurations change with CNN models and applications irregularly, it

ALGORITHM 5: Configuration generation algorithm

```

1: Test Set 1 :  $B_s = 128, R_o = C_o = 32, K = 3$ 
2: for  $N_o = 64; N_o \leq 384; N_o += 64$  do
3:   for  $N_i = 64; N_i \leq 384; N_i += 64$  do
4:      $CONV(B_s, N_i, N_o, R_o, C_o, K)$ 
5:   end for
6: end for
7: Test Set 2 :  $B_s = 128, N_i = N_o = 128, K = 3$ 
8: for  $R_o = C_o = 8; R_o \leq 128; R_o += R_o, C_o += C_o$  do
9:    $CONV(B_s, N_i, N_o, R_o, C_o, K)$ 
10: end for
11: Test Set 3 :  $B_s = 128, N_i = N_o = 128, R_o = C_o = 64$ 
12: for  $K = 3; K \leq 11; K += 2$  do
13:    $CONV(B_s, N_i, N_o, R_o, C_o, K)$ 
14: end for
15: Test Set 4 :  $N_i = N_o = 128, R_o = C_o = 64, K = 128$ 
16: for  $B_s = 32; B_s \leq 512; B_s *= 2$  do
17:    $CONV(B_s, N_i, N_o, R_o, C_o, K)$ 
18: end for

```

Table 3. Specifications of SW26010 and K40m/K80m GPU

Specifications		SW26010	NVIDIA K40m	NVIDIA K80m
Release Year		2014	2013	2014
TDP		250W	235W	375W
Number of Cores		260	2880 (15 SM ¹)	4992 (26 SM)
Memory	Capacity	32GB	12GB	24GB
	Bandwidth	144GB/s	288GB/s	480GB/s
Peak Perf.	float	3.02TFlops	4.29TFlops	8.74TFlops
	double	3.02TFlops	1.43TFlops	2.91TFlops

¹ Each SM (Streaming Multiprocessor) has 192 CUDA cores.

is unnecessary to traverse all possibilities. Therefore, we derive the test cases according to Alg. 5, where 4 sets of test cases are generated targeting different values of N_i/N_o , $Ro(Co)$ K and B_s separately.

Tab. 3 lists the specifications of SW26010 and NVIDIA K40/K80 GPUs. Taking the peak performance in both single and double-precision into consideration, we choose K40m GPU as a comparison to SW26010 in our evaluation. We run the test cases using our implementation and the convolution subroutine of cuDNN (v5.1) on NVIDIA K40m GPU. The evaluation results are summarized into four categories to show that how the performance changes with different configurations as shown in Fig. 6 and Fig. 7.

As we can see from Fig. 6(a), the performance of our implementation is more sensitive to the value of N_i . As discussed in Section 3.4.2, in each step of the register communication process, $\frac{N_i}{8}$ kernel tasks are executed. Therefore, larger N_i will lead to a longer process with consecutive kernel tasks, so that can provide better performance. Fig. 7(a) shows that the performance with small value of R_o (and C_o) is relatively low, which is because we use double buffering design to achieve the overlap of the data access from main memory to LDM and the core computation. The design can be considered as a pipeline and there is a starting phase at the beginning of the process. Small R_o and C_o will shorten the pipeline, and therefore, lower the overall performance. The performance with different N_o and different K is relatively stable according to Fig. 6(b) and Fig 7(b).

In Fig. 7(c), small B_s (such as 32 and 64) can not take fully use of the 8×8 CPE mesh, so that the performance penalty of the proposed implementation is quite apparent. For large B_s (such as 256 and 512) the performance improvement is also apparent since large B_s will benefit the performance of the innermost matrix multiplication computation in our design.

Considering all test cases, the performance of our implementation ranges from 1.3TFlops to 2.0TFlops, and the average performance is about 1.68TFlops, which is about 56% of the peak performance of SW26010. For the evaluation of cuDNN on K40m GPU, the average performance is about 0.47TFlops. The peak double-precision performance of K40m is 1.43TFlops, so the efficiency of cuDNN is about 32%. Compared with cuDNN, our work can achieve about 3.6 times speedup on performance and about 24% improvement on hardware efficiency.

To illustrate the effectiveness of the optimization methods proposed in this paper, we show the performance of the implementations after adopting different optimization methods in Fig. 8. In our work, we take the implementation of Alg. 2 as the basic version and follow the steps of adopting vectorization design, LDM-related optimization, register-related optimization and instruction-related optimization successively, which forms an optimization process guided by the performance model. Finally, we propose the 4-CG parallelization design and introduce the implementation based on Alg. 4. As we can see, in the optimization process, distinct performance improvement can be achieved in each step and 48 times speedup is achieve in total.

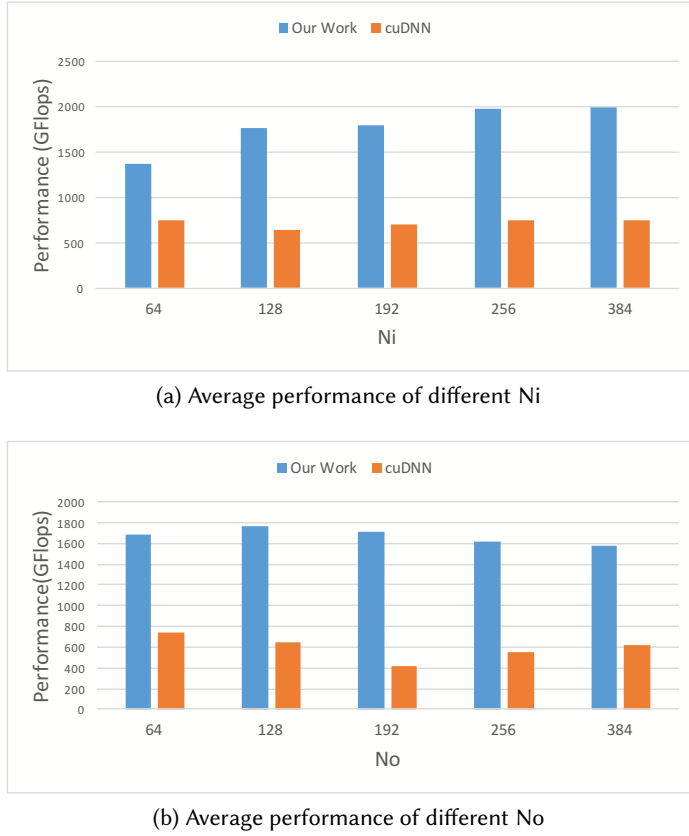
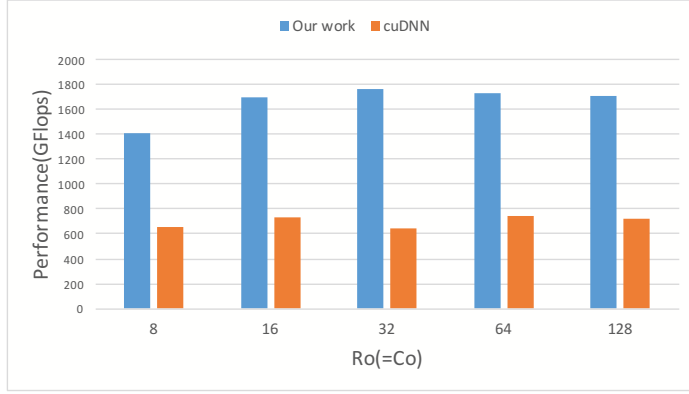


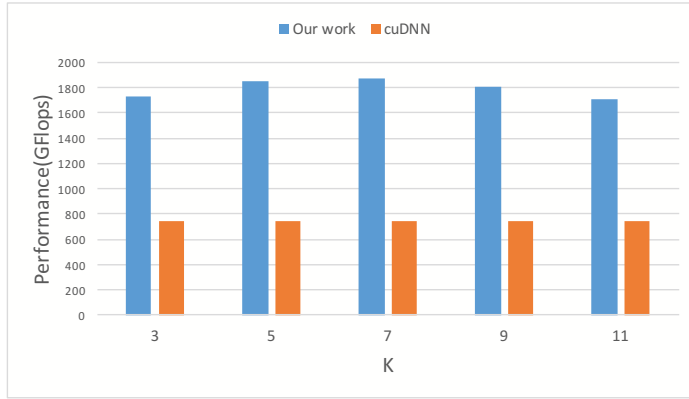
Fig. 6. Performance evaluation on Ni and No (v.s. cuDNNv5.1 on K40m GPU in double-precision)

Generally speaking, some of the proposed optimization techniques, such as the vectorization, register blocking and instruction-related optimization, are also applicable to other heterogeneous architectures like GPU and Intel Xeon Phi. In our work, we consider the features of SW26010 many-core architecture and customize the practical optimization strategies for these general optimization techniques. Besides, other optimization techniques, such as the LDM utilization and register communication, are specific to SW26010 architecture. In summary, both the architecture-specific optimization techniques and the architecture-customized strategies for general optimization techniques are considered as *architecture-oriented* optimization methods, which can also be general for other application and algorithm optimization problems on SW26010 many-core architecture.

We further evaluate the performance of our convolution implementation based on single-precision data representation, which is generally used in the practical training process of CNN models. As discussed in Section 3.7, a straightforward implementation and a float2double implementation are proposed. In the experiment, we train VGG-16 [21] with both float and double data precision based on our work on SW26010 and cuDNN on K40m. There are 13 convolutional layers with different configurations in VGG-16. We show the average performance and hardware efficiency of the convolutional layers in Fig. 9.



(a) Average performance of different Ro(Co)



(b) Average performance of different K



(c) Average performance of different Bs

Fig. 7. Performance evaluation on Ro(Co), K and Bs (v.s. cuDNNv5.1 on K40m GPU in double-precision)

Considering the performance on SW26010, the float2double implementation has about 10% improvement on the hardware efficiency over the straightforward implementation. Therefore,

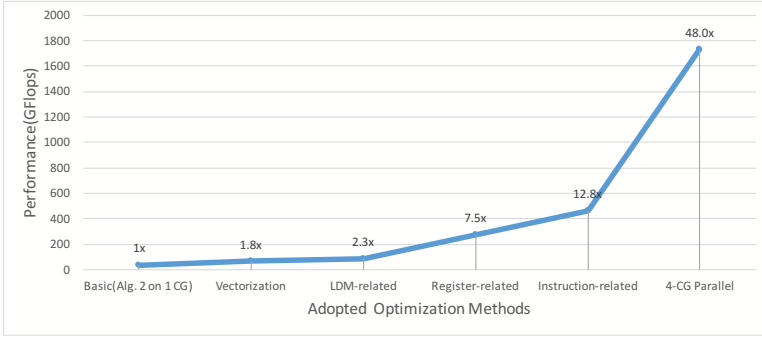


Fig. 8. Performance improvement after adopting different optimization methods

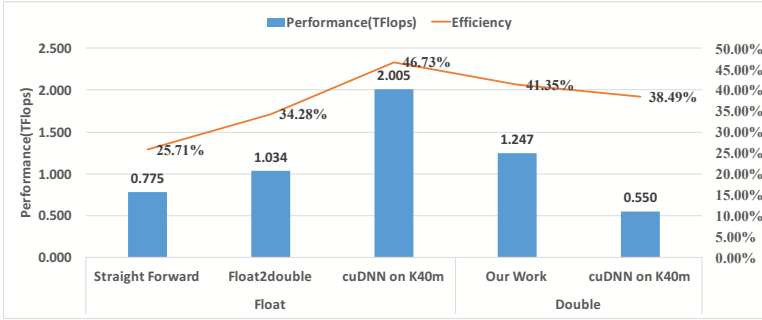


Fig. 9. Average performance and efficiency in float/double-precision (training VGG-16 model)

we adopt the float2double implementation when training CNN models with single-precision. To support more efficient computation in lower data precisions, such as single or half precision, further improvement to the SW26010 architecture is necessary and should target two main aspects. First is to provide optimized SIMD operations for lower data precisions, such as $8\times$ SIMD in single-precision or $16\times$ SIMD in half-precision, which can be realized by using the current 256-bit vector registers and adopting hardware optimization for the ALU part in CPEs. The second aspect is to support more completed low precision instructions (such as *vldr* for single/half precision), so as to improve the overlapping of computation and data access (or data transferring).

4 TRAINING PROCESS OPTIMIZATION

Based on the proposed algorithm optimization methods in Section 3, we further design the swDNN library and an customized Caffe framework, called swCaffe, which can provide a complete solution to train CNN models on Sunway TaihuLight supercomputer.

4.1 swDNN Library

Section 3 focuses on detailed algorithm and code optimization methods for convolution algorithm, which is the most computational intensively part in a CNN. In order to provide a high-performance solution for the complete training process of CNN on the Sunway TaihuLight supercomputer, we further put efforts on optimizing the computation of all kinds of layers with all possible conditions, as well as the backward propagation process to support practical CNN models.

Firstly, we consider different conditions for a convolutional layer. The proposed implementation performs well when the numbers of input and output channels are large enough to assign the tasks to all CPEs. Usually for the first few layers in a practical CNN, the number of channels is small. In this case, the performance of the proposed implementation is poor, so we provide an alternate implementation based on a time-domain transformation method proposed by [13], which contains a *Img2Col* function to transform the input maps to a matrix and a general matrix-matrix multiplication (GEMM) function to do the computation. The GEMM implementation on SW26010 shares the same core computation with the proposed convolution algorithm, so we can skip over the optimization details for brevity. Different implementations are chosen under different conditions, together to support all kinds of convolutional layers.

The fully-connected layer, which is realized through matrix-multiplication, involves the second largest amount of computation in a CNN. The implementation is also based on GEMM.

Besides the computation intensive layers like convolutional and fully-connected layers, other layers can be considered as memory intensive layers, such as pooling layers, normalization layers, and activation function layers. Memory access bandwidth is the key factor that affects the performance of these layers. As shown in Alg. 1 and Alg. 3, the original data layout is (B_s, N_o, R_o, C_o) , and the optimized data layout is (R_o, C_o, N_o, B_s) . Therefore, as discussed in Section 3.1, we propose parallel implementations using CPEs with task partition along the output channel dimension (N_o), which, in both data layout cases, can guarantee a successive memory access with large granularity, so as to take fully advantage of the memory bandwidth. Similarly, data transformation operations, such as the data layout and *Img2Col* transformation can also be accelerated using CPEs.

Usually the output layer of a CNN is a softmax layer. The algorithm of softmax is hard to be parallelized and the computation amount of the softmax layer is rather small. Therefore, there's no need to design a CPE-based implementation for the softmax layer.

Each iteration of the training process contains a forward process and a backward process. Above implementations are focused on the forward process. The output of a forward process is the classification results given by the current model. In the backward process, we first evaluate the *error* of the output results referring to the true labels of the input samples. Then we propagate the error back from the output layer to the input layer and adjust the weights in the layers to minimize the error. In each layer, the backward process shares similar computation patterns with the forward process, but involves approximately two-fold of computation operations for both error propagation and weight update. Therefore, the algorithm design and optimization for the backward process of a layer is similar to the forward process but has different input/output data. We implement CPE-based backward process for each layer to provide a highly-efficient backward propagation in the training process.

Table 4. Summary of swDNN Library

Layers	Conditions	Using CPE	Parallel Strategy	
			1-CG	4-CG
Convolution Layer	Ni and No ≥ 64	YES	Data Transform + Proposed methods	
	Ni or No < 64	YES	Img2Col + GEMM	On batch size
Fully-connected Layer		YES	GEMM	On batch size
Pooling Layer	Max/Min/Avg	YES	On output channel	On batch size
Activation Function	ReLU, Tanh, etc.	YES	On output channel	On batch size
Normalization		YES	On output channel	On batch size
Softmax Layer		NO	None(Only MPE)	On batch size

Integrating the above implementations for different layers and corresponding data transformation functions, we present a library for accelerating deep neural networks on SW26010 many-core architecture, called **swDNN**. A summary of swDNN library is shown in Tab. 4. For each subroutine in swDNN, we provide two implementations. The basic implementation utilizes 1-CG of a SW26010. For the training process of large CNN models, we provide 4-CG parallel implementation to take advantage of the all-shared memory. The 4-CG parallel strategy is to adopt the task partition along the outermost dimension of the data. For the convolutional layers with optimized data layout, the outermost dimension is R_o , as introduced in Section 3.6. For other layers with the original data layout, the outermost dimension is batch size (B_s).

4.2 swCaffe Framework

To support more efficient CNN model development and training task deployment, we port Caffe, an open-source deep learning framework, onto Sunway TaihuLight supercomputer. The original Caffe calls the BLAS library to do the arithmetic computation. On Sunway TaihuLight, swBLAS is one of the fundamental libraries which provide CPE-based implementations on 1 CG. We consider the Caffe framework depending on swBLAS as the basic version, which has no specialized optimization for the CNN models.

Based on the basic version, we propose three optimization methods to customize the Caffe framework for the SW26010 many-core architecture, and finally we present **swCaffe**.

First, we implement swDNN-based layers, as listed in Tab. 4, to substitute for the original implementations of different layers in Caffe.

Table 5. Computation and data transformation time of swCaffe with/without data transformation layer (one iteration of training VGG-16 with $B_s = 128$)

		Computation	Data Trans.			Total
			Weights	Feature Maps	Total	
Without Data Trans. Layer	time(s)	13.55	1.86	3.16	5.02	18.57
	percentage	73%	10%	17%	27%	100%
With Data Trans. Layer	time(s)	13.49	1.83	0.73	2.56	16.05
	percentage	84%	11%	5%	16%	100%

Second, we add new data transformation layer to swCaffe. In most CNN models, there are consecutive convolutional layers and pooling layers that can be accelerated with optimized data layout, such as the 2nd to 5th convolutional layers in AlexNet [14] and the 2nd to 13th convolutional layers in VGG-16. Here we take the convolutional layers in VGG-16 as examples. If we do data transformation for input/output feature maps and weights in each layer, the data transformation time is about 27% of the total execution time of all convolutional layers, as listed in Tab. 5. We add a dedicated data transformation layer into swCaffe, so that for the consecutive convolutional layers, the data transformation of input/output feature maps is performed only once. The data transformation time is reduced to 16% of the total execution time of all convolutional layers. Specifically, the data transformation time for feature maps is reduced to about one fourth (from 3.16s to 0.73s).

Third, we extend swCaffe to support 4-CG parallel in the complete training process. A straightforward way is to utilize 4-CG implementations in swDNN for each layer and the parallel process is shown in Fig. 10(a). As we can see, the training process is started on CG 0. For layers that have 4-CG parallel implementations in swDNN, we call *pthread_create* to start 3 computing threads on CG1-CG3. After the computation finished, we call *pthread_join* to release the computing threads

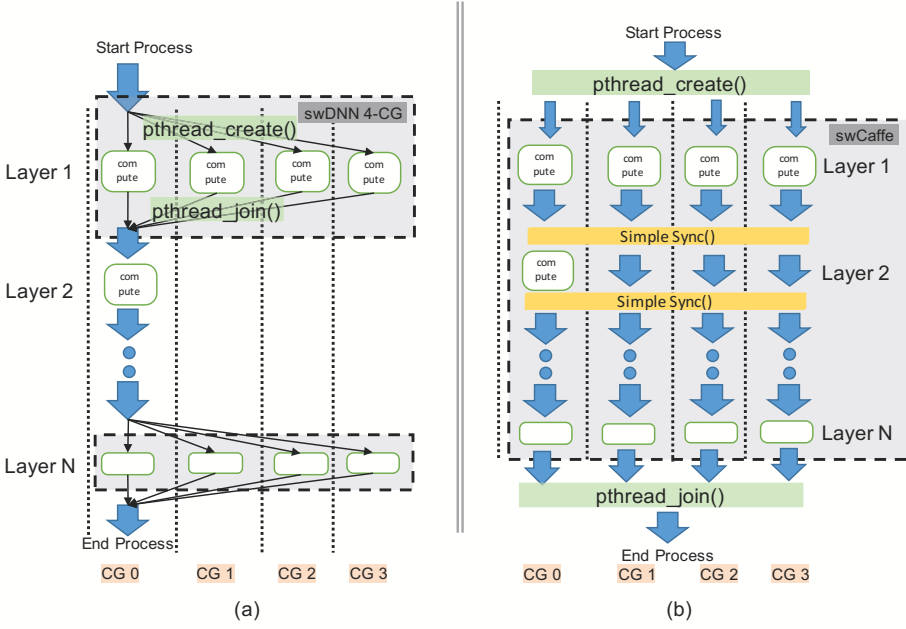


Fig. 10. (a) Caffe using 4-CG implementation in swDNN; (b) 4-CG design for swCaffe

and continue the process on the main thread. In a CNN model, when most of the layers are based on swDNN, calling `pthread_create` and `pthread_join` repeatedly will lead to a relatively large overhead for creating or releasing the thread context.

Addressing the above problems, we propose a framework level parallelization design as shown in Fig. 10(b). At the beginning of the process, we call `pthread_create` to start 4 threads on 4 CGs, all of which will be activated during the whole training process. For layers that can be implemented in parallel, such as Layer 1 in Fig. 10, computation can be done in 4 CGs without extra overhead. For layers that can not be implemented in parallel, such as Layer 2, we firstly call a simple synchronization function to guarantee that all 4 threads are at the same stage, then do the computation on CG 0, and finally call the synchronization function again to continue the process on 4 CGs. Alg. 6 describes the synchronization function (`Simple_Sync()`), which is based on an handshake (initiation-confirmation) strategy through the semaphore (`Signal[NThread]` and `Respond[NThread]`) stored in the shared memory.

4.3 Evaluation

The performance of a complete training process is evaluated based on training VGG-16 model, which is one of the typical and widely-used CNN models. To show the performance improvement obtained from different framework-level optimization methods, we train VGG-16 using three versions of Caffe on Sunway TaihuLight, including:

- Caffe-swBLAS: the basic swBLAS based Caffe, utilizing only 1 CG on SW26010.
- Caffe-swDNN: Caffe with swDNN based layer implementations, utilizing 4 CGs on SW26010.
- swCaffe: swDNN based Caffe with customized data transformation layers and framework parallelization design for 4 CGs.

ALGORITHM 6: Description of Synchronization Function

```

1: //Initialization
2: set  $NThread = 4$ 
3: //Signal[ $NThread$ ] is used to initiate synchronization
4: //Respond[ $NThread$ ] is used to confirm synchronization
5: for  $i := 0 : 1 : NThread$  do
6:   set  $Signal[i] = 0$ 
7:   set  $Respond[i] = 0$ 
8: end for
9: //Function definition
10: define Simple_Sync():
11: set  $thread\_id = get\_thread\_id()$ 
12: if  $thread\_id == 0$  then
13:   for  $i := 1 : 1 : NThread$  do
14:     set  $Signal[i] = 1$  //initiating synchronization on CG 0
15:   end for
16:   set  $nRespond = NThread - 1$ 
17:   while  $nRespond > 0$  do
18:     //waiting for the confirmation from CG 1, 2, 3
19:     for  $i := 1 : 1 : NThread$  do
20:       if  $Respond[i] == 1$  then
21:         set  $nRespond = nRespond - 1$ 
22:         set  $Respond[i] = 0$ 
23:       end if
24:     end for
25:   end while
26: else
27:   //waiting for synchronization on CG 1,2,3
28:   while  $Signal[thread\_id] \neq 1$  do
29:     //waiting for synchronization signal
30:   end while
31:   set  $Signal[thread\_id] = 0$ 
32:   set  $Respond[thread\_id] = 1$  //set the confirmation
33: end if

```

For comparison, we provide the performance of training VGG-16 with Caffe on Intel multi-core CPUs (2×E5-2670v3, 24 cores, with 128GB memory) and NVIDIA K40m GPU. The training dataset is the ImageNet (ILSVRC) 2012 image classification dataset. We use *sample per second (sample/s)* as the metric to show the average training speed. Results are shown in Fig. 11.

As we can see, for the single-precision based training process, the proposed swCaffe framework can achieve about 4.6 times speedup over the basic swBLAS based framework, mainly because of the utilization of 4 CGs. Besides, the optimization targeting data transformation layers and 4-CG parallelization can provide about 20% (speedup from 3.8× to 4.6×) performance improvement. Overall, the proposed optimization methods are proved to be effective.

Compared with CPU and GPU results, swCaffe is 4.6 times more efficient than two 12-core CPUs (based on OpenBLAS) and is nearly half the performance of K40m (based on cuDNNv5.1). As a supplementary, performance of double-precision based training process is also provided. The double-precision performance of swCaffe is even higher than the single-precision performance on SW26010, and is 8.9 and 1.8 times of CPUs and K40m GPU, respectively.

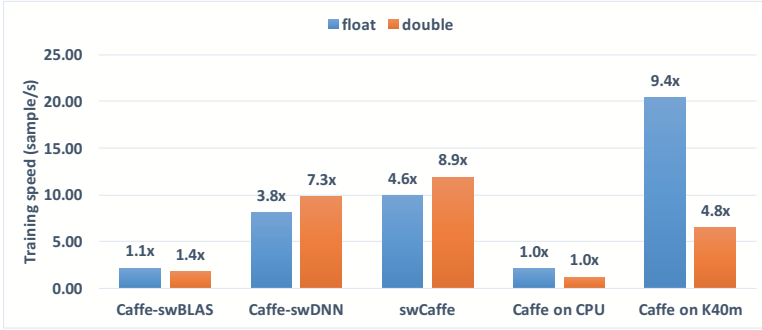


Fig. 11. Performance evaluation of training VGG-16 model

5 CONCLUSIONS

In this paper, we present our work on optimizing the convolutional neural network on SW26010 many-core processor. We propose architecture-oriented optimization methods for the algorithm implementation and framework parallelization. Based on the proposed optimization methods, we develop a customized deep learning library (swDNN) and a customized Caffe framework (swCaffe).

Evaluation results show that the proposed optimization methods can bring 48 times performance improvement to the convolution routine in swDNN compared with the basic implementation. The optimized swCaffe framework achieves 4 times performance improvement for the complete training process of the VGG-16 network, compared with original Caffe with swBLAS. Moreover, the proposed convolution routine in swDNN and the swCaffe framework show nearly half the performance of cuDNN library (on K40m GPU) in single-precision, while achieving 3.6 times and 1.8 times speedup over cuDNN (on K40m GPU) in double-precision, respectively.

The presented work can provide highly-efficient solutions for training CNN models with SW26010 many-core processor. Moreover, it proves the capability of deploying large-scale deep learning applications on Sunway TaihuLight supercomputer.

REFERENCES

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. 2016. Tensorflow: Large-scale Machine Learning on Heterogeneous Distributed Systems. *arXiv preprint arXiv:1603.04467* (2016).
- [2] Kumar Chellapilla, Sidd Puri, and Patrice Simard. 2006. High Performance Convolutional Neural Networks for Document Processing. In *Tenth International Workshop on Frontiers in Handwriting Recognition*. Suvisoft.
- [3] Tianshi Chen, Zidong Du, Ninghui Sun, Jia Wang, Chengyong Wu, Yunji Chen, and Olivier Temam. 2014. Diannao: A Small-Footprint High-Throughput Accelerator for Ubiquitous Machine-Learning. In *ACM Sigplan Notices*, Vol. 49. ACM, 269–284.
- [4] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. 2015. MXNet: A Flexible and Efficient Machine Learning Library for Heterogeneous Distributed Systems. *Statistics* (2015).
- [5] Yunji Chen, Tao Luo, Shaoli Liu, Shijin Zhang, Liqiang He, Jia Wang, Ling Li, Tianshi Chen, Zhiwei Xu, Ninghui Sun, et al. 2014. Dadiannao: A Machine-Learning Supercomputer. In *Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture*. IEEE Computer Society, 609–622.
- [6] Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, Jonathan Cohen, John Tran, Bryan Catanzaro, and Evan Shelhamer. 2014. cuDNN: Efficient Primitives for deep learning. *arXiv preprint arXiv:1410.0759* (2014).
- [7] Ronan Collobert, Samy Bengio, and Johnny Marthoz. 2002. Torch: A Modular Machine Learning Software Library. *Idiap* (2002).
- [8] George E Dahl, Dong Yu, Li Deng, and Alex Acero. 2011. Context-Dependent Pre-Trained Deep Neural Networks for Large-Vocabulary Speech Recognition. *IEEE Transactions on Audio Speech & Language Processing* 20, 1 (2011), 30–42.

- [9] Zidong Du, Robert Fasthuber, Tianshi Chen, Paolo Ienne, Ling Li, Tao Luo, Xiaobing Feng, Yunji Chen, and Olivier Temam. 2015. ShiDianNao: Shifting Vision Processing Closer to the Sensor. In *ACM SIGARCH Computer Architecture News*, Vol. 43. ACM, 92–104.
- [10] Jiarui Fang, Haohuan Fu, Wenlai Zhao, Bingwei Chen, Weijie Zheng, and Guangwen Yang. 2017. swDNN: A Library for Accelerating Deep Learning Applications on Sunway TaihuLight. In *Parallel and Distributed Processing Symposium*. 615–624.
- [11] Haohuan Fu, Junfeng Liao, Jinzhe Yang, Lanning Wang, Zhenya Song, Xiaomeng Huang, Chao Yang, Wei Xue, Fangfang Liu, Fangli Qiao, et al. 2016. The Sunway TaihuLight Supercomputer: System and Applications. *Science China Information Sciences* 59, 7 (2016), 072001.
- [12] Geoffrey Hinton, Li Deng, Dong Yu, George E. Dahl, Abdel Rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, and Tara N. Sainath. 2012. Deep Neural Networks for Acoustic Modeling in Speech Recognition: The Shared Views of Four Research Groups. *IEEE Signal Processing Magazine* 29, 6 (2012), 82–97.
- [13] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. 2014. Caffe: Convolutional Architecture for Fast Feature Embedding. In *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 675–678.
- [14] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet Classification with Deep Convolutional Neural Networks. In *Advances in neural information processing systems*. 1097–1105.
- [15] Andrew Lavin. 2015. maxDNN: An Efficient Convolution Kernel for Deep Learning with Maxwell GPUs. *arXiv preprint arXiv:1501.06633* (2015).
- [16] Andrew Lavin and Scott Gray. 2016. Fast Algorithms for Convolutional Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4013–4021.
- [17] Daofu Liu, Tianshi Chen, Shaoli Liu, Jinhong Zhou, Shengyuan Zhou, Olivier Teman, Xiaobing Feng, Xuehai Zhou, and Yunji Chen. 2015. Pudiannao: A Polyvalent Machine Learning Accelerator. In *ACM SIGARCH Computer Architecture News*, Vol. 43. ACM, 369–381.
- [18] Michael Mathieu, Mikael Henaff, and Yann LeCun. 2013. Fast Training of Convolutional Networks through FFTs. *arXiv preprint arXiv:1312.5851* (2013).
- [19] Jiantao Qiu, Jie Wang, Song Yao, Kaiyuan Guo, Boxun Li, Erjin Zhou, Jincheng Yu, Tianqi Tang, Ningyi Xu, Sen Song, et al. 2016. Going Deeper with Embedded FPGA Platform for Convolutional Neural Network. In *Proceedings of the 2016 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 26–35.
- [20] D Silver, A. Huang, C. J. Maddison, A Guez, L Sifre, den Driessche G Van, J Schrittwieser, I Antonoglou, V Panneershelvam, and M Lanctot. 2016. Mastering the Game of Go with Deep Neural Networks and Tree Search. *Nature* 529, 7587 (2016), 484.
- [21] Karen Simonyan and Andrew Zisserman. 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition. *CoRR* abs/1409.1556 (2014). <http://arxiv.org/abs/1409.1556>
- [22] Yi Sun, Xiaogang Wang, and Xiaoou Tang. 2014. Deeply Learned Face Representations are Sparse, Selective, and Robust. (2014), 2892–2900.
- [23] Nicolas Vasilache, Jeff Johnson, Michael Mathieu, Soumith Chintala, Serkan Piantino, and Yann LeCun. 2014. Fast Convolutional Nets with fbfft: A GPU Performance Evaluation. *arXiv preprint arXiv:1412.7580* (2014).
- [24] Matthew D. Zeiler and Rob Fergus. 2014. Visualizing and Understanding Convolutional Networks. 8689 (2014), 818–833.
- [25] Chen Zhang, Peng Li, Guangyu Sun, Yijin Guan, Bingjun Xiao, and Jason Cong. 2015. Optimizing FPGA-Based Accelerator Design for Deep Convolutional Neural Networks. In *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. ACM, 161–170.
- [26] Chen Zhang, Di Wu, Jiayu Sun, Guangyu Sun, Guojie Luo, and Jason Cong. 2016. Energy-Efficient CNN Implementation on a Deeply Pipelined FPGA Cluster. In *Proceedings of the 2016 International Symposium on Low Power Electronics and Design*. ACM, 326–331.
- [27] Wenlai Zhao, Haohuan Fu, Wayne Luk, Teng Yu, Shaojun Wang, Bo Feng, Yuchun Ma, and Guangwen Yang. 2016. F-CNN: An FPGA-based framework for training Convolutional Neural Networks. In *IEEE International Conference on Application-Specific Systems, Architectures and Processors*. 107–114.