

# Accelerating Deconvolution on Unmodified CNN Accelerators for Generative Adversarial Networks – A Software Approach

Kaijie Tu

Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

tukaijie@ict.ac.cn

## ABSTRACT

Generative Adversarial Networks (GANs) are the emerging machine learning technology that can learn to automatically create labeled datasets in massive application domains such as speech, image, video and texts. A GAN typically includes a generative model that is taught to generate any distribution of data, and a discriminator trained to distinguish the synthetic data from real-world data. Both convolutional and deconvolutional layers are the major source of performance overhead for GANs and directly impacts the efficiency of GAN-based systems. There are many prior works investigating specialized hardware architectures that can accelerate convolution and deconvolution simultaneously, but they entail intensive hardware modifications to the existing CNN accelerators or processors that focus on convolution acceleration. In contrast, this work proposes a novel deconvolution layer implementation with a software approach and enables fast and efficient generative network inference on the legacy Convolutional Neural Networks (CNNs) accelerators. Our proposed method reorganizes the computation of deconvolutional layer and allows the CNN accelerators to treat it as the standard convolutional layer after we split the original deconvolutional filters into multiple small filters. The proposed data flow is implemented on representative CNN accelerators including dot-production array and regular 2D PE array architectures. Compared to the prior baseline acceleration scheme, the implemented acceleration scheme achieves 2.4X - 4.3X performance speedup and reduces the energy consumption by 27.7% - 54.5% on a set of realistic benchmarks.

## 1. INTRODUCTION

Deep neural networks are making continuous breakthroughs in massive research territories over the years. Among them, GANs have been demonstrated to be superior in a broad domain of content-generation applications and unsupervised learning, because they are able to mimic any distribution of data by learning from a small amount of labeled dataset. Unlike the standard Convolutional Neural Networks, GANs are composed of a discriminative model and a generative model. The generative model randomly samples the latent space and generates data for approximation, while the discriminative model distinguishes if the data is produced by the generator. The two models compete with each other in a zero-sum game framework and produce a stronger discriminator and a generator. Typically, the generator often involves convolutional layers and also deconvolutional layers for up-sampling. Both layers are compute-intensive and are the performance bottleneck of GANs. Therefore, it is demanded to accelerate GANs, especially the generative networks on end-devices for real-time and low power operation.

To accelerate GANs with customized hardware, researchers have tried a number of approaches from distinct angles. The

authors in [4] opted to build independent accelerator engines for convolution and deconvolution respectively. This approach requires a large number of hardware resources and chip area. An intuitive improved approach is to reuse the convolution accelerator and build a unified fully convolutional accelerator for both convolution and deconvolution operations. The input data of deconvolution can be reorganized by dynamically padding zero activations and then treat the deconvolution as conventional convolution as presented in Figure 1. Figure 1(a) is the classic deconvolutional operation example with the stride of 2 while Figure 1(b) is converted equivalent convolutional operation with stride set to be 1. Eventually, the deconvolution can be mapped to the convolution accelerator without hardware modification. However, the zero activations induce considerable redundant computing and degrade the performance of GANs which is illustrated in [6]. Different from the above two approaches, the authors in [5] and [6] proposed to revisit the convolutional accelerator and change the micro-architecture to support both convolution and deconvolution efficiently in a unified accelerator. In addition, these methods also need dedicated data flow scheduling to make use of the computing engine. The advantage is the competitive performance and hardware reuse, while the disadvantage is the additional engineering cost or inefficient deconvolutional operations in off-the-shelf CNN accelerators without deconvolution support such as Diannao [2] and TPU [3].

Inspired by the prior work, we seek to support fast and efficient deconvolution layer implementation on general CNN accelerators like Diannao and Eyeriss. For these classic CNN accelerators, many zero-value activations must be padded to the feature map in order to map the deconvolution layers on to them. Instead of zero-padding that induces numerous redundant computing operations, we tailor a novel implementation of deconvolution layer from the software angle, and pre-partition the deconvolutional filters into multiple small convolutional filters, so that the deconvolution operations are converted and efficiently implemented on any CNN accelerator without hardware modification. In our evaluation on classic CNN accelerators, the performance as well as the energy efficiency of our deconvolution implementation remains competitive compared to prior work of specialized GAN accelerators.

In summary, our contributions are listed as follows:

- We proposed a novel filter partitioning and reorganization approach to convert the deconvolution to standard convolution operations such that the deconvolution can be implemented efficiently on general convolution accelerator without hardware modification.
- We explored the use of the deconvolution reorganization method on typical CNN accelerators such as CNN with dot-production/regular 2D computing array and demonstrate the applicability of the method on general CNN accelerators.
- We evaluated the proposed deconvolution performance on a

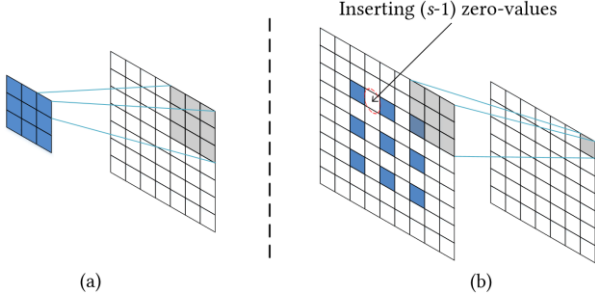


Figure 1: Computational process of (a) deconvolution and (b) deconvolution with inserted zero-values.

set of representative benchmarking networks with comprehensive experiments, the experiments show that the proposed approach achieves competitive performance over the state-of-the-art deconvolution accelerators on a general CNN accelerator.

## 2. RELATED WORK

Compared to the general purposed processors, convolution neural network accelerators outperform on both performance and energy efficiency and attract a large number of resources devoted to the CNN accelerator design and optimization [1, 2, 10, 11]. However, GANs that consist of both compute-intensive convolution and deconvolution operators cannot be fitted to the conventional CNN accelerators directly. Thereby, GANs acceleration have gained massive attention recently. Zhang X et al. in [4] proposed to optimize deconvolution with reverse looping and stride hole skipping. Despite the good performance, combining independent convolution and deconvolution components in an accelerator induce a large chip area and power consumption. Amir Y et al. in [5] proposed a unified MIMD-SIMD accelerator and it avoids the redundant computing brought by the inserted zero activations by adding a set of distributed on-chip buffers. Base on [5], the authors further developed an end-to-end template-based solution, which can generate the optimized synthesizable unified accelerator from a high-level specification of GANs in [8]. Instead of adding zeros to input feature map, Xu et al. in [6] proposed a unified FCN accelerator on top of a bi-direction systolic array. The FCN accelerator performs the computing on original input features. The weight and data of adjacent PEs are shared and passed periodically by taking advantage of the small column buffers added to the 2D PE array.

Some of the researchers attempted to seek approximate approaches for reusing conventional CNN hardware for deconvolution. Shi et al. [7] presented an example of the transformation from deconvolution to convolution with zero-padding to the input feature maps. The zero-padding causes considerable errors because the output feature maps in each layer

will incur some errors when they are cropped [6] and the errors will accumulate with deeper networks. Moreover, the low-precision neural network may overflow frequently during the computing which deteriorates the computing errors. In addition, this work has not gone through the peer review and it is still lack of sufficient experiments. Chang et al. [9] mainly targeted at image super-resolution problems and also proposed an approximate conversion approach. They utilized approximate filter deformation to convert the deconvolutional layers to convolution layers. For super-resolution image reconstruction problem that can tolerate the errors, it works fine, but it causes considerable computing errors in cases such as GANs. More importantly, the approach proposed in [9] opts to rearrange the deconvolutional result on CPU instead of CNN hardware. It is suitable for image super-resolution which has only a single deconvolutional layer, but it will cause massive data communication and synchronization between CPU and the accelerator which usually involves multiple deconvolutional layers. Unlike the above two deconvolution conversion approaches, we proposed a new universal deconvolution partition and reorganization approach to convert the deconvolution to standard convolution without compromising model accuracy or inference performance. This approach allows efficient deconvolution on unmodified convolution accelerators.

## 3. TYPICAL CNN ACCELERATOR

There are two types of CNN accelerators used in this paper. Figure 2 shows a typical CNN accelerator architecture which consists of a dot-production PE array. The array has  $D_{out}$  neural processing units and each neural process unit includes  $D_{in}$  multipliers and an adder tree.  $D_{in}$  input activations and weights are fed to each processing unit every cycle and a dot production can be obtained each cycle because of the pipelined processing unit architecture. All the processing units share the same  $D_{in}$  activations while the weights are different. The results of each processing unit are output elements of a different output channel. When the on-chip buffers cannot accommodate all the data, tiled neural networks can be fitted to the same architecture. This architecture has been adopted in [2, 10].

Another typical CNN accelerator architecture with regular 2D PE array is illustrated in Figure 3. Each PE in the array has a multiplier and an accumulator and it performs all the operations required to yield an output activation. The weights are fed from the first column of the array and flow through PEs from left to right, while the input activations are directly broadcasted to all the PEs in the same column each cycle. Meanwhile, the overlapped input activations will flow across to neighboring PEs on the left when the window slides. Each row of the PE array produces the output activations of one output feature map and each column PE produces the output activations belonging to different output feature maps.

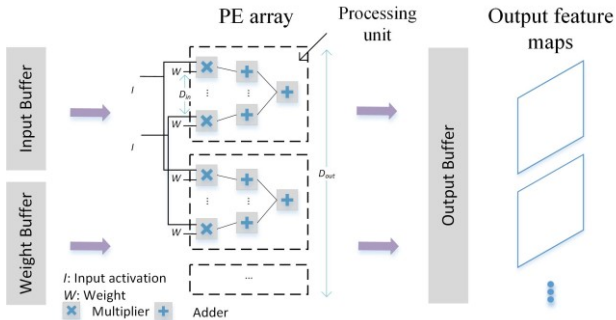


Figure 2: Dot-production based CNN accelerator [2, 10]

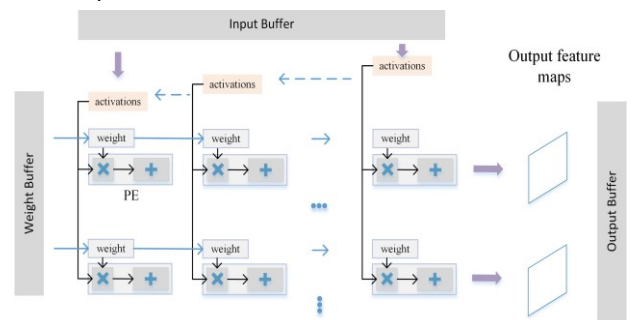


Figure 3: Regular 2D array CNN accelerator [1, 3, 6]

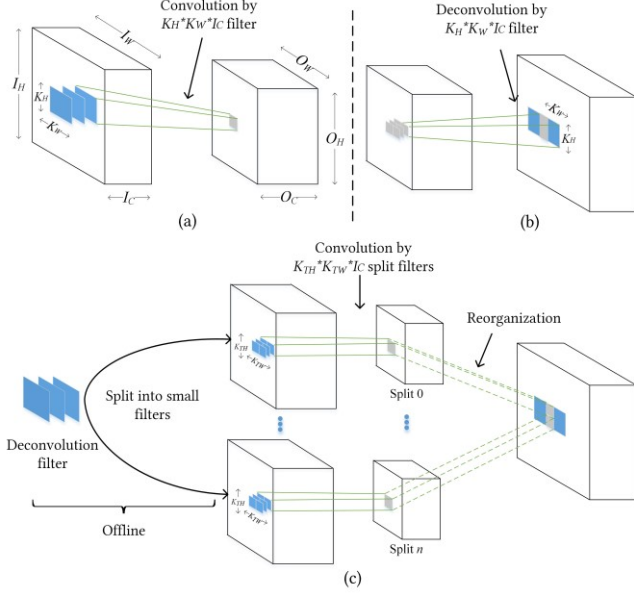


Figure 4: (a) Convolutional layer (b) Deconvolutional layer (c) Split deconvolution that converts a deconvolution layer to multiple convolution layers

## 4. PROPOSED SPLIT DECONVOLUTION

In Section 4.1, we analyze the correlation between the convolution and deconvolution and brief the idea of converting a deconvolution operation to generic convolutions. Then we present the detailed conversion steps from generic deconvolution operations to standard convolution operations in Section 4.2.

### 4.1 Correlation between Convolution and Deconvolution

Convolution and deconvolution are the major sources of overhead in GANs. Figure 4(a) and 4(b) show the basic computing patterns of the two operations. In convolution i.e. Figure 4(a), windows of input features are convolved with the corresponding filters first. Then the results are added up to obtain an output element of the output feature. In deconvolution i.e. Figure 4(b), each element of the input feature maps is multiplied to each weight matrix first. Then the production in the overlapped position will be accumulated as the final output activation. By definition, convolution and deconvolution is completely different.

In order to reuse the conventional CNN accelerators for deconvolution operations, we further look into the computing patterns of convolution and deconvolution. The pseudocode of the two operations are presented in Algorithm 1. Note that  $I_C$  and  $O_C$  indicate the input and output channel of the feature map.  $I_H$  and  $I_W$  denote the length and width of the input feature map, and  $O_H$ ,  $O_W$  are the length and width of the output feature map.  $K_H$  and  $K_W$  is the length and width of the filter.  $s$  refers to stride. The notations will be used through this paper. Basically, both convolution and deconvolution can be computed with elementwise approach. Each output activation of convolution i.e.  $output(o_h, o_w, o_c)$  is the accumulation of production of input feature windows ( $[o_h \times s, o_h \times s + K_H), [o_w \times s, o_w \times s + K_W)$ ) with corresponding filters. For deconvolution, each output activation is also the accumulation of production of input feature map window and a set of weights. The only difference is the filter coordination. In convolution, consecutive weight matrices will be used directly. In deconvolution, the filters are selected with stride  $s$ .

With this observation, we proposed a split deconvolution

### Algorithm 1: Convolution and Deconvolution

```

for ( $o_c = 0$ ;  $o_c < O_C$ ;  $o_c++$ )
  for ( $o_h = 0$ ;  $o_h < O_H$ ;  $o_h++$ )
    for ( $o_w = 0$ ;  $o_w < O_W$ ;  $o_w++$ )
      Convolution ( $o_h, o_w, o_c$ )
      Deconvolution ( $o_h, o_w, o_c$ )

Convolution ( $o_h, o_w, o_c$ ):
  for ( $i_c = 0$ ;  $i_c < I_C$ ;  $i_c++$ )
    for ( $k_h = 0$ ;  $k_h < K_H$ ;  $k_h++$ )
      for ( $k_w = 0$ ;  $k_w < K_W$ ;  $k_w++$ )
         $output(o_h, o_w, o_c) += input(o_h \times s + k_h, o_w \times s + k_w, i_c)$ 
           $\times weight(k_h, k_w, i_c, o_c)$ 

Deconvolution( $o_h, o_w, o_c$ ):
   $lo\_bound\_y = \max(0, \text{ceil}((o_h - K_H)/s))$ 
   $hi\_bound\_y = \min(I_H - 1, lo\_bound\_y + \text{ceil}(K_H/s))$ 
   $lo\_bound\_x = \max(0, \text{ceil}((o_w - K_W)/s))$ 
   $hi\_bound\_x = \min(I_W - 1, lo\_bound\_x + \text{ceil}(K_W/s))$ 
  for ( $i_c = 0$ ;  $i_c < I_C$ ;  $i_c++$ )
    for ( $i_h = lo\_bound\_x$ ;  $i_h \leq hi\_bound\_x$ ;  $i_h++$ )
      for ( $i_w = lo\_bound\_y$ ;  $i_w \leq hi\_bound\_y$ ;  $i_w++$ )
         $output(o_h, o_w, o_c) += input(i_h, i_w, i_c) \times$ 
           $weight(o_h - i_h \times s, o_w - i_w \times s, i_c, o_c)$ 

```

approach which divides the deconvolution filters into multiple smaller filters with stride  $s$ . In this case, the split filters become consecutive and each deconvolution operation is converted to multiple standard convolution operations as illustrated in Figure 4(c). Accordingly, deconvolution can be deployed on conventional CNN accelerators without any hardware modification. While we need to reorganize the filters, detailed conversion approach will be elaborated in the next subsection.

### 4.2 Generic Deconvolution Conversion

Following the above idea, we convert generic deconvolution operation to a set of independent convolution operations. The conversion roughly consists of four steps, as shown in Figure 5.

The first step is the weight preprocessing in which the original deconvolutional filters will be expanded with zeros on the top and left side when its length and width is not divisible by stride  $s$ . It ensures that the deconvolution can be converted to multiple identical convolution operations. The padded zeros will expand the output accordingly while the orientation of the padded zeros

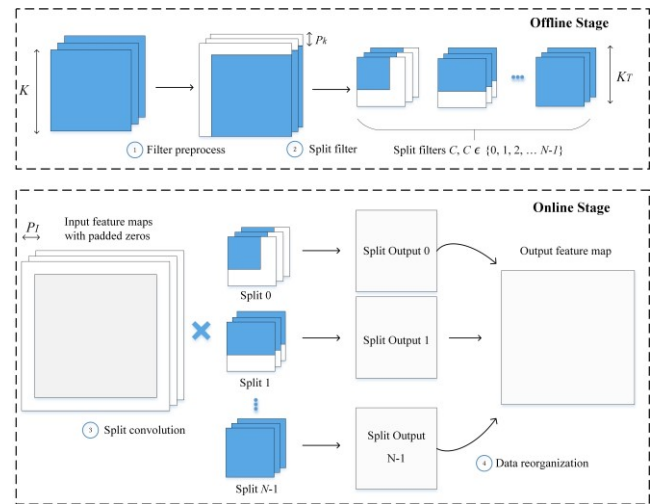
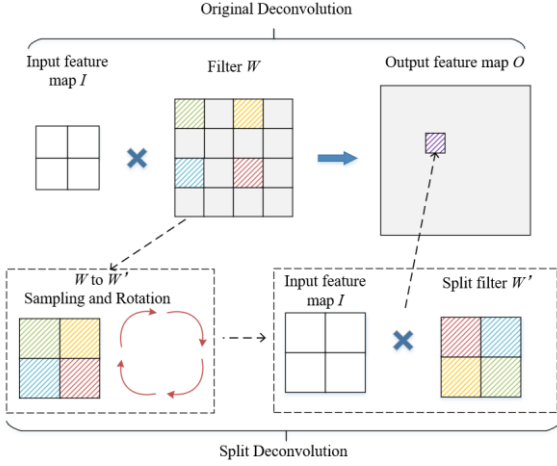


Figure 5: Conversion steps from deconvolution to convolution, it consists of four steps. 1) The filter is expanded when the filter size is not divisible by the stride. 2) Split the deconvolution filters to multiple small filters according to Equations (4-6). 3) The padded input feature maps convolve with the split filters. 4) Reorganize the split deconvolution results to construct the expected deconvolution output by Equations (10) and (11).



**Figure 6: A study case of weight distribution for an output activation in original deconvolution and split deconvolution where the filter is 4 by 4 and the stride is 2.**

guarantees that the center of the expanded output covers the standard deconvolution output. The expanded length and width  $P_K$  can be calculated with Equation (1) where  $K_T$  is the split filter size (assuming it is square) and can be obtained from Equation (2).

$$P_K = s \times K_T - K \quad (1)$$

$$K_T = \text{ceil}(K / s) \quad (2)$$

The second step is to split the deconvolution filters into multiple small filters with sampling and rotation. Figure 6 illustrates the coordinate distribution of filters before and after the conversion with a small but representative example. To compute an output deconvolution activation with standard convolution operations, filters need to be sampled with stride  $s$  and reorganized into new filters. In addition, each sampled filter needs to be rotated 180 degree to ensure correct computing. Equation (3) presents the generic conversion. Each deconvolution will be split into  $s^2$  convolution operations. The stride of the split convolution operations is constant 1. Without loss of generality, suppose  $W_n$  is the  $n$ th convolutional filter. It can be obtained with Equation (4-8) where  $W$  is the deconvolution filter,  $(y, x)$  is the original filter coordinate and  $(y_n, x_n)$  is the new coordinate.

$$N = s^2 \quad (3)$$

$$n = s \times \text{mod}(y, s) + \text{mod}(x, s) \quad (4)$$

$$W_n[y_n][x_n] = W[y][x] \quad (5)$$

$$x_n = K_T - \text{ceil}(x / s) \quad (6)$$

$$y_n = K_T - \text{ceil}(y / s) \quad (7)$$

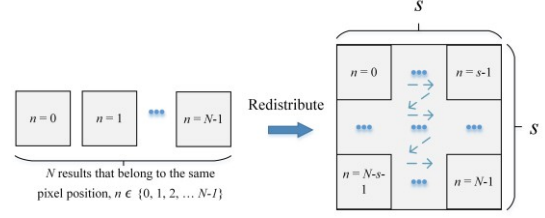
where

$$\begin{aligned} 0 \leq x < K + P_K \\ 0 \leq y < K + P_K \end{aligned} \quad (7)$$

$$\begin{aligned} 0 \leq x_n < K_T \\ 0 \leq y_n < K_T \end{aligned} \quad (8)$$

$$n \in \{0, 1, 2, \dots, N-1\}$$

Step 1 and Step 2 basically split the deconvolution filters to multiple small convolution filters. This needs to be done only once and can be reused. Therefore, they can be done offline with software approach. Unlike the first two steps, Step 3 and 4 are performed on the CNN accelerators for each input feature map. In step 3, the input feature maps also need to be padded with zeros to obtain equivalent deconvolution output. Otherwise, the output



**Figure 7: A demonstration of the redistributed output activations.**

activations on the edge will be ignored.  $P_I$  columns/rows of zeros will be added where  $P_I$  is obtained from Equation (9).

$$P_I = K_T - 1 \quad (9)$$

Finally, the  $N$  split convolution outputs need to be merged to form the deconvolution output. The reorganization pattern is illustrated in Figure 7 and formulated in Equations (10-11). Contrary to the filter splitting process, we pick an element of each convolution output to construct an  $s \times s$  window in the deconvolution output. Note that  $\text{ConvO}_n[x_i][y_i]$  represents the  $n$ th split convolution output and  $O[x_f][y_f]$  refers to the expected deconvolution output. Suppose  $(y_i, x_i)$  is coordinate of a split convolution output and  $(y_f, x_f)$  is the coordinate of deconvolution output. The reorganization here does not need additional hardware as long as the partial convolution output can write the buffers with stride  $s$  which is usually allowed in generic CNN accelerators supporting tiling.

$$O[x_f] = \text{ConvO}_n[x_i] \times s + \text{mod}(n / s) \quad (10)$$

$$O[y_f] = \text{ConvO}_n[y_i] \times s + \text{floor}(n / s) \quad (11)$$

where

$$0 \leq x_i < I + 2P_I - K_T + 1 \quad (12)$$

$$0 \leq y_i < I + 2P_I - K_T + 1$$

$$0 \leq x_f < (I + 2P_I - 1) \times s + K + P_K \quad (13)$$

$$0 \leq y_f < (I + 2P_I - 1) \times s + K + P_K$$

With the above four steps, we can convert generic deconvolution operations to split convolution operations and apply deconvolution on an unmodified CNN accelerator. In spite of the hardware compatibility, the proposed split deconvolution approach may extend the filters and input feature maps, which will induce additional computing overhead. On the other hand, the padding are zeros and can be potentially skipped by the conventional CNN accelerator optimizations. Detailed evaluation on realistic benchmarks will be discussed in the experiments.

## 5. EVALUATION

The section evaluates the performance and energy consumption of the proposed split deconvolution comprehensively with a set of representative benchmark networks first. Then we will demonstrate its use on both DCGAN and Fast-Style-Transfer applications.

### 5.1 Experimental setup

Both the 8-bit dot-production PE array and the 2D PE array are implemented and synthesized with Synopsys Design Compiler (DC) under TSMC 65nm library. For the dot-production array, there are 16 processing units and each unit performs dot production on 16 input activations and weights. The 2D PE array is set to be 32 by 7.

On top of the two different CNN accelerator architectures, we have the proposed split deconvolution implemented. On the 2D PE array architecture, we implemented three different split deconvolution approaches. (1) Split deconvolution on standard





Figure 8: Performance comparison of the deconvolutional layers in the dot-production PE array

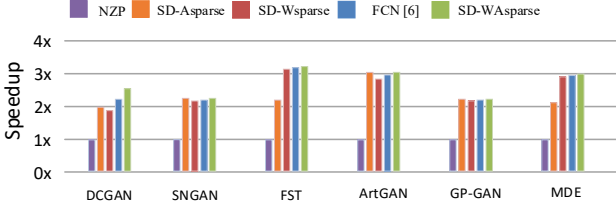


Figure 9: Performance comparison of the deconvolutional layers in the regular 2D PE array

CNN accelerator with the expanded filters (SD-Asparse). The filter’s length or width is not divisible by stride  $s$  so that zeros need to be padded on the top and left side of the original filter to guarantee the split filters are identical, which is conducive to hardware implementation. (2) Split deconvolution on weight-sparsity aware CNN accelerator with the expanded input (SD-Wsparse). The split deconvolution produces only the center area of the original deconvolution output feature maps by default and we add zero padding to the input feature maps to obtain equivalent deconvolution output feature maps. (3) Split deconvolution on weight/activation sparsity aware accelerator with the standard output (SD-WAsparse) which has optimized the influence of expanded filter on a sparsity accelerator and the effect of expanded input by skipping the multiply-add operations of padded zeros. We have the deconvolution approaches realized on two different architectures i.e. generic CNN accelerator with naïve zero-padding (NZP) and FCN-engine for comparison [6]. As the accelerator with dot-production PE array used in this paper is a non-weight-sparsity architecture, we have only split deconvolution with expanded input (SD) and SD-Asparse evaluated. FCN is designed for the 2D PE array, so we only include NZP for comparison.

To evaluate the different deconvolution approaches, we selected a set of advanced neural networks including ArtGAN [13] on Cifar 10 (ArtGAN), DCGAN [12] on Large-scale CelebFaces Attributes Dataset (DCGAN), Spectral Normalization for GAN [14] on Cifar 10 (SNGAN), Unsupervised Monocular Depth Estimation of FCN on KITTI and Cityscapes [15] (MDE), GP-GAN on Transient Attributes Database [16] (GP-GAN), and Fast-Style-Transfer [17] on CoCo2014 as our benchmark (FST).

## 5.2 Experimental results

### 5.2.1 Performance

Since the performance of convolutional layers remains the same, we focus on comparing the performance of the deconvolutional layers. Figure 8 depicts the normalized performance of three acceleration schemes on the dot-production PE array. NZP incurs 75% computing redundancy on average on the benchmark neural networks when converting the deconvolution to convolution. Unlike the NZP, split deconvolution has only marginal zero paddings on the boundary in some corner cases. Therefore, it has much less computing redundancy, which is projected in the 2.5x performance boost of

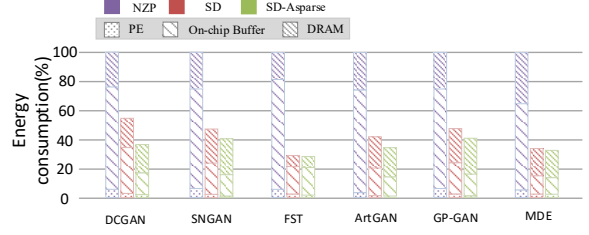


Figure 10: Energy consumption of the deconvolutional layers in the dot-production PE array

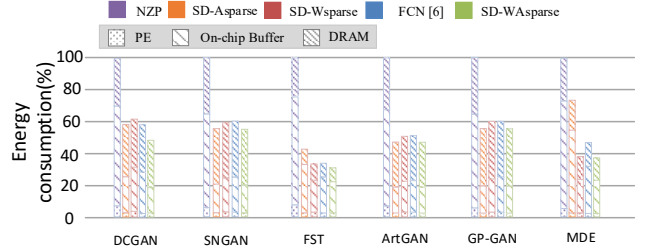


Figure 11: Energy consumption of the deconvolutional layers in the regular 2D PE array

SD over NZP. When the specified input activation lines can be skipped to generate standard deconvolution output, the performance can further be improved. Particularly, SD-Asparse on DCGAN improves by 1.4x. The major reason lies in the fact that the DCGAN has less network layers and smaller input feature maps. As a result, the computing redundancy caused by the padding affects the overall performance more significantly.

On the 2D PE array CNN accelerator as shown in Figure 9, SD-Asparse and SD-Wsparse in the experiments show the influence of the filter expansion and the input expansion respectively. Although SD-Wsparse induces some redundant computation due to padding to the input feature maps, most of the convolution accelerators support zero skipping and can squeeze the computing redundancy automatically. Compared to SD-Wsparse, SD-WAsparse that enables the zero skipping reduces 22% redundant computation on average. Similarly, SD-Asparse has zero padding added to the weights and the redundant computing can also be eliminated on a sparse convolution accelerator architecture. For workloads like DCGAN, FST and MDE, the filters need to be expanded. In these cases, SD-WAsparse reduces 75% ~ 80% computing redundancy with zero skipping. When the split deconvolution is deployed on optimized CNN accelerators, the performance of SD-WAsparse is on par with that of FCN in all the benchmark neural networks. The deconvolution approach presented in FCN-engine [6] adopts a bi-directional data flow. It has the input activations multiplied with each filter and then accumulates the overlapped production. By taking advantage of the column buffers, it can transmit the partial results for accumulation efficiently. However, the output feature maps on the edge are redundant and need to be cropped, which inevitably induces computing overhead especially for smaller deconvolution layers. Therefore, SD-WAsparse outperforms FCN-engine on some of the neural networks like DCGAN as shown in Figure 9.

### 5.2.2 Energy consumption

Figures 10 and 11 present the relative energy consumption distribution of the different deconvolution approaches on the dot-production PE array and regular 2D PE array respectively. Compared to NZP, the average energy consumption of SD-Asparse and SD-WAsparse reduce by 36.15% and 43.63% respectively on the two CNN architectures. Unlike the performance comparison, the energy consumption comparison is

TABLE 1. SSIM VAULE COMPARED WITH STANDARD OUTPUT

Benchmarks	Larger output	Shi [7]	Chang [9]
DCGAN [12]	1	0.534	0.568
FST [17]	0.989	0.939	0.742

less significant. In general, the deconvolution energy consumption roughly consists of three parts i.e. PE, on-chip buffer and DRAM. According to the estimation with CACTI [20], the energy is mostly consumed by the DRAM access and the on-chip buffer access. While the amount of DRAM access of the different deconvolution approaches is about the same, the DRAM consumption has little difference across these approaches. Despite the dramatic difference of the PE activity and energy consumption, the PE energy consumption is too small to affect the overall deconvolution energy consumption. As a result, the energy consumption difference is essentially determined by the amount of on-chip buffer access. This explains all the energy consumption difference. For example, SD-Asparse induces relatively more weight reading and thus higher energy consumption. Similarly, FCN requires additional on-chip buffer to support the unified convolution and deconvolution, so the overall energy consumption is clearly higher than that of SD-WAsparse in all the benchmark networks, though their performance is quite close to each other.

### 5.2.3 Application analysis

Some of the CNN accelerators may not be able to discard part of the output activations and generate the expanded output in split deconvolution and the approaches of [7] and [9] lack versatility for generic GANs which generate unprecise results, we further evaluate the generated images of these three circumstances on two real applications using SSIM metric [19] which is widely used to measure similarity between images. One of them is a typical generative task using DCGAN and the other one is a fast style transfer network. For split deconvolution, the expanded output feature map will not lose any information based on DCGAN. In contrast, the generated images from [7] and [9] are structurally different from the original, as shown by SSIM metric in Table 1. In the application of fast style transfer (FST), the SSIM is 0.989 of expanded output since there are instance normalization layers [18] in the networks, the expanded output feature map will lead to different computing result. Even so, the images we generated are still superior to others in image similarity.

## 6. CONCLUSION

Deconvolution is a critical computing-intensive operation in GANs that have been widely deployed in many generative tasks. Unlike prior deconvolution acceleration work which may either require intensive hardware modification of existing CNN accelerators or bring in large amount of redundant computing, we proposed a novel approach to convert the deconvolution to multiple standard convolution such that the deconvolution can be accelerated efficiently on legacy CNN accelerators without hardware modification. Basically, it reorganizes the deconvolutional layers and splits the original deconvolution filter into multiple smaller filters. Afterwards, we can perform the convolution with the split filters on input feature data. With comprehensive experiments, we demonstrate the deconvolution

approach on both a dot-production PE array and regular 2D PE array. The split deconvolution achieves 2.4X – 4.3X performance speedup over the naïve zero padding methods and is on par with the prior optimized implementation on modified fully convolution neural network accelerator on a set of representative benchmark networks.

## 7. REFERENCES

- [1] Chen Y H, Emer J, Sze V. Eyeriss: A spatial architecture for energy-efficient dataflow for convolutional neural networks[C]//ACM SIGARCH Computer Architecture News. IEEE Press, 2016, 44(3): 367-379.
- [2] Chen T, Du Z, Sun N, et al. Diannao: A small-footprint high-throughput accelerator for ubiquitous machine-learning[J]. ACM Sigplan Notices, 2014, 49(4): 269-284.
- [3] Jouppi N P, Young C, Patil N, et al. In-datacenter performance analysis of a tensor processing unit[C]//Computer Architecture (ISCA), 2017 ACM/IEEE 44th Annual International Symposium on. IEEE, 2017: 1-12.
- [4] Zhang X, Das S, Neopane O, et al. A Design Methodology for Efficient Implementation of Deconvolutional Neural Networks on an FPGA[J]. arXiv preprint arXiv:1705.02583, 2017.
- [5] Yazdanbakhsh A, Falahati H, Wolfe P J, et al. GANAX: A Unified MIMD-SIMD Acceleration for Generative Adversarial Networks[J]. arXiv preprint arXiv:1806.01107, 2018.
- [6] Xu D, Tu K, Wang Y, et al. FCN-engine: accelerating deconvolutional layers in classic CNN processors[C]//Proceedings of the International Conference on Computer-Aided Design. ACM, 2018: 22.
- [7] Shi W, Caballero J, Theis L, et al. Is the deconvolution layer the same as a convolutional layer?[J]. arXiv preprint arXiv:1609.07009, 2016.
- [8] Yazdanbakhsh A, Brzozowski M, Khaleghi B, et al. FlexiGAN: An End-to-End Solution for FPGA Acceleration of Generative Adversarial Networks[C]//2018 IEEE 26th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM). IEEE, 2018: 65-72.
- [9] Chang J W, Kang S J. Optimizing FPGA-based convolutional neural networks accelerator for image super-resolution[C]//Proceedings of the 23rd Asia and South Pacific Design Automation Conference. IEEE Press, 2018: 343-348.
- [10] Song L, Wang Y, Han Y, et al. C-brain: a deep learning accelerator that tames the diversity of CNNs through adaptive data-level parallelization[C]//Design Automation Conference (DAC), 2016 53rd ACM/EDAC/IEEE. IEEE, 2016: 1-6.
- [11] Wang Y, Xu J, Han Y, et al. DeepBurning: automatic generation of FPGA-based learning accelerators for the neural network family[C]//Proceedings of the 53rd Annual Design Automation Conference. ACM, 2016: 110.
- [12] Radford A, Metz L, Chintala S. Unsupervised representation learning with deep convolutional generative adversarial networks[J]. arXiv preprint arXiv:1511.06434, 2015.
- [13] Tan W R, Chan C S, Aguirre H E, et al. ArtGAN: Artwork synthesis with conditional categorical GANs[C]//Image Processing (ICIP), 2017 IEEE International Conference on. IEEE, 2017: 3760-3764.
- [14] Miyato T, Kataoka T, Koyama M, et al. Spectral normalization for generative adversarial networks[J]. arXiv preprint arXiv:1802.05957, 2018.
- [15] Godard C, Mac Aodha O, Brostow G J. Unsupervised monocular depth estimation with left-right consistency[C]//CVPR. 2017, 2(6): 7.
- [16] Wu H, Zheng S, Zhang J, et al. Gp-gan: Towards realistic high-resolution image blending[J]. arXiv preprint arXiv:1703.07195, 2017.
- [17] Fast Style Transfer in TensorFlow, [https:// github.com/lengstrom/faststyle-transfer](https://github.com/lengstrom/faststyle-transfer)
- [18] Ulyanov D, Vedaldi A, Lempitsky V S. Instance normalization: the missing ingredient for fast stylization. CoRR abs/1607.0 (2016)[J].
- [19] Wang Z, Bovik A C, Sheikh H R, et al. Image quality assessment: from error visibility to structural similarity[J]. IEEE transactions on image processing, 2004, 13(4): 600-612.
- [20] Li S, Chen K, Ahn J H, et al. CACTI-P: Architecture-level modeling for SRAM-based structures with advanced leakage reduction techniques[C]//Proceedings of the International Conference on Computer-Aided Design. IEEE Press, 2011: 694-701.