

Proj.2

Zexian Deng, Yaotian Liu
519021910279, 519021910090

I. INTRODUCTION

FOR Proj.2, we implement a CNN accelerating architecture based on *Optimizing the Convolution Operation to Accelerate Deep Neural Networks*[2].

The research question of the paper relies on the fact that prior works lack fully studying the convolution loop optimization before the hardware design phase, resulting in accelerator which hardly exploit the data reuse and manage data movement efficiently. So the author of this paper studies CNN-related techniques, e.g. loop unrolling, tiling and interchange, by quantitatively analyzing and optimizing the design objectives of the CNN accelerator based on multiple design variables. At the same time, based on conclusions above, the author proposed a hardware CNN accelerator which concerns:

- limited computational resources
- storage capacity
- off-chip communication

providing simultaneous maximization of resource utilization and data reuse, and minimization of data communication.

The reason we choose this paper is that, it provides not only a smart new design, but presents a very detailed and in-depth analysis of the basis of convolution loop optimization as well, making it a perfect choice for us to study.

Below is the scheme proposed:

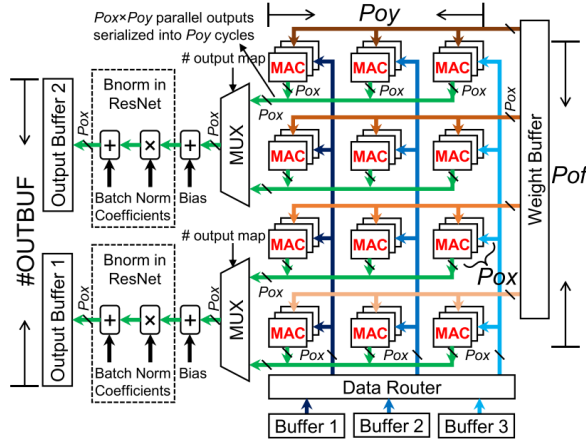


Fig. 1. Convolution acceleration architecture with $Pox \times Poy \times Pof$ MAC units.

During the process of our reproduction, we encounter some difficulties and also discover some flaws that should be corrected or could be further optimized.

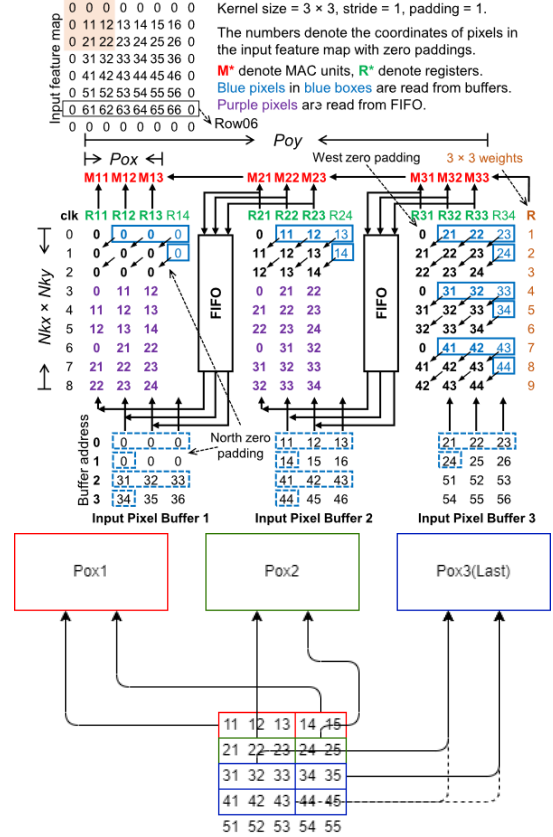


Fig. 2. BUF2PE and the mapping between activation data and PE

II. DIFFICULTIES & ANALYSIS

A. BUF2PE

BUF2PE is where the convolution computation is done, thus the cardinal module in the whole design.

Its design is actually simple and elegant. A MAC unit MAC_{yx} is responsible for calculating the (x, y) pixel in the output feature map¹. Each unit is decoupled from one another, so no adder trees are required.

We define Pox as a compound of multiple MACs which calculates the output pixels in the same row, and Poy as a compound of multiple Pox . So within the Pox , each MAC and its neighbour are using neighbouring pixels in the same row of the input map. Hence for one MAC, its stream of input can be constructed from that of other MACs by offsetting the time

¹To be more precise, its relative position in the $Pox \times Poy$ block currently being processed

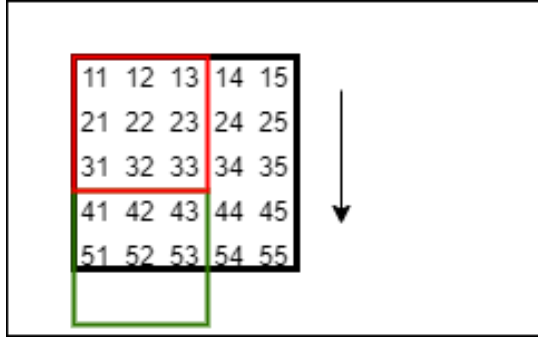


Fig. 3. The current region of output pixels to be produced (in red), the loaded data in buffer (in black), and the next region to be computed (in green)

cycles. Similarly within the *Poy* the neighbouring *Poxes* are using neighbouring rows. This observation is crucial because it gives us a insight on how the data reuse can be done.

a) For cycle 0 \rightarrow $kernel_size - 1$: All MACs are first initialized with their corresponding rows read from buffer. Then the pixels are all shifted to the MAC to the left. From the observation above, we know the new pixels are exactly the new input required. A history of the input pixels are pushed into a FIFO for each *Pox*.

b) For cycle $kernel_size \rightarrow kernel_size * kernel_size - 1$: The last *Pox* continues its practice, while for the rest, their input can be sourced from their neighbouring *Pox*. And so we do. We pop the data from the FIFO and take it as the new input for the remainder of *Pox*.

In the implementation, for each *MAC* a Mux is employed to decided from where a new pixel is fetched. As demonstrated, there can be three choice: **BUFFER**, **SHIFT** and **FIFO**. The state of the Mux is given according to the current row and column, and we use two interdependent counters to obtain them. Aligning the time series is the most tedious work during the design as a result.

B. MUX

The MUX is a module, acting like a regular multiplexer but with a more developed function, providing a further serialization of the outputs from MAC. However, relevant introduction in the paper is far from accurate and detailed, and the figure is also a little deceptive.

According to the paper, "Since $Poy < Nkx \times Nky \times Nif$ for all the layers, we serialize the $Pox \times Poy \times Pof$ MAC outputs into *Poy* cycles." But actually, if

$$N \times Poy < \text{Min}(Nkx \times Nky \times Nif)$$

where Nkx , Nky , Nif are the shape of filters, then there can be a N times folded serialization.

So in our final design, according to our design variables, where N can be equal to Pof , only one MUX is used to greatly save the use of following calculation units.

C. BN & ReLU

According to the paper, ReLU is applied after batch normalization. This is also the original proposition of BN[1]. However, in recent study, this may not be promising compared to applying ReLU before BN.

This question is discussed in Adrian Rosebrock's *Deep Learning for Computer Vision with Python*[3]. In Adrian's view, using BN ahead of ReLU does not make sense.

Under the original setting, a BN layer is normalizing the features coming out of a CONV layer, causing nearly half of the features to be negative due to normalization. Then these negative features will be clamped by a nonlinear activation such as ReLU under our design. That is to say, no matter what comes from CONV layer, even all of the outputs are positive, there are always nearly half of total features be clamped to zero.

Instead, if we place the BN after ReLU, we normalize the positive features without statistically biasing them with features that would have otherwise not make it to the next CONV layer.

What's more funny, according to François Chollet:

"I can guarantee that recent code written by Christian [from the BN paper] applies ReLU before BN."

According to some posts online, they also find out using ReLU before BN can yield a better result.

Of course, there are other opinions supporting using BN first. In this way, the parameters of CNN, like bias, can be merged into BN, resulting in a faster inference. But from my point of view, I prefer a better logic to "tricks" to speed up inference.

Thus, in our reproduction, we apply ReLU before BN.

$$\begin{aligned}
 x' &= x + b \\
 x'' &= \gamma \cdot \frac{\text{ReLU}(x') - E[x]}{\sqrt{\text{Var}[x] + \epsilon}} + \beta \\
 \Rightarrow x'' &= \frac{\gamma}{\sqrt{\text{Var}[x] + \epsilon}} \cdot \text{ReLU}(x) + (\beta - \gamma \cdot \frac{E[x]}{\sqrt{\text{Var}[x] + \epsilon}}) \\
 x'' &= K \cdot \text{ReLU}(x) + B
 \end{aligned}$$

$$\text{where } K = \frac{\gamma}{\sqrt{\text{Var}[x] + \epsilon}}, \quad B = \beta - \gamma \cdot \frac{E[x]}{\sqrt{\text{Var}[x] + \epsilon}}$$

D. Fixed Point

Another major issue we encounter is fixed point. As the paper says, the proposed architecture uses 16-bit fixed point precision pixels and weights.

At first, we suppose the fixed-point multiplication should not be a big problem. With a little deeper thought, we realize the condition of its overflow could be really tricky. After trying classical fixed-point multiplication algorithm such as Booth, we still find it difficult to determine the overflow. One path to this is check whether the upper bits contains valid information, that is to say, all 0 for positive numbers or all 1 for negatives. At last, we turn to a third party fixed point multiplier with overflow detection to solve this issue.

How to choose the right position of decimal point is also a problem. After we extract pixels and weights from our trained model, we find most parameters' absolute value is around 10^{-3} to 1, but some, for instance, one of fused parameters in BN, is more than 500 due to the small variation. We guess the reason behind it might be ReLU, which eliminates all negatives, contributing to a smaller variation.

Back to our topic, we have to choose a suitable fixed point setting. If we try to cover the bigger part, most pixels and weights around 10^{-3} might suffer serious detriment, for the fact that 10 bits for decimal can only provide $2^{-10} \approx 10^{-3}$ precision. Even if we only use 4 bits for integer, there are only $16 - 1 - 4 = 11$ bits left for decimal, which is highly insufficient.

During my lunch time, a walk to my take-out, an idea came to me. We can use a dynamic fixed point notation, with an extra $\log(\#bits)$ to indicate the position of decimal point for each fixed point number. The initial parameters for this extra indicator comes from software, say, the driver. The driver examines every weight and pixel and provides a suitable, custom setting of decimal point. For weights and pixels produced in the process, the dynamic decimal point is determined by both operands. In addition, the decimal point is inferred by the bigger one. In multiplication, the decimal point is inferred by the sum of these two.

Even this process seems time and space consuming, but I believe, in this way, using less bits for fixed point numbers can maybe result much less accuracy degradation.

Sad story, I do not have enough time to do more experiment on my dynamic fixed point design..

III. TEST

A. Test Model

To simplify our verification, we only use a rather small CNN model on MNIST dataset. This CNN model is based on the classical LeNet-5 with BN but without implicit linear layer. The model can achieve a 98% Acc. easily.

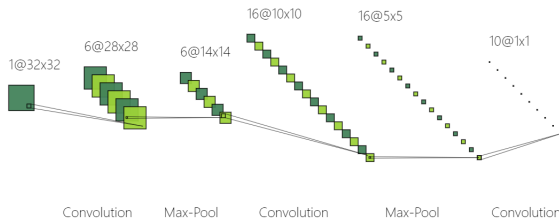


Fig. 4. Out test model

After each convolution layer, we apply a ReLU and BN.

Then we analyze the distribution of parameters and pixels under inference.

As Fig. 5 shows, parameters has a much smaller variation, with most value inside $(-0.25, 0.25)$, and the mean value is about 0. But the distribution of pixels is much more complicated. Due to ReLU, most negatives are clamped to

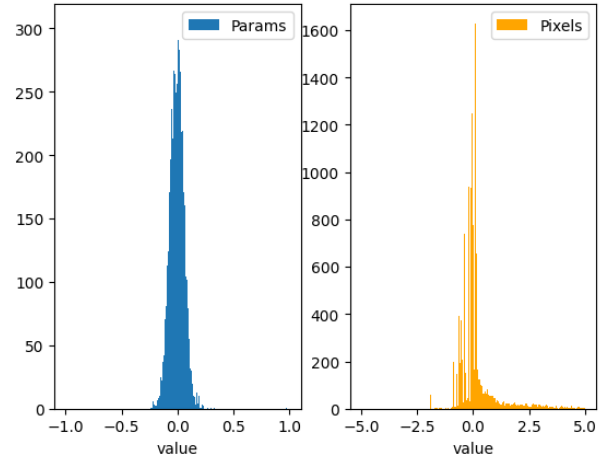


Fig. 5. Distribution of parameters and pixels

0. At the same time, the variation is larger, with noticeable number of pixels extended to positive 5, actually the biggest value can reach more than 20 during inference, causing the fixed point setting difficult and tricky.

B. Test Platform

The test platform is constructed in cocotb. This Python package provides a friendly interface for interacting with the signals of out dut, while unleashing Python's potential for data analysis.

Before initializing the dut, we first construct feature maps and weights as the testcase. This is generated through numpy's rand function, from which a matrix of randomized elements can be easily created. We further use the convolution function from Homework 1 to calculate the real result of our testcase.

The timing sequence of the input requires extreme discretion. First, we supply the dut with a forked clock stimulus and wait some cycles before pulling down the reset signal, initiating the circuit. We extract a part of the feature map which the dut would handle, padding it if necessary. Then we extract the data from it, and drive the ports as described in II-A. When the computation is complete, we must save the accumulated result and reset the circuit to avoid contamination to the next cycle. While designing the inner logic of BUF2PE to ensure the correctness of data flow is indeed no easy job, aligning the counters, activation and weight input is a true nightmare.

The obtained result is around, but not precisely equal to the result from what we obtained from the function. We conjecture this is due to a loss of precision, as while the function computes the result in floating point arithmetics, the circuit use fix-point arithmetics whose multiplication involves shifting and rounding.

REFERENCES

- [1] Sergey Ioffe and Christian Szegedy. *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*. Mar. 2, 2015. DOI: 10.48550/

arXiv.1502.03167. arXiv: 1502.03167 [cs]. URL: <http://arxiv.org/abs/1502.03167> (visited on 09/23/2022).

- [2] Yufei Ma et al. “Optimizing the Convolution Operation to Accelerate Deep Neural Networks on FPGA”. In: *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 26.7 (July 2018), pp. 1354–1367. ISSN: 1557-9999. DOI: 10.1109/TVLSI.2018.2815603.
- [3] Adrian Rosebrock. *Deep Learning for Computer Vision with Python*. September 2017.