

A Report on

# **Low-Power VLSI Implementation of CNN on Accelerated 2D Systolic Array**

Team RizzNet (CE)

*Krish Mehta*

*Aryan Devrani*

*Anoushka Saraswat*

*Kumar Divij*



ECE 284

University of California, San Diego (UCSD)

16th Dec. 2023

# 1. Motivation

2D Systolic Architecture enables modern high throughput AI/ML computation. We aim to implement an optimized 2D Systolic Array Design, with the goal of learning the ins and outs of its implementation, and various possible enhancements to the baseline architecture. The improvements can be broadly categorized into:

- Parallelizing the Systolic Array Computation Stages
- Reducing FIFO Depth from 64 to 16
- Sparsity Aware Clock Gating
- Scalable RTL Design for Multi-channel PEs
- Additional Analysis/Verification

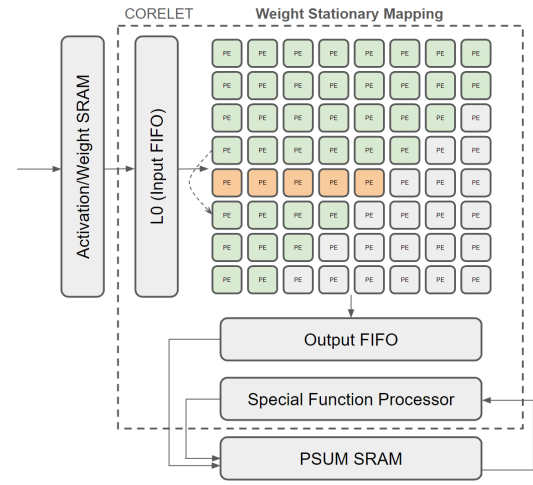
# 2. Baseline RTL Design & Testbench

Our design is laid out as follows: The “core\_tb” is the testbench, which is also the control logic that operates the top level module, the “core”. The core comprises of two SRAMs, one for storing Activation and Kernel values (“xmem”) and one for storing the Partial Sums (“pmem”). The systolic array portion of the design is under the “corelet” which is part of the core. The corelet has the Input and Output FIFOs (L0 and OFIFO) and the MAC Array. MAC Array is made up of 8 MAC Rows, each containing 8 PEs.

The underlying PE in our MAC Array is written such that it can handle the MAC operation of 1 or 2 input channels (in one cycle) as dictated by a design parameter. This is described in detail in section 7.

The control logic in the testbench is well optimized in order to orchestrate the computation in an efficient pipelined and parallelized manner. This is described in detail in the next section. We have implemented a weight stationary mapping. The data flows from the model (inputs/weights) to SRAM -> L0 -> MAC Array -> OFIFO -> SRAM -> SFU -> Output.

We have verified our design at multiple steps along the way. During the development process we verified that the data written to SRAM and the data being written to L0 are matching. We verified that the PSUMs corresponding to each Output position for each Output channel are matching the ones expected by the model. And ultimately we also verified that the end-to-end Convolution + ReLU result matches the result generated by pytorch.



# 3. Parallelizing Systolic Array Computation

In the baseline model, all Kernel Loadings & MAC Operations are processed in a serial manner. In order to reduce the total number of execution cycles, these operations were pipelined to execute parallelly by modifying the control logic of the MAC Array Architecture:

## 3.1 Simultaneous Read/Write for Input & Output FIFOs

The L0 (IFIFO) and OFIFO can perform both Read & Write operations in the same cycle simultaneously, thus reducing the number of cycles taken to serially perform the two tasks.

## 3.2 Pipelining of MAC Load & Execute Operations

Instead of waiting for all Kernel Loadings to complete, when PE(x,y) is loading its weights, PE(x, y-2) starts its execution as it has already received its inputs. This increases the overall utilization of the MAC array as it stays idle for a much lesser time.

## 3.3 No Reset for L0 and OFIFO between Weight Change

This allows for new weights to start loading (from SRAM->L0->MAC) while the previous PSUMs are still being written to Memory.

## Improvements:

- Total Execution Cycles: **1973 -> 917 (54% reduction)**
- Full MAC Utilization: **17% -> 28% (1.65X increase)**

## 4. Reducing L0 & OFIFO Depth from 64 to 16

After achieving this level of parallelism in the MAC computations, the FIFOs do not require a 64 depth. Hence, the FIFO depths were reduced to 16 by using a single 16:1 mux (instead of 4).

### Improvements:

1. Number of Logical Elements: **17063 -> 7199 (67% reduction)**
2. Core Dynamic Power Consumption: **31.86 mW -> 18.68 mW (40% reduction)**
3.  $F_{max}$  for Slow 100C Model: **129.95 MHz -> 131.6 MHz**
4. TOP/s-W : **0.522 -> 0.900**
5. Total Execution time: **15.18  $\mu$ s (vanilla) -> 7.06  $\mu$ s (parallelized) -> 6.97  $\mu$ s (parallelized+fifo\_16)**

The following graph shows the flow of Pipelined MAC Operations across multiple iterations:

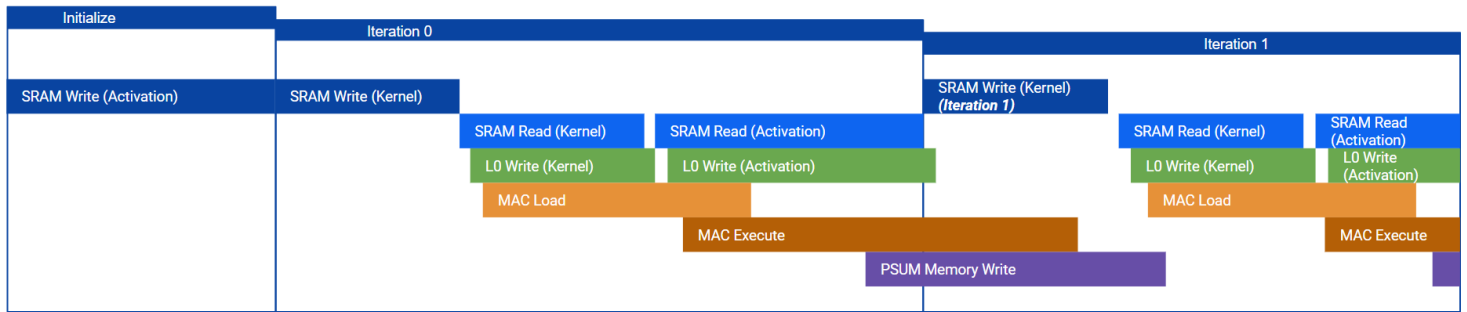


Fig.1 Timeline of Pipelined MAC Operations

In order to analyze and prove the enhancements achieved by these steps, we performed a Utilization Analysis for SRAM, L0 and MAC operations to observe the achieved parallelism:

Total Execution Time	1973 cycles (Serial) 917 cycles (Parallel)
Stage	Utilization (% of total execution time)
MAC Active (At least 1 PE Loading/Executing)	34 % (Serial) 52 % (Parallel)
MAC Fully Active (Producing 8 PSUMs in a cycle)	17 % (Serial) 28 % (Parallel)
L0 FIFO Reading+Writing	0 % (Serial) 41 % (Parallel)

```
[ 883] [14th] After RELU sfp_out: 000000410000011400000008;
14-th output featuremap Data matched! :D
[ 900]
[ 900] [15th] After RELU sfp_out: 000000280000011d000000a5;
15-th output featuremap Data matched! :D
[ 917]
[ 917] [16th] After RELU sfp_out: 0000004a000000b5000d0031;
16-th output featuremap Data matched! :D
##### No error detected #####
##### Project Completed !! #####

===== Systolic Array Statistics =====
SRAM Read   : 45 percent
SRAM Write  : 11 percent
L0 Read     : 43 percent
L0 Write    : 43 percent
L0 Rd/Wr    : 41 percent
MAC Active  : 52 percent
MAC Fully Active : 28 percent
Accumulation : 29 percent
=====

[ 917] Last cycle
./verilog/core_tb.v:646: $finish called at 917000 (1ps)
PS C:\Users\DEVVRANI\Desktop\ece284fa23\project> |
```

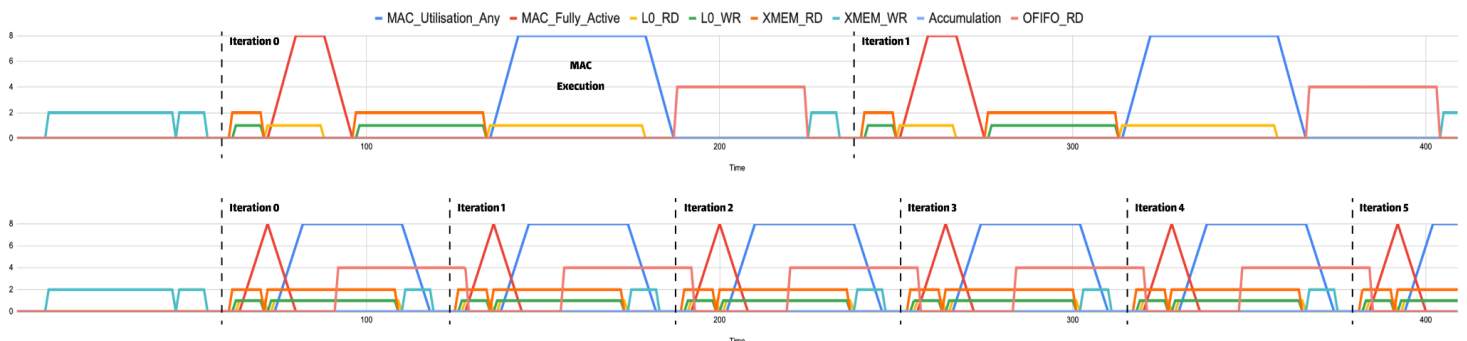


Fig.2(a) Statistics for achieved % Utilization across different stages; 2(b) Verilog output for Utilization; 2(c) Utilization Graphs for Baseline (serialized) and the Alpha (parallelized) models.

## 5. Sparsity Aware Clock Gating

In order to reduce Dynamic Power consumption for a core, we can reduce the amount of toggling operations that are made even when the inputs are sparse.

### 5.1 VGGNet Training & Structured Pruning

The VGG16 model was first trained with 4-bit Quantization Aware training, and its 27th Conv layer squeezed to 8x8. The model achieved 92.07% accuracy. Then, structured pruning on the 4-bit Quantized VGGNet model allowed us to introduce high sparsity levels (>50%), and we were still able to recover >90% model accuracy on the sparse inputs with finetuning. This was a good indicator to try and exploit the sparse inputs by reducing the amount of toggling operations for trivial calculations. Note that only structured pruning could have been useful to our case, as we need an entire MAC row to be 0.

Model Type	Accuracy	Psum Error
4-bit Quantized VGGNet (with 8x8 Conv Layer)	92.07%	$1.0455 \times 10^{-7}$
4-bit Quantized VGGNet (with 8x8 conv + 50% structured pruning)	91.72%	$1.722 \times 10^{-7}$
4-bit Quantized VGGNet (with 8x8 conv + 70% structured pruning)	88.35%	$1.3864 \times 10^{-7}$

Fig.3 Model Accuracies and Psum Error with quantization and pruning

### 5.2 Input-Channel Wise Clock Gating

- There are 8 clock lines going to each of the 8 mac\_row modules, from there it passes to each of 8 mac\_tile modules. Both the free running clock and gated clock are sent to every mac\_tile.
- Gated clocks are generated if ALL weights in a row are 0 (Zero-Condition).
- If Zero-Condition occurs, gated clock prevents new data latching and thereby multiplication calculations, by freezing the horizontal movement of data.
- Free running clock still ensures proper movement of partial sum data from north to south.

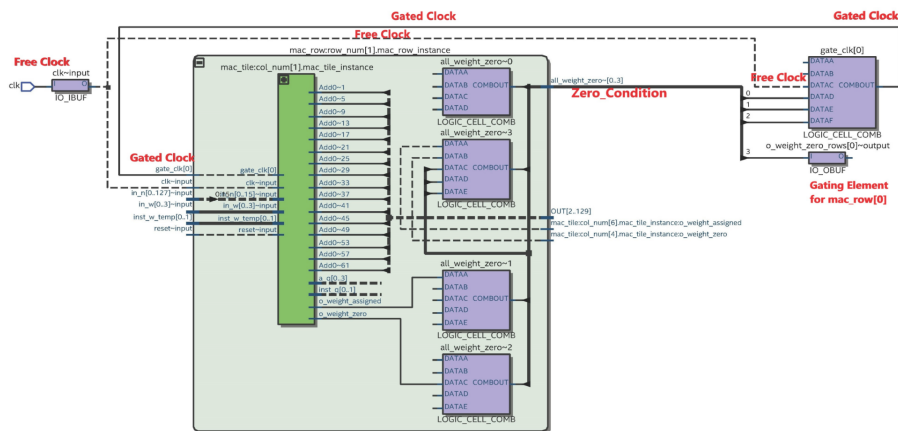


Fig.4 Structure of Input-Channel Wise Clock Gating

**Improvements:** Core Dynamic Power Consumption: **31.86 mW (Vanilla) -> 18.68mW (Parallelized + FIFO16) -> 10.90mW (Parallelized + FIFO16 + Clock Gated)**

## 6. Hardware Mapping, Timing & Power Analysis

We performed hardware mapping & analysis for 4 different models, using Quartus Prime (Cyclone IV):

1. **VGGNet (Vanilla)** - MAC Operations are serialized.
2. **VGGNet Parallelized (Control Only)** - MAC Operations are pipelined, FIFO depth is 64
3. **VGGNet Parallelized (Control + FIFO Depth 64->16)** - Pipelined + FIFO Depth reduced to 16
4. **VGGNet Parallelized + Clock Gated** - Pipeline + FIFO Depth 16 + Input-channel Clock Gating

Power Analyzer Summary	
<<Filter>>	
Power Analyzer Status	Successful - Fri Dec 15 09:48:42 2023
Quartus Prime Version	20.1.0 Build 711 06/05/2020 SJ Lite Edition
Revision Name	corelet
Top-level Entity Name	corelet
Family	Cyclone IV GX
Device	EP4CGX150DF3117AD
Power Models	Final
Total Thermal Power Dissipation	322.13 mW
Core Dynamic Thermal Power Dissipation	18.68 mW
Core Static Thermal Power Dissipation	119.28 mW
I/O Thermal Power Dissipation	184.18 mW
Power Estimation Confidence	Low: user provided insufficient toggle rate

Power Analyzer Summary	
<<Filter>>	
Power Analyzer Status	Successful - Sun Dec 03 15:17:34 2023
Quartus Prime Version	20.1.0 Build 711 06/05/2020 SJ Lite Edition
Revision Name	corelet
Top-level Entity Name	corelet
Family	Cyclone IV GX
Device	EP4CGX150DF3117AD
Power Models	Final
Total Thermal Power Dissipation	335.42 mW
Core Dynamic Thermal Power Dissipation	31.86 mW
Core Static Thermal Power Dissipation	119.58 mW
I/O Thermal Power Dissipation	183.99 mW
Power Estimation Confidence	Low: user provided insufficient toggle rate

Power Analyzer Summary	
<<Filter>>	
Power Analyzer Status	Successful - Sat Dec 16 13:39:53 2023
Quartus Prime Version	20.1.0 Build 711 06/05/2020 SJ Lite Edition
Revision Name	corelet
Top-level Entity Name	corelet
Family	Cyclone IV GX
Device	EP4CGX150DF3117AD
Power Models	Final
Total Thermal Power Dissipation	309.48 mW
Core Dynamic Thermal Power Dissipation	10.90 mW
Core Static Thermal Power Dissipation	119.28 mW
I/O Thermal Power Dissipation	179.30 mW
Power Estimation Confidence	Low: user provided insufficient toggle rate data

Fig.5 Power Analysis for (a) Baseline (b) VGGNet Parallelized (Control+FIFO) & (c) VGGNet Parallelized + Clock Gated At their respective  $F_{max}$  frequencies

Fitter Summary	
<<Filter>>	
Fitter Status	Successful - Fri Dec 15 12:07:28 2023
Quartus Prime Version	20.1.0 Build 711 06/05/2020 SJ Lite Edition
Revision Name	corelet
Top-level Entity Name	corelet
Family	Cyclone IV GX
Device	EP4CGX150DF3117AD
Timing Models	Final
Total logic elements	17,063 / 149,760 (11 %)
Total registers	12098
Total pins	452 / 508 (89 %)
Total virtual pins	0
Total memory bits	0 / 6,635,520 (0 %)
Embedded Multiplier 9-bit elements	0 / 720 (0 %)
Total GXB Receiver Channel PCS	0 / 8 (0 %)
Total GXB Receiver Channel PMA	0 / 8 (0 %)
Total GXB Transmitter Channel PCS	0 / 8 (0 %)
Total GXB Transmitter Channel PMA	0 / 8 (0 %)
Total PLLs	0 / 8 (0 %)

Fitter Summary	
<<Filter>>	
Fitter Status	Successful - Fri Dec 15 09:44:43 2023
Quartus Prime Version	20.1.0 Build 711 06/05/2020 SJ Lite Edition
Revision Name	corelet
Top-level Entity Name	corelet
Family	Cyclone IV GX
Device	EP4CGX150DF3117AD
Timing Models	Final
Total logic elements	7,199 / 149,760 (5 %)
Total registers	4354
Total pins	453 / 508 (89 %)
Total virtual pins	0
Total memory bits	0 / 6,635,520 (0 %)
Embedded Multiplier 9-bit elements	0 / 720 (0 %)
Total GXB Receiver Channel PCS	0 / 8 (0 %)
Total GXB Receiver Channel PMA	0 / 8 (0 %)
Total GXB Transmitter Channel PCS	0 / 8 (0 %)
Total GXB Transmitter Channel PMA	0 / 8 (0 %)
Total PLLs	0 / 8 (0 %)

Fitter Summary	
<<Filter>>	
Fitter Status	Successful - Sat Dec 16 13:37:37 2023
Quartus Prime Version	20.1.0 Build 711 06/05/2020 SJ Lite Edition
Revision Name	corelet
Top-level Entity Name	corelet
Family	Cyclone IV GX
Device	EP4CGX150DF3117AD
Timing Models	Final
Total logic elements	7,584 / 149,760 (5 %)
Total registers	4546
Total pins	453 / 508 (89 %)
Total virtual pins	0
Total memory bits	0 / 6,635,520 (0 %)
Embedded Multiplier 9-bit elements	0 / 720 (0 %)
Total GXB Receiver Channel PCS	0 / 8 (0 %)
Total GXB Receiver Channel PMA	0 / 8 (0 %)
Total GXB Transmitter Channel PCS	0 / 8 (0 %)
Total GXB Transmitter Channel PMA	0 / 8 (0 %)
Total PLLs	0 / 8 (0 %)

Fig.6 Synthesis Fitter Summary for (a) Baseline (b) Parallelized (Control+FIFO) & (c) VGGNet Parallelized + Clock Gated

(NOTE: Hardware mapping results are identical for Vanilla & Parallelized (Control Only), as only the control logic has changed, there is no change in RTL components. They differ in total execution cycles as summarized in the table below.)

Parameter	VGGNet (Vanilla)	VGGNet Parallelized (Control only)	VGGNet Parallelized (Control + FIFO Depth 16 )	VGGNet Parallelized + Clock Gated
$f_{max}$ (Slow 100C Model)	129.95 MHz	129.95 MHz	131.6 MHz	80.78 MHz
Dynamic Power (at $f_{max}$ )	31.86 mW	31.86 mW	18.68 mW	10.90 mW
Power at matched freq. (80 Mhz)	19.76 mW	19.76 mW	11.41 mW	10.90 mW
TOP/s	0.0166	0.0166	0.0168	0.0103
TOP/s-W	0.522	0.522	0.900	0.945
Logical Elements	17063	17063	7199	7584
Total Execution Cycles	1973	917	917	917
Total Execution Time	15.18 $\mu$ s	7.06 $\mu$ s	6.97 $\mu$ s	11.35 $\mu$ s

● -> Worst ● -> Middle ● -> Best

- > Upon parallelizing the Control logic, total execution cycles are reduced by **54%**.
- > In the parallelized model, reducing FIFO depth to 16 allows higher  $f_{\max}$ , reduces no. of logical elements by **67%**, reduces Dynamic power consumption by **40%**, and increases TOP/s-W by **72%**.
- > By introducing Clock Gating to the above model, the dynamic power further reduces by **40%**.
- > The tradeoff between power consumption & execution time depends on what the designer wants to optimize. However, if we look at the TOP/s-W metric which combines both Power and Time, the Clock Gated model clearly outperforms the rest.

## 7. Scalable RTL Design for Multichannel PEs

## 8. Additional Analysis/Verification

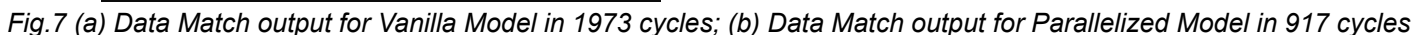


Fig. 8 Power Analysis at common  $F=80$  MHz for (a) Vanilla; (b) Parallelized(control+fifo16) & (c) Parallelized + Clock Gated

5