# Tiled Matrix Multiplication using Systolic Array

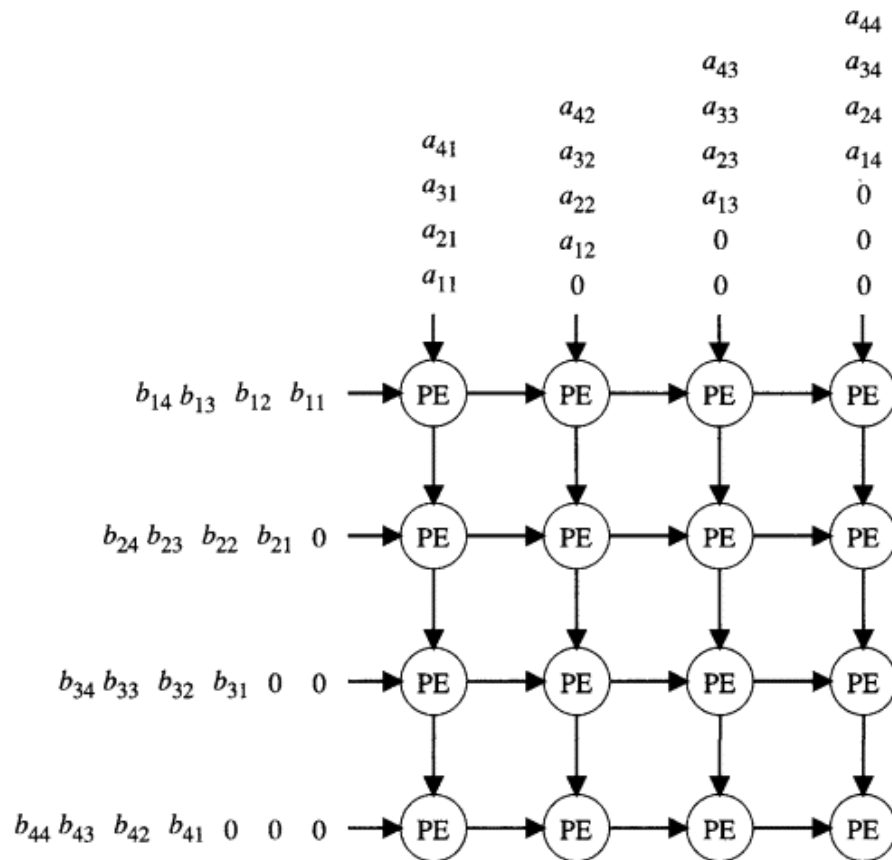By Krithika Kandasamy, Ram Jaswanth and Disha Gulur

Georgia Tech

# Problem Statement

- The sequential implementation of matrix multiplication is a very time-consuming task especially for large matrices.

- Systolic Array implementation of matrix multiplication is one approach to gain computational speed-up.

- Our project aims at accelerating matrix multiplication operation of large input matrices – 512x512 matrix multiplication.

# Systolic Array Design

## Design on paper



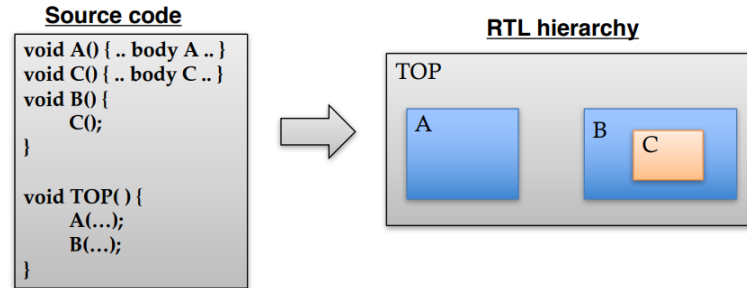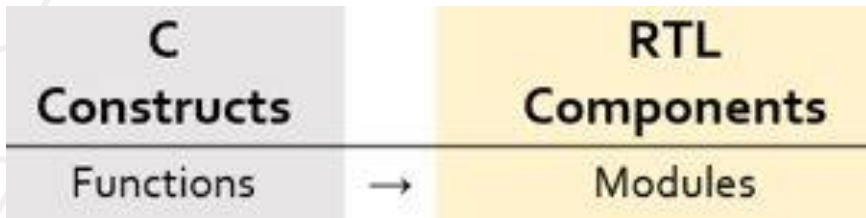Reference- *Evolutionary Mapping Techniques for Systolic Computing System*

## HLS Implementation?

- How to implement the processing elements?

- How is data transferred between the processing elements?

- How to overlap the read, write and compute operations of the processing element?
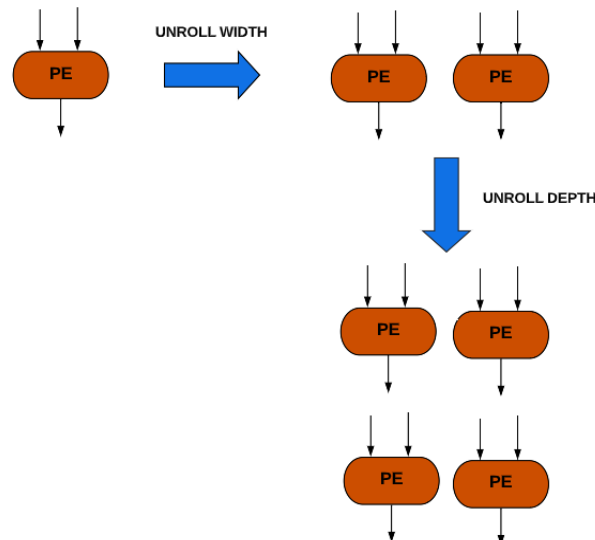
# HLS Implementation
## How to implement the processing elements?

- C functions translate to modules.



- Unrolling is scaling!



Has more meaning in FPGA than instruction-based architectures – Every "instruction" is a separate hardware.
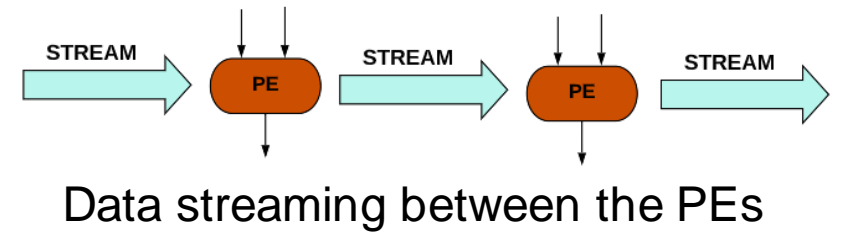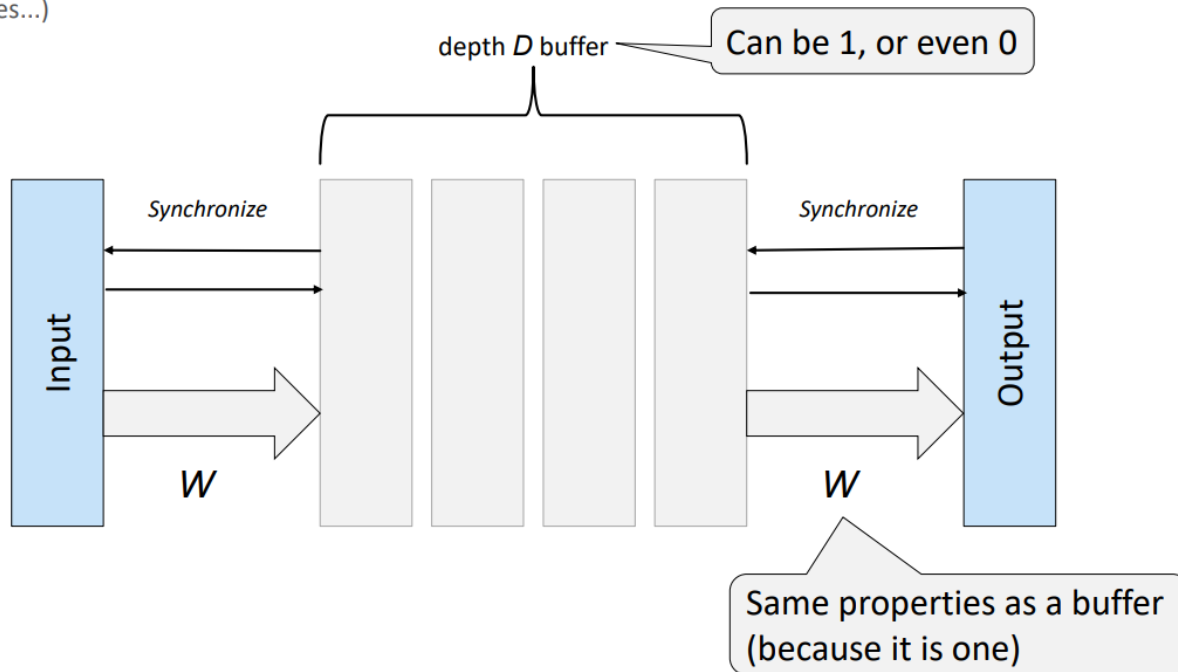
# HLS Implementation
## Data transfer between the processing elements

- Synchronous interface
- **Asynchronous interface**
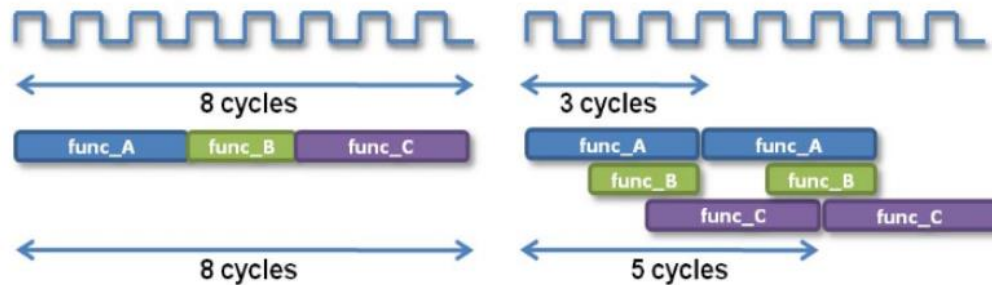


Data streaming between the PEs

# HLS Implementation
## How to overlap execution in the processing element?
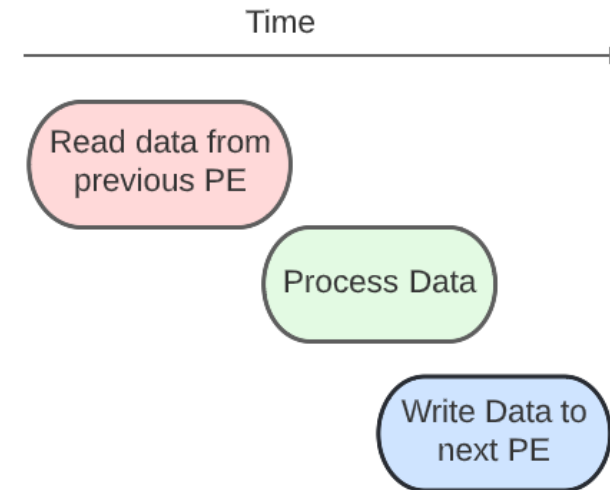
- Data flow to the rescue!

```
void top (a,b,c,d) {
...
  func_A(a,b,i1);        func_A
  func_B(c,i1,i2);       func_B
  func_C(i2,d)           func_C

  return d;
}
```

Inside the processing element

Time

Read data from previous PE

Process Data

Write Data to next PE

8 cycles

func_A  func_B  func_C

8 cycles

(A) Without Dataflow Pipelining

3 cycles

func_A  func_A
  func_B  func_B
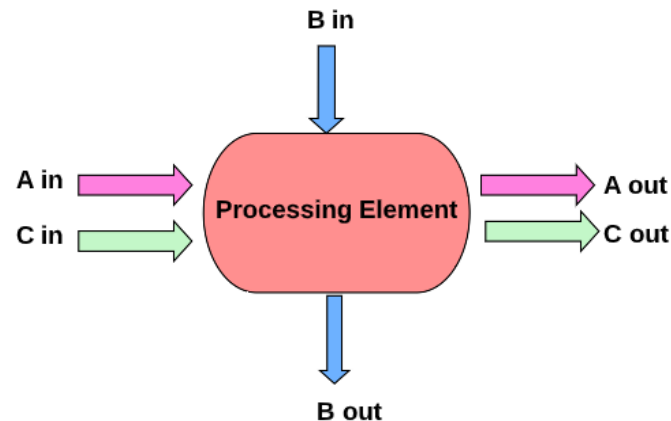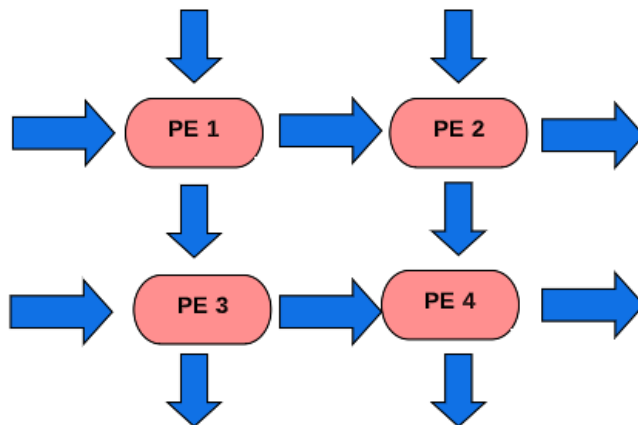    func_C    func_C

5 cycles

(B) With Dataflow Pipelining

Pipelining in the function level

# HLS Implementation

```
ComputeLoop:
for (int i = 0; i < N; i++)
  {
    #pragma HLS UNROLL

    for (int j = 0; j < N; j++)
      {
        #pragma HLS UNROLL
        //Creates nine instances of processing elements in hardware (Considering N=3)
        //Format is ProcessingElement(a_in, a_out, b_in, b_out, c_in, c_out)
        ProcessingElement(a_pipes[i][j], a_pipes[i][j + 1], b_pipes[i][j], b_pipes[i + 1][j], c_pipes[i][j], c_pipes[i][j+1], i, j);
      }
  }
```
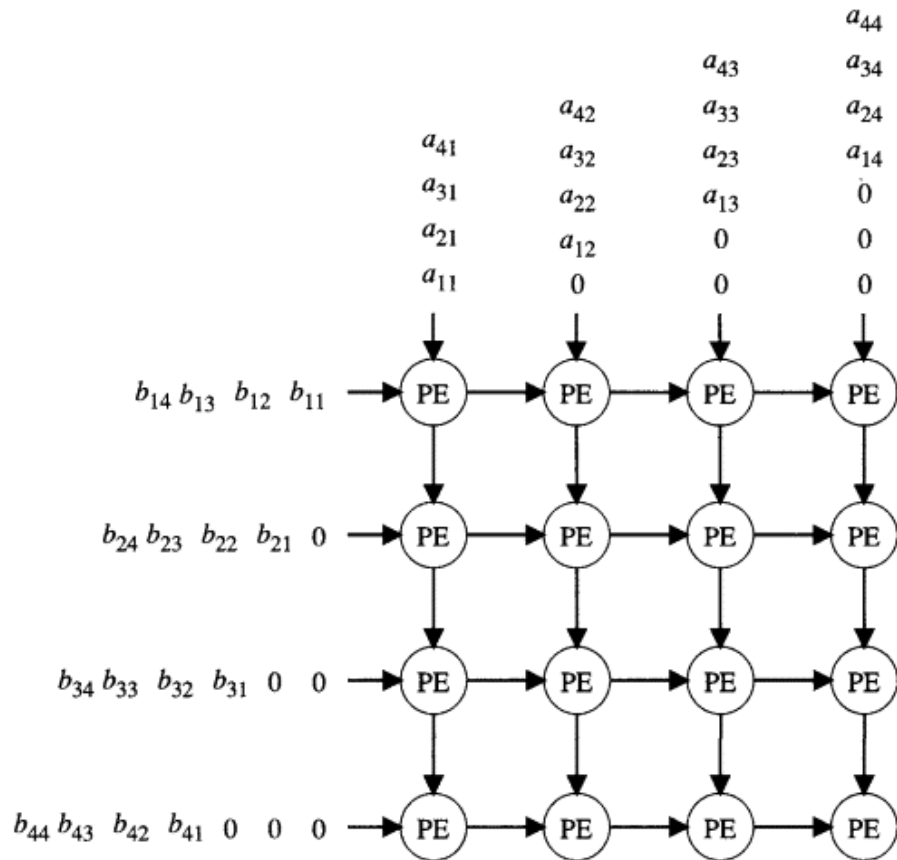
- Each PE is a function.
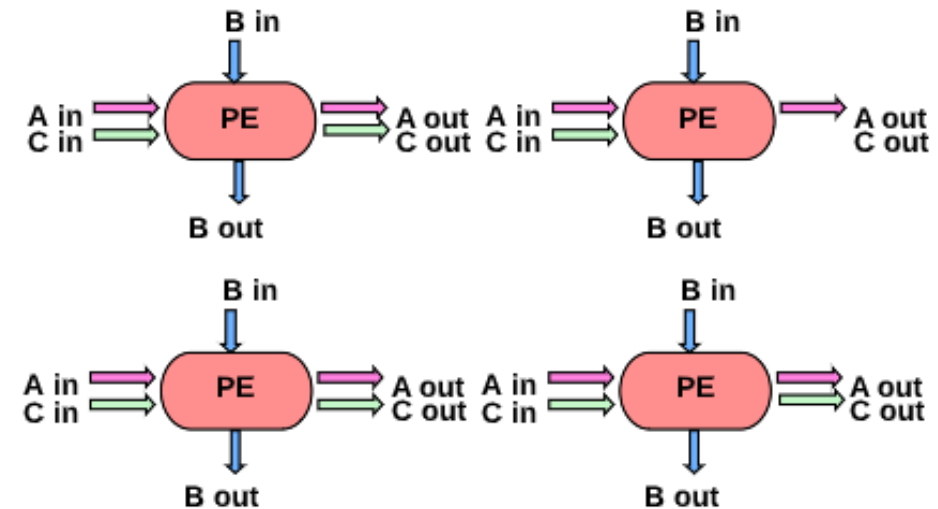- 2D array of streams for data transfer.

# Systolic Array Design – Putting it all together

## Design on paper (4x4 systolic array)



Reference- *Evolutionary Mapping Techniques for Systolic Computing System*

## HLS Implementation (2x2 systolic array)

# Tiling Matrix Multiplication

- Output stationary!



Matrix $A$ · Matrix $B$ = Matrix $C$

tile size $T$

matrix size $N$

Outer loop over tiles · Inner loop over elements · Temporary result tile

Current tile in outer loop · Current element in inner loop

# BASELINE

```
455  void matrix_mul(int MatA_DRAM[M][N], int MatB_DRAM[N][K], int MatC_DRAM[M][K])
456  {
457  #pragma HLS interface m_axi depth = 100000 port = MatA_DRAM offset = slave bundle = memA
458  #pragma HLS interface m_axi depth = 100000 port = MatB_DRAM offset = slave bundle = memB
459  #pragma HLS interface m_axi depth = 100000 port = MatC_DRAM offset = slave bundle = memC
460  #pragma HLS interface s_axilite port = return
461      int MatA[TILE][TILE];
462      int MatB[TILE][TILE];
463      int MatC[TILE][TILE];
464      int MatC_partial[TILE][TILE];
465      InitC(MatC);
466      for (int mt = 0; mt < M; mt = mt + TILE)
467      {
468          for (int kt = 0; kt < K; kt = kt + TILE)
469          {
470            //  InitC(MatC);
471              for (int nt = 0; nt < N; nt = nt + TILE)
472              {
473                  ReadInputs(MatA_DRAM, MatB_DRAM, MatA, MatB, mt, nt, kt);
474                  systolic_array(MatA, MatB, MatC_partial);
475                  addC(MatC, MatC_partial);
476              }
477              writeC(MatC, MatC_DRAM, mt, kt);
478          }
479      }
480  }
```

## LATENCY --------> 1.8 seconds

```
void ProcessingElement(hls::stream<int>& a_in, hls::stream<int>& a_out, hls::stream<int>& b_in, hls::stream<int>& b_out,
    hls::stream<int>& c_in, hls::stream<int>& c_out, int pe_row, int pe_col)
{
PE_Compute:
    int a_buffer, b_buffer;
    int c_buffer = 0;
    int count = 0;
    int c_prev;
    if (pe_col == 0)
    {
        c_prev = c_in.read();
    }
    // 3*N-2
    for (int k = 0; k < pe_row + pe_col + TILE; k++)
    {
#pragma HLS PIPELINE II = 1
        // Read a_buffer value from a_pipe_in
        if (pe_col == 0)
        {
            if (a_in.empty())
            {
                a_buffer = 0;
            }
            else if (count < pe_row)
            {
                a_buffer = 0;
                a_out.write(a_buffer);
            }
            else
            {
                a_buffer = a_in.read();
                a_out.write(a_buffer);
            }
```

```
            c_buffer += a_buffer * b_buffer;
            count++;
            if (!c_in.empty())
            {
                c_out.write(c_in.read());
            }
        }
    }
    c_out.write(c_buffer);
} // End of processing element function
```

- Each PE element performs a matrix multiplication and accumulation operation.

- Within the PE these operations are pipelined across different cycles.

# Optimization – II
## Pipeline + Loop Fusion

```
void ReadInputs(int MatA_DRAM[][N], int MatB_DRAM[][K], int MatA[][TILE], int MatB[][TILE], int m0, int n0, int k0)
{
#pragma HLS pipeline off
    for (int i = 0; i < TILE; i++)
    {
#pragma HLS LOOP_TRIPCOUNT min = 16 max = 16
        for (int j = 0; j < TILE; j++)
        {
#pragma HLS PIPELINE
#pragma HLS LOOP_TRIPCOUNT min = 16 max = 16
            MatA[i][j] = MatA_DRAM[i + m0][j + n0];
            MatB[i][j] = MatB_DRAM[i + n0][j + k0];
        }
    }
}
```

- Pipelining is also performed during memory access operations
- Since we use 512 x 512 matrices both the arrays can be accessed in the same loop

**LATENCY --------> 0.348 sec**

Georgia Tech

# Optimization – III
# Ping Pong Buffers



```
449  void matrix_mul(int MatA_DRAM[M][N], int MatB_DRAM[N][K], int MatC_DRAM[M][K])
450  {
451  #pragma HLS interface m_axi depth = 100000 port = MatA_DRAM offset = slave bundle = memA
452  #pragma HLS interface m_axi depth = 100000 port = MatB_DRAM offset = slave bundle = memB
453  #pragma HLS interface m_axi depth = 100000 port = MatC_DRAM offset = slave bundle = memC
454  #pragma HLS interface s_axilite port = return
455      int MatAO_pong[TILE][TILE];
456      int MatBO_pong[TILE][TILE];
457      int MatCO_pong[TILE][TILE];
458      int MatAO_ping[TILE][TILE];
459      int MatBO_ping[TILE][TILE];
460      int MatCO_ping[TILE][TILE];
461      int MatC[TILE][TILE];
462  #pragma HLS array_partition variable=MatC dim=1 complete
463  #pragma HLS array_partition variable=MatC dim=2 complete
464  #pragma HLS array_partition variable=MatCO_ping type= complete
465  #pragma HLS array_partition variable=MatCO_pong type= complete
466      InitC(MatC);
```

```
ReadInputs(MatA_DRAM, MatB_DRAM, MatAO_ping, MatBO_ping, mt, 0, kt);

for (int nt = 0; nt < N; nt = nt + TILE)
{
    if (nt % 2 == 0)
    {
        //Execute from ping buffer
        systolic_array(MatAO_ping, MatBO_ping, MatCO_ping);

        addC(MatC, MatCO_ping);
        //Read pong buffer
        if (nt +3*TILE < N )
        {
            ReadInputs(MatA_DRAM, MatB_DRAM, MatAO_pong, MatBO_pong, mt, nt + 2*TILE, kt);
        }
    }
    else
    {
        //Execute from pong buffer
        systolic_array(MatAO_pong, MatBO_pong, MatCO_pong);

        addC(MatC, MatCO_pong);
        //Read ping buffer
        if (nt +3*TILE < N )
        {
            ReadInputs(MatA_DRAM, MatB_DRAM, MatAO_ping, MatBO_ping, mt, nt + 2 * TILE, kt);
        }
    }
}
writeC(MatC, MatC_DRAM, mt, kt);
```

- Ping Pong Buffers are used to overlap read operation(for the next tile) and systolic array (current tile) operation.

**LATENCY --------> 0.269 sec**

Georgia Tech

- Start two systolic arrays at the same time.

- Integrate Ping Pong buffer operation.

**LATENCY --------> 0.222 sec (2 Systolic Arrays)**

# Optimization – V
## 4 Parallel Systolic Arrays

```
int MatC[TILE][TILE];
InitC(MatC);
for (int mt = 0; mt < M; mt = mt + TILE)
{
    for (int kt = 0; kt < K; kt = kt + TILE)
    {
        //  InitC(MatC);
        for (int nt = 0; nt < N; nt = nt + 4*TILE)
        {

            ReadInputs(MatA_DRAM, MatB_DRAM, MatA0, MatB0, mt, nt , kt);
            ReadInputs(MatA_DRAM, MatB_DRAM, MatA1, MatB1, mt, nt + TILE, kt);
            ReadInputs(MatA_DRAM, MatB_DRAM, MatA2, MatB2, mt, nt + 2*TILE, kt);
            ReadInputs(MatA_DRAM, MatB_DRAM, MatA3, MatB3, mt, nt + 3* TILE, kt);

            systolic_array(MatA0, MatB0, MatC0);
            systolic_array(MatA1, MatB1, MatC1);
            systolic_array(MatA2, MatB2, MatC2);
            systolic_array(MatA3, MatB3, MatC3);

            addC(MatC, MatC0, MatC1, MatC2, MatC3);
        }
        writeC(MatC, MatC_DRAM, mt, kt);
    }
}
```

- Start 4 systolic_arrays at the same time.

- Removed the Ping Pong Buffer implementation as the resource utilization is above 100%

**LATENCY --------> 0.189 sec (4 Systolic Arrays)**

# Results

| Modules & Loops | Issue Type | Violation Type | Distance | Slack | Latency(cycles) | Latency(ns) |
|---|---|---|---|---|---|---|
| ▼ ● matrix_mul | | | | - | 22177093 | 2.220E8 |
| ▶ ● matrix_mul_Pipeline_VITIS_LOOP_414_1_VITIS_LOOP_416_2 | | | | - | 258 | 2.580E3 |
| ▶ ● ReadInputs | | | | - | 465 | 4.650E3 |
| ▶ 🔷 systolic_array | | | | - | 322 | 3.220E3 |
| ▶ 🔷 systolic_array | | | | - | 322 | 3.220E3 |
| ▶ ● writeC | | | | - | 449 | 4.490E3 |
| ● addC | | | | - | 8 | 80.000 |
| ▶ 🔁 VITIS_LOOP_475_1 | | | | - | 22176832 | 2.220E8 |

| Name | BRAM_18K | DSP | FF | LUT | URAM |
|---|---|---|---|---|---|
| DSP | - | - | - | - | - |
| Expression | - | - | 0 | 80 | - |
| FIFO | - | - | - | - | - |
| Instance | 12 | 1536 | 346892 | 360786 | 0 |
| Memory | 8 | - | 2048 | 256 | 0 |
| Multiplexer | - | - | - | 6497 | - |
| Register | - | - | 24865 | - | - |
| Total | 20 | 1536 | 373805 | 367619 | 0 |
| Available | 4032 | 9216 | 2592000 | 1296000 | 960 |
| Available SLR | 1344 | 3072 | 864000 | 432000 | 320 |
| Utilization (%) | ~0 | 16 | 14 | 28 | 0 |
| Utilization SLR (%) | 1 | 50 | 43 | 85 | 0 |

➢ Memory access is the bottleneck!

# Results

PEs in the initial columns

| | | | |
|---|---|---|---|
| ▼ ⚛ systolic_array | - | **322** | 3.220E3 |
| ▶ ◉ read_inputAB11 | - | **258** | 2.580E3 |
| ◉ read_inputC12 | - | 0 | 0.0 |
| ▶ ◉ ProcessingElement13 | - | 22 | 220.000 |
| ▶ ◉ ProcessingElement14 | - | 22 | 220.000 |
| ▶ ◉ ProcessingElement29 | - | 23 | 230.000 |
| ▶ ◉ ProcessingElement15 | - | 23 | 230.000 |

PEs in the final rows and columns

| | | | | | |
|---|---|---|---|---|---|
| ▶ ◉ ProcessingElement26 | | 45 | 450.000 | | 45 |
| ▶ ◉ ProcessingElement25 | - | 50 | 500.000 | - | 50 |
| ▶ ◉ ProcessingElement26 | - | 50 | 500.000 | - | 50 |
| ▶ ◉ ProcessingElement26 | - | **51** | 510.000 | - | 51 |
| ▶ ◉ write_resultC269 | - | **258** | 2.580E3 | - | 258 |

Georgia Tech

# Results

```verilog
matrix_mul_ProcessingElement131 ProcessingElement131_U0(
    .ap_clk(ap_clk),
    .ap_rst(ap_rst),
    .ap_start(ProcessingElement131_U0_ap_start),
    .ap_done(ProcessingElement131_U0_ap_done),
    .ap_continue(ProcessingElement131_U0_ap_continue),
    .ap_idle(ProcessingElement131_U0_ap_idle),
    .ap_ready(ProcessingElement131_U0_ap_ready),
    .a_in_07125_dout(a_pipes_V_7_6_dout),
    .a_in_07125_empty_n(a_pipes_V_7_6_empty_n),
    .a_in_07125_read(ProcessingElement131_U0_a_in_07125_read),
    .a_out_022390_din(ProcessingElement131_U0_a_out_022390_din),
    .a_out_022390_full_n(a_pipes_V_7_7_full_n),
    .a_out_022390_write(ProcessingElement131_U0_a_out_022390_write),
    .b_in_037638_dout(b_pipes_V_7_6_dout),
    .b_in_037638_empty_n(b_pipes_V_7_6_empty_n),
    .b_in_037638_read(ProcessingElement131_U0_b_in_037638_read),
    .b_out_054908_din(ProcessingElement131_U0_b_out_054908_din),
    .b_out_054908_full_n(b_pipes_V_8_6_full_n),
    .b_out_054908_write(ProcessingElement131_U0_b_out_054908_write),
    .c_in_0691155_dout(c_pipes_V_7_6_dout),
    .c_in_0691155_empty_n(c_pipes_V_7_6_empty_n),
    .c_in_0691155_read(ProcessingElement131_U0_c_in_0691155_read),
    .c_out_0841412_din(ProcessingElement131_U0_c_out_0841412_din),
    .c_out_0841412_full_n(c_pipes_V_7_7_full_n),
    .c_out_0841412_write(ProcessingElement131_U0_c_out_0841412_write)
);
```

*A simple instantation of a PE module in the matrix_mul.v module (Verilog)*

Georgia Tech

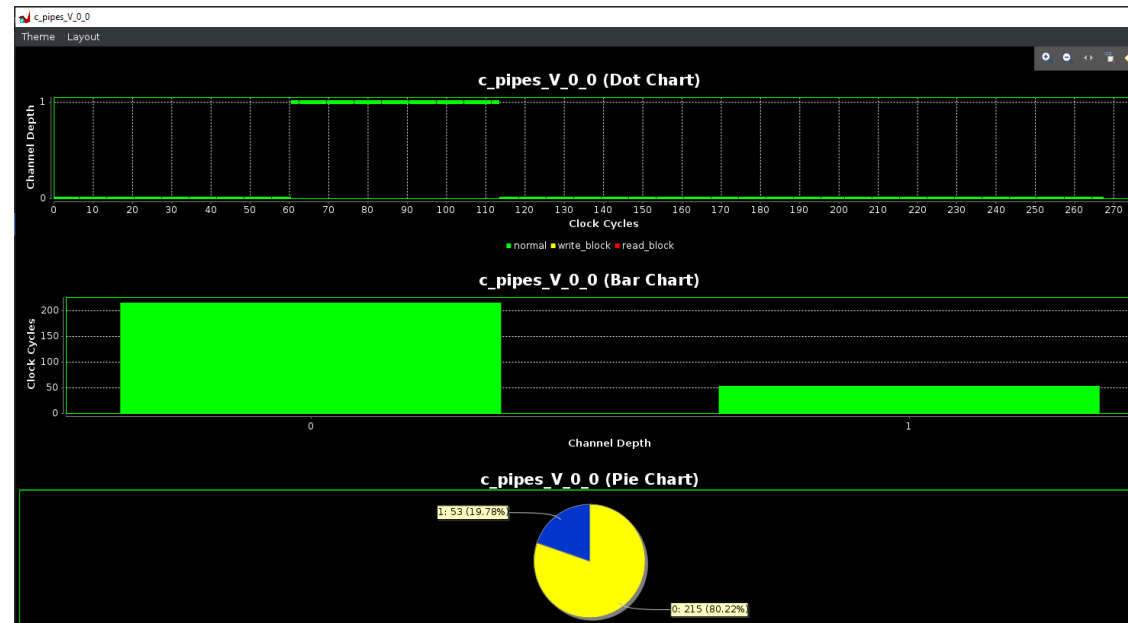# FPGA Implementation

- Implemented 3x3 systolic array



Data flow viewer

# FPGA Implementation

## Data flow properties

Waveform viewer for one of the PEs

Activity time Of the PEs

# FPGA Implementation  - Points to note

- Each stream should have one consumer and one producer.
- Insufficiently sized FIFOs can cause deadlocks



- Blocking statement calls – Stall the flow (Can check the Cosim stall time)

# Scope for future work

- Address the memory accessing bottleneck using techniques like **data packing, ping pong buffers, etc.**

- Try implementing a larger systolic array to increase the resource utilization as it is available.

- Scale our design to multiply **asymmetric matrices.**

- We would like to extend our systolic array implementation to convolution as a follow-on work.

Georgia Tech