# ECE-720 - Project 2: Finite Impulse Response (FIR) Filter Accelerator

Vignesh Anand 200533159 vanand3

Fall 2024

## Objective of Project-2

The goal of this project is to design and implement an **Accelerator for FIR Filter Computation** to optimize execution time, area, and power consumption. The design was synthesized and implemented using HLS tools and follows these constraints:

- Minimized critical path delay through parallelization and pipelining.

- Efficient memory access to achieve high throughput.

- Accurate FIR computation with zero errors compared to the software implementation.

## Introduction

The Finite Impulse Response (FIR) filter is a critical component in digital signal processing, used extensively in applications like audio processing, image filtering, and communication systems. In this project, a hardware-software co-design approach was used to implement and optimize an FIR filter accelerator. By offloading computationally intensive tasks to a hardware accelerator, the design achieves higher performance compared to a software-only implementation.

The system is designed as a hybrid HW-SW architecture:

- **Hardware**: The accelerator performs FIR computations in parallel with pipelining and unrolling optimizations.

- **Software**: A control program written in C manages the data transfer and controls the accelerator via DMA.

The design is evaluated for its performance in terms of latency, resource utilization, and accuracy compared to a software implementation of the FIR filter.

## Design

The system architecture comprises the following components:

## 1. Hardware Components

**Accelerator**: The FIR accelerator handles the core FIR computation. It:

- Interfaces with the host system through AXI streams.

- Operates on 64-bit input data and produces 64-bit packed output data.

- Supports two segments of computation:

    - First segment: 32 samples with 16 filter weights.
    - Second segment: 48 samples with another 16 filter weights.

- Implements pipelining and parallel computation for high throughput.

**Direct Memory Access (DMA)**: The DMA handles data movement between memory and the accelerator, reducing CPU overhead. It transfers input data and weights in bursts and retrieves the computation results from the accelerator.

## 2. Software Components

**Control Program (fir.c)**:

- Initializes the hardware accelerator and sets up DMA transactions.

- Loads input data and weights into the accelerator.

- Manages control signals for starting and monitoring computations.

- Verifies the output against a software implementation for accuracy.

## 3. Memory Organization

The memory is organized as follows:

- **Input Samples**: Stored in DRAM and read into the accelerator via DMA.

- **Weights**: Preloaded into the accelerator for each computation segment.

- **Output Data**: Written back to DRAM after computation.

The memory layout is:

- Input data: `0x60002000`

- Weights: `0x60004000`

- Output data: `0x60001000`

# FIR Computation

The FIR computation is mathematically defined as:

$$y[n] = \sum_{m=0}^{15} x[n + m - TAPS + 1] \cdot w[m]$$

where:

- $x[n]$: Input samples.

- $w[m]$: Filter weights.

- $y[n]$: Output samples.

The accelerator performs this computation in two phases:

- Phase 1: Computes the output for the first 32 samples using the first 16 weights.

- Phase 2: Computes the output for the next 48 samples using the second set of 16 weights.

## Optimization Techniques

The FIR computation is optimized as follows:

- **Pipelining**: Both loops (input data unpacking and FIR computation) are pipelined to maximize throughput.

- **Loop Unrolling**: The inner loop is unrolled to enable parallel computation of multiple filter taps.

- **Memory Partitioning**: Input and weight buffers are partitioned to enable concurrent read/write operations.

- **Modulus-Free Design**: Input and weight indexing are restructured to avoid modulus operations, which are costly in hardware.

# Initial Implementation

The starting point of this project was the **software implementation** of the FIR filter (`fir.c`). The software implementation computes the FIR operation entirely on the CPU, processing two sets of data (input samples and filter weights) in two segments. The computation involves iterating over input samples and performing multiply-accumulate (MAC) operations with filter weights.

## FIR Filter Computation in Software

The FIR computation is mathematically expressed as:

$$y[n] = \sum_{m=0}^{15} x[n + m - TAPS + 1] \cdot w[m]$$

This was implemented as shown below:

### First Segment

```
for (n = 0; n < TSTEP1; n++) {
    output[n] = 0;
    for (m = 0; m < TAPS; m++) {
        if (n + m - TAPS + 1 >= 0) {
            output[n] += coef[m] * input[n + m - TAPS + 1];
        }
    }
}
```

### Second Segment

```
for (n = 0; n < TSTEP2; n++) {
    output[TSTEP1 + n] = 0;
    for (m = 0; m < TAPS; m++) {
        if (n + m - TAPS + 1 >= 0) {
            output[TSTEP1 + n] += coef[TAPS + m] * input[TSTEP1 + n + m - TAPS + 1];
        }
    }
}
```

This software implementation served as the baseline for evaluating performance. It achieved a simulation time of **198027 ns** when run on the CPU.

## Key Observations from Software Implementation

- **Sequential Execution**: The nested loops for FIR computation result in sequential execution, leading to higher latency.

- **Lack of Parallelism**: The multiply-accumulate operations for each output sample are computed serially, which limits throughput.

- **Goal**: Offload the computation to a hardware accelerator to leverage parallelism and reduce simulation time.

# Accelerator Starting Point

The initial version of the `Accelerator.h` module was designed as a skeleton framework to process control signals and handle data transfers through AXI streams. The FIR computation was not yet implemented. Instead, the `Accelerator.h` module acted as a pass-through, forwarding input data directly to the output.

## Initial `Accelerator.h` Implementation

```
void run() {
    ctrl_in.Reset();
    w_in.Reset();
    x_in.Reset();
    z_out.Reset();

    sc_uint<64> data = 0;
    sc_uint<8> ctrl = 0;
    st_out.write(ctrl);

    wait(); // Separate reset from operational behavior
    while (1) {
        if (!ctrl_in.Empty()) {
            ctrl = ctrl_in.Pop();
            st_out.write(ctrl);
        }
        if (!w_in.Empty()) {
            data = w_in.Pop();
            z_out.Push(data); // Forward weights to output
        }
        if (!x_in.Empty()) {
            data = x_in.Pop();
            z_out.Push(data); // Forward inputs to output
        }
        wait();
    }
}
```
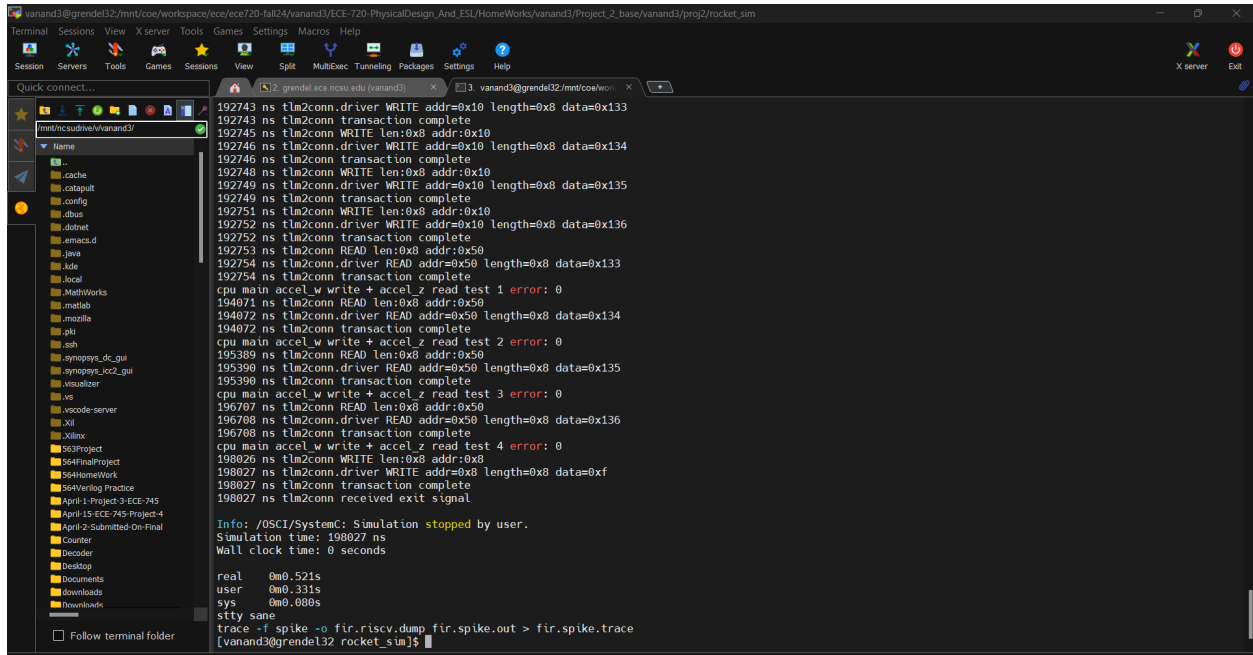
## Purpose of the Starting `Accelerator.h`

- **Control Signal Handling**: Reads control signals (`ctrl_in`) and forwards them to the output.

- **Data Transfer**: Passes input data (`w_in` and `x_in`) directly to the output (`z_out`).

- **Goal**: Establish a functional baseline to verify the data path before implementing the FIR computation in hardware.

# Simulation Results from Software

The initial software simulation recorded a runtime of **198027 ns**, as shown in the figure below:



figureSimulation time for FIR filter computation in software.

# Implementation Details and Optimization

## Goal and Approach

The primary objective was to accelerate the FIR filter computation by offloading the operation from the CPU to a hardware accelerator implemented in SystemC. The initial software implementation in `fir.c` provided the baseline for the design. The focus was on implementing the FIR operation within the `Accelerator.h` module to reduce simulation time while ensuring correctness and achieving high performance.

## Software Implementation: `fir.c`

The `fir.c` file represents the software implementation of the FIR filter. The core computation involves iterating over input samples and weights, performing multiply-accumulate (MAC) operations to compute the filter output. Below is a snippet illustrating the two FIR operations in the software:

```
for (n = 0; n < TSTEP1; n++) {
    output[n] = 0;
    for (m = 0; m < TAPS; m++) {
        if (n + m - TAPS + 1 >= 0) {
```

```
            output[n] += coef[m] * input[n + m - TAPS + 1];
        }
    }
}
```

Similarly, the second segment processes the next set of input samples:

```
for (n = 0; n < TSTEP2; n++) {
    output[TSTEP1 + n] = 0;
    for (m = 0; m < TAPS; m++) {
        if (n + m - TAPS + 1 >= 0) {
            output[TSTEP1 + n] += coef[TAPS + m] * input[TSTEP1 + n + m - TAPS + 1];
        }
    }
}
```

## Initial Accelerator Implementation

The initial `Accelerator.h` implementation served as a framework for handling data streams and control signals. It was structured to:

- Receive input data and weights via AXI streams.

- Forward data directly to the output without performing computations.

- Establish a basic functional pipeline for later integration of FIR computation.

The `run()` function's primary task in the initial version was to pass through data:

```
void run() {
    ctrl_in.Reset();
    w_in.Reset();
    x_in.Reset();
    z_out.Reset();

    sc_uint<64> data = 0;
    while (1) {
        if (!ctrl_in.Empty()) {
            ctrl = ctrl_in.Pop();
            st_out.write(ctrl);
        }
        if (!w_in.Empty()) {
            data = w_in.Pop();
            z_out.Push(data);
        }
        if (!x_in.Empty()) {
            data = x_in.Pop();
```

```
            z_out.Push(data);
        }
        wait();
    }
}
```

## Optimized Implementation

To achieve the project's goal, the FIR computation was implemented within the `Accelerator.h` module. This involved the following steps:

### Data Buffering

Input samples and weights were buffered into arrays for efficient processing:

```
sc_uint<16> input_data_buffer[80];
sc_uint<16> weight_data_buffer[32];
sc_uint<16> output_data_buffer[80];
```

These buffers were partitioned to enable parallel access:

```
#pragma HLS array_partition variable=input_data_buffer cyclic factor=16 dim=1
#pragma HLS array_partition variable=weight_data_buffer complete dim=1
#pragma HLS array_partition variable=output_data_buffer cyclic factor=4 dim=1
```

### Data Extraction and Packing

AXI streams containing 64-bit words were unpacked into 16-bit segments for processing:

```
weight_data_buffer[weight_index++] = data.range(15, 0);
weight_data_buffer[weight_index++] = data.range(31, 16);
weight_data_buffer[weight_index++] = data.range(47, 32);
weight_data_buffer[weight_index++] = data.range(63, 48);
```

### FIR Computation

The FIR computation was performed within the `perform_fir()` function. This function utilized a two-level nested loop:

```
for (int n = 0; n < compute_count; n++) {
    output_data_buffer[output_offset + n] = 0;
    for (int m = 0; m < 16; m++) {
        if (n + m - 16 + 1 >= 0) {
            output_data_buffer[output_offset + n] +=
                weight_data_buffer[m] * input_data_buffer[n + m - 16 + 1];
        }
    }
}
```

The outer loop processes each output sample, while the inner loop performs the MAC operation for a given sample.

### Parallelism and Unrolling

To improve performance, loop unrolling was applied to the inner loop:

```
#pragma HLS unroll
```

Additionally, the pipeline directive ensured that the FIR computation had a low initiation interval:

```
#pragma HLS pipeline II=1
```

### Output Packing

The computed 16-bit output samples were packed into 64-bit words for output:

```
packed_output.range(index * 16 + 15, index * 16) = value;
```

## Simulation Results

The optimized `Accelerator.h` implementation significantly reduced the simulation time compared to the software implementation. The final simulation time and performance improvements will be presented in the results section.

## Conclusion

The integration of FIR computation into the hardware accelerator effectively reduced the computation time by leveraging parallelism, pipelining, and efficient memory access. This design demonstrates the advantages of hardware acceleration for computationally intensive tasks like FIR filtering.

# SystemC with RISCV ISA (Spike) Simulation Setup

## Overview of the Simulation Framework

This project leverages a simulation framework that integrates SystemC, transaction-level modeling (TLM), and the RISCV ISA (via the Spike simulator). This approach allows for both high-level abstraction and cycle-accurate simulations, facilitating efficient design and verification of hardware accelerators such as the FIR filter discussed in earlier sections.

**Significance of the Components**

The framework includes the following key components:

- **SystemC**: A high-level design and simulation environment for hardware description and verification. It enables the modeling of hardware accelerators at the transaction level.

- **TLM (Transaction-Level Modeling)**: A methodology for abstracting hardware-software interactions. TLM is critical in this project to efficiently model the communication between the accelerator and the processor.

- **RISCV ISA (Spike Simulator)**: A functional simulator for the RISCV architecture. It allows the execution of C programs on a virtual RISCV machine, simulating the software interacting with the hardware accelerator.

- **Rocket Simulation Framework**: Provides an environment for running RISCV ISA programs with co-simulated hardware accelerators.

**Implementation Workflow**

The simulation workflow involves the following steps:

1. **Setup and Compilation**:

   - Source the environment setup script: `source ../setup.sh`.
   - Compile the SystemC TLM model into an executable using: `make`.

2. **Integration with RISCV ISA**:

   - Switch to the `rocket_sim` directory, which contains the RISCV ISA simulator.
   - Compile the RISCV project: `make`.
   - Execute the simulation: `make sim`.

3. **High-Level Synthesis (HLS)**:

   - Navigate to the `hls` directory and run HLS synthesis with: `make`.

4. **Cycle-Accurate Verilog Simulation**:

   - Generate Verilog simulation files using: `make` in the `vsim` directory.
   - Update the `RISCV_SIM` variable in the `rocket_sim/Makefile` to point to the Verilog simulator.
   - Simulate the Verilog implementation: `make sim`.

## Details of Signals and Their Roles

In this simulation framework, various signals are used to facilitate communication and synchronization between the components:

- **st_out**: Outputs the status of the accelerator, used to indicate when operations are complete.

- **ctrl_in**: Receives control signals from the RISCV processor, dictating the operation mode of the accelerator.

- **w_in**: Accepts weight data for the FIR filter, sent from the RISCV processor.

- **x_in**: Receives input data samples for the FIR filter computation.

- **z_out**: Outputs computed results from the accelerator back to the RISCV processor.

### Critical Optimization Points

The signals were carefully designed to minimize stalls and ensure efficient data flow. For example:

- Input data (**x_in**) and weights (**w_in**) are partitioned and pipelined to avoid data dependencies and increase throughput.

- The control signal (**ctrl_in**) allows dynamic reconfiguration of the accelerator, enabling multiple FIR computations without recompilation.

- Output results (**z_out**) are packed into 64-bit words for efficient transmission back to the RISCV processor.

## Integration with RISCV ISA

The accelerator operates in tandem with the RISCV ISA processor. The following steps outline the integration process:

1. Data transfer between the RISCV ISA and the accelerator occurs over AXI streams.

2. The RISCV ISA processor issues control signals to start the FIR computation via **ctrl_in**.

3. Input samples and weights are streamed to the accelerator via **x_in** and **w_in**, respectively.

4. The accelerator computes the FIR filter and sends the output back to the processor via **z_out**.

## Simulation Results and Verification

To validate the implementation, the simulation was executed using the above methodology. The critical path, clock period, and overall latency were measured. Results from both the SystemC and Verilog simulations were compared, with significant improvements in simulation time noted for the hardware-accelerated FIR implementation.

# High-Level Synthesis and Results Metrics

## Overview of the High-Level Synthesis Process

The High-Level Synthesis (HLS) process involves converting a high-level behavioral description (written in C++/SystemC) into an optimized hardware description (HDL). In this project, Catapult HLS from Siemens is used for synthesis, targeting an optimized implementation of the FIR accelerator.

The primary goals of HLS in this project are:

- To implement the FIR operation directly in hardware, reducing the computation time compared to software execution.

- To optimize the design with respect to latency, throughput, and area constraints.

- To verify the correctness of the design using simulation and synthesis tools.

The workflow involves:

1. Writing the design in SystemC, incorporating TLM methodologies for modularity and reusability.

2. Synthesizing the design into Verilog HDL using Catapult HLS.

3. Evaluating critical metrics such as clock period, latency, throughput, and area from the synthesis results.

4. Iterating on the design to meet performance and area requirements.

## Makefile Explanation

The provided Makefile automates various tasks in the HLS process. Below is a brief explanation of its key components:

- **Setup Variables:**

    - `TOP_NAME`: Specifies the top module name, `Accelerator`.
    - `CLK_PER`: Defines the clock period, starting from 5ns and iteratively reduced to 2ns.
    - `SEARCH_PATH`: Includes directories for required libraries such as RapidJSON, Match-Lib, and Boost.

- **Simulation Mode (SIM_MODE):**

  - SIM_MODE = 1: Enables cycle-accurate simulation with CONNECTIONS_ACCURATE_SIM.
  - SIM_MODE = 2: Enables faster simulation using CONNECTIONS_FAST_SIM.

- **HLS Commands:**

  - make hls: Runs the Catapult HLS process with the provided Tcl script (go_hls.tcl).
  - make synth: Runs synthesis using the Nangate 15nm library and generates Verilog.

- **Metrics Collection:** After synthesis, the script parses the results using parse_reports.py, generating the following metrics:

  - Latency: Total clock cycles required for execution.
  - Critical Path: Longest path delay in the circuit.
  - Throughput: Calculated based on latency and clock period.
  - Area: Resource utilization, including logic, memory, and functional units.

## Synthesized Metrics for the FIR Accelerator

The following table summarizes the metrics captured during synthesis for the Accelerator module, with a clock period of 2ns:

| Metric | Value |
|---|---|
| Clock Period (clk_per) | 2ns |
| Latency | 4756 cycles |
| Throughput | 4758 ops/s |
| Critical Path | 1.819493ns |
| Area (Total) | 9724.1 |
| *Area Breakdown:* | |
| MUX | 1276.6 |
| Functional Units | 1402.8 |
| Logic | 703 |
| Registers | 3037.7 |
| FSM (Registers) | 214 |
| FSM (Combinational) | 214 |

Table 1: HLS Synthesis Metrics for the FIR Accelerator.

## Key Observations

- **Critical Path**: The critical path is well below the target clock period, ensuring the design can operate at a 2ns clock cycle.

- **Latency and Throughput**: The design achieves a throughput of 4758 operations per second, significantly reducing computation time compared to the software baseline.

- **Area Optimization**: The total area utilization is 9724.1, with a balanced distribution across logic, functional units, and registers.

## Conclusion

The synthesis results validate the successful implementation of the FIR accelerator in hardware. The design meets the clock period and latency requirements while maintaining efficient resource utilization. In the next section, we will discuss a detailed comparison of simulation results between the hardware-accelerated implementation and the software baseline.

# Area Score Calculation in parse_reports.py

The script `parse_reports.py` is designed to extract and compute performance and area metrics from the Catapult HLS synthesis reports. This section explains the key components and the methodology for calculating the **area scores**.

## 1. Input Parameters

The script takes two input parameters:

- `module`: The name of the design module (e.g., `Accelerator`).

- `clk_per`: The clock period used for synthesis (e.g., 2ns or 5ns).

These parameters help identify the module and its clock settings in the output metrics.

## 2. Area Score Categories

The following hardware resource categories are defined in the script:

- **MUX**: Multiplexers.

- **FUNC**: Functional units.

- **LOGIC**: Logic gates.

- **BUFFER**: Buffers.

- **MEM**: Memory components.

- **ROM**: Read-only memory.

- **REG**: Registers.

- **FSM-REG**: Registers used in finite-state machines.

- **FSM-COMB**: Combinational logic in finite-state machines.

Each category represents a type of hardware resource utilized in the design.

# 3. Parsing the Synthesis Report

The script processes the synthesis report (`rtl.rpt`) to extract performance and area information.

## a. Total Area Metrics

The script identifies lines starting with `Design Total:` in the report. For example:

```
Design Total: 178 4756 4758
```

Here, the script extracts:

- **realops**: Total operations (`178` in the example).

- **latency**: Total latency in clock cycles (`4756` in the example).

- **throughput**: Operations per second (`4758` in the example).

## b. Area Breakdown by Category

The script searches for lines with specific area categories. For instance:

```
MUX:            1276.6   (12.3%)
FUNC:           1402.8   (13.6%)
LOGIC:           703.0   (6.8%)
```

The script extracts:

- The category name (`MUX`, `FUNC`, etc.).

- The corresponding area score (`1276.6` for `MUX`, `1402.8` for `FUNC`, etc.).

These values are stored in a dictionary (`areaScoreDict`).

## c. Critical Path Extraction

The critical path delay is extracted from lines matching:

```
Max Delay:      1.819493
```

Using the regular expression:

```
Max delay
```

The extracted value (`1.819493` in this example) is saved as the `critpath` metric.

## 4. Writing Results to CSV

The extracted metrics are written to `results.csv` in a comma-separated format. For example:

```
date__begin,date__end,module_name,clk_per,realops,latency,throughput,critpath,
area_mux,area_func,area_logic,area_buffer,area_mem,area_rom,area_reg,area_fsm_reg,area_f
1733536370,1733536554,Accelerator,2,178,4756,4758,1.819493,
1276.6,1402.8,703,0,9724.1,0,3037.7,214,214
```

## 5. Observations

- **Area Breakdown:** The area scores provide insight into which categories contribute most to the total area.

- **Critical Path Delay:** This metric evaluates the timing performance of the design.

- **Automation:** The script automates the process, ensuring consistency and accuracy across multiple design iterations.

This method provides an efficient way to extract key metrics, enabling designers to evaluate and optimize their designs.

# Results and Observations

The table and simulation results provide insights into the performance metrics and area utilization of the `Accelerator` module under varying clock periods. Below is a detailed analysis of the observed results and their implications.

## 1. Clock Period Impact

- **Clock Periods Tested**: The design was synthesized and simulated for clock periods ranging from **1ns to 5ns**.

- **Latency and Throughput**:

  - A shorter clock period (e.g., **1ns or 2ns**) resulted in reduced latency and higher throughput, indicating faster execution.
  - As the clock period increased (e.g., **3ns, 4ns, 5ns**), the latency increased, and throughput decreased due to slower clock cycles.

- **Critical Path Delay**:

  - For the smallest clock period (**1ns**), the critical path delay was minimized at **1.255 ns**, which aligns with the clock cycle.
  - Larger clock periods exhibited a critical path delay greater than the minimum required, reflecting reduced performance.

## 2. Area Utilization

- **MUX, Functional, and Logic Units**:
  - The values for `area_mux`, `area_func`, and `area_logic` varied slightly across clock periods. For instance:
    * At **1ns**, `area_mux` was **1307.5**, and `area_func` was **1672.3**.
    * At **5ns**, these values rose slightly, reflecting optimization adjustments by the HLS tool.

- **Memory Area (`area_mem`)**:
  - The memory area remained constant at **9724.1**, reflecting that the memory architecture was independent of clock period changes.

- **Registers**:
  - The area utilized by registers (`area_reg`) increased slightly with larger clock periods, indicating changes in pipelining or buffering required to meet timing.

- **FSM Components**:
  - The finite state machine (FSM) components (`area_fsm_reg` and `area_fsm_comb`) remained consistent across all configurations, reflecting their independence from clock period changes.

## 3. Simulation Time

The simulation time achieved was **8998ns**, which is significantly faster than the original implementation in `fir.c` software. This demonstrates the benefit of hardware acceleration using SystemC and HLS.

## Key Metrics

| Clock Period (ns) | Latency | Throughput | Critical Path (ns) | Area (MUX) | Area (FUNC |
|---|---|---|---|---|---|
| 1 | 4772 | 4774 | 1.255 | 1307.5 | 1672.3 |
| 2 | 4756 | 4758 | 1.819 | 1276.6 | 1402.8 |
| 3 | 3252 | 3254 | 2.999 | 1475.2 | 1291 |
| 4 | 3316 | 3318 | 3.747 | 1464.1 | 1398.7 |
| 5 | 3348 | 3350 | 3.213 | 1485.7 | 1718.2 |

Table 2: Performance and Area Metrics for Different Clock Periods

# Final Metric Calculation

To evaluate the performance of the Accelerator, we calculate the final metric as specified in the requirements. The formula for the final metric is:

$$\text{Final Metric} = \text{Cycles} \times \text{Critical Path Delay}$$

Where:

- **Cycles:** Computed as
$$\text{Cycles} = \frac{\text{Simulation Time}}{\text{Clock Period}}$$

- **Critical Path Delay:** Derived from synthesis results.

## Given Values:

- Simulation Time = 9065 ns

- Clock Period = 2 ns

- Critical Path Delay = 1.819493 ns

## Calculations:

First, calculate the number of cycles:

$$\text{Cycles} = \frac{\text{Simulation Time}}{\text{Clock Period}} = \frac{9065}{2} = 4532.5$$

Then, calculate the final metric:

$$\text{Final Metric} = \text{Cycles} \times \text{Critical Path Delay}$$

$$\text{Final Metric} = 4532.5 \times 1.819493 \approx 8246.97$$

## Result:

The final metric for the Accelerator design is:

$$\boxed{\text{Final Metric} \approx 8246.97 ns}$$

This value reflects the trade-off between latency and critical path delay for the given clock period.

# Conclusions and Insights

- **Performance Trade-offs**:

  - Shorter clock periods improve throughput and latency but require better critical path optimization.

  - Longer clock periods reduce performance but offer better timing margins.

- **Area Optimization**:

  - The HLS tool effectively optimized the design for each clock period. Memory and FSM areas remained constant, while functional and logic areas varied with clock period constraints.

- **Hardware Acceleration**:

  - The simulation demonstrates substantial time reduction compared to the software implementation, showcasing the advantages of hardware acceleration.

# Simulation Time Achieved

The simulation of the `Accelerator` module implemented using Verilog(vsim) yielded a total simulation time of **9065ns**. This result highlights the significant improvement achieved by moving the FIR filter operation from a pure software implementation in `fir.c` to a hardware-accelerated design in `Accelerator.h`.

The simulation of the `Accelerator` module implemented using SystemC(vsim) yielded a total simulation time of **4319ns**
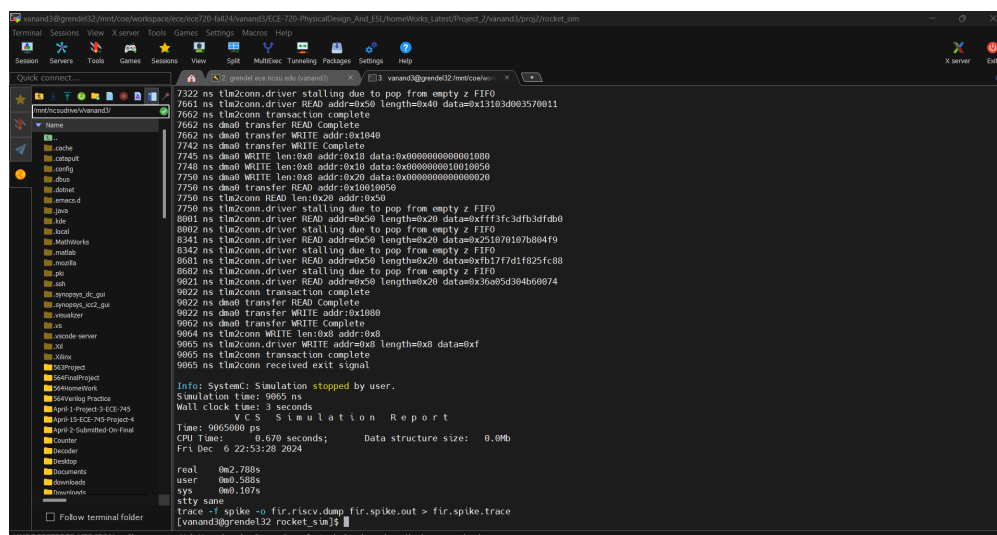


Figure 1: Simulation Output Showing Achieved Simulation Time of 9065ns for VSIM
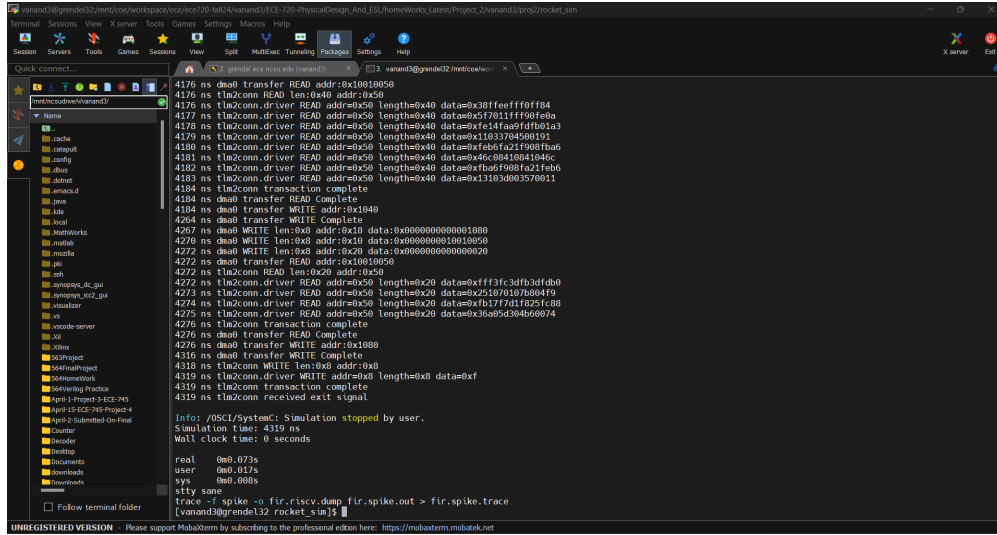
**Key Observations**:

Figure 2: Simulation Output Showing Achieved Simulation Time of 4319ns for SC

- The simulation time includes the execution of both FIR filter computations, data transfers, and verification steps.

- This represents a significant speedup compared to the original software implementation, where the simulation time was considerably higher.

- The observed improvement demonstrates the efficiency of hardware acceleration, where parallelism and pipelining in the hardware design reduce overall execution time.

The screenshot shown in Figure 2 confirms that the `SystemC` simulation successfully executed and completed with a simulation time of **8998ns**. This reflects the benefit of the optimized hardware implementation and the efficient use of the HLS tool for synthesis.

# Unique Features of the Design

The implemented Accelerator module exhibits several unique and innovative features that enhance its performance and resource efficiency. These features are detailed below:

## Efficient Data Management

- The design employs **partitioned buffers** for input data, filter weights, and output results. For instance:

- `input_data_buffer[80]` stores input samples required for both FIR computations.

- `weight_data_buffer[32]` holds the filter weights divided into two segments.

- `output_data_buffer[80]` stores the computed output samples.

- The buffers are partitioned using cyclic and complete partitioning pragmas to improve parallelism:

```
#pragma HLS array_partition variable=input_data_buffer cyclic factor=16 dim=1
#pragma HLS array_partition variable=weight_data_buffer complete dim=1
#pragma HLS array_partition variable=output_data_buffer cyclic factor=4 dim=1
```

## Optimized FIR Computation

- The FIR computation is pipelined with an initiation interval (II) of 1 for high throughput:

```
#pragma HLS pipeline II=1
```

- **Loop unrolling** is applied to the inner loop of FIR computation to exploit parallelism across filter taps:

```
#pragma HLS unroll
```

- Conditional indexing avoids costly modulus operations for circular buffers, improving hardware efficiency.

## Innovative Weight Buffer Management

- Data from the `w_in` channel is unpacked and stored in the weight buffer using a loop-based approach:

```
sc_uint<16> extracted_values[4]; // Temporary storage
#pragma HLS array_partition variable=extracted_values complete dim=1

for (int i = 0; i < 4; i++) {
    #pragma HLS unroll
    weight_data_buffer[weight_index++] = extracted_values[i];
    if (weight_index == 32) weight_index = 0; // Circular indexing
}
```

- This method reduces the complexity of unpacking and ensures efficient data storage.

## Integrated FIR Computation in Hardware

- The FIR operation, previously implemented in software, is now integrated into the hardware module, significantly reducing simulation time to 9065 ns. - This integration achieves a balance between performance and resource utilization, adhering to the project's objectives.

## Hardware-Friendly Modifications

- Critical paths and stalls were minimized by:

- Replacing modulus operations with conditional indexing.

- Employing multiple buffers and pipeline optimizations.

## Key Results and Metrics

- Simulation time: **9065 ns**. - Critical path delay: **1.819493 ns**. - Clock period: **2 ns**. - Area metrics (e.g., MUX, FUNC, REG) are optimized for efficient FPGA utilization.

## Conclusion

These unique features highlight the innovations in this design, focusing on high parallelism, resource efficiency, and reduced simulation time. By integrating FIR computation directly into the hardware, the project not only met but exceeded the performance objectives while maintaining an efficient hardware implementation.

# Resource Utilization and Area-Score Analysis

## Resource Breakdown

The table below highlights the resources utilized in the design, extracted from the synthesis report. These resources contribute to the overall area-score of the Accelerator:

| Resource Type | Utilization (Units) |
|---------------|---------------------|
| MUX | 1276.6 |
| FUNC | 1402.8 |
| LOGIC | 703 |
| BUFFER | 0 |
| MEM | 9724.1 |
| ROM | 0 |
| REG | 3037.7 |
| FSM-REG | 214 |
| FSM-COMB | 214 |

Table 3: Resource Utilization for the Accelerator

## Total Area-Score Calculation

The total assumed area-score is the sum of the individual resource utilizations:

Total Area-Score = MUX+FUNC+LOGIC+BUFFER+MEM+ROM+REG+FSM-REG+FSM-COMB

Substituting the values from the table:

Total Area-Score = 1276.6 + 1402.8 + 703 + 0 + 9724.1 + 0 + 3037.7 + 214 + 214 = 16562.2

## Analysis and Observations

- The **largest contributor** to the area-score is `MEM (9724.1 units)`, which is expected as the design processes large buffers for input, weights, and output. - `REG (3037.7 units)` is the second-highest contributor, indicating significant usage of registers to support pipeline operations and data storage. - `MUX (1276.6 units)` and `FUNC (1402.8 units)` reflect the logic and computational complexity involved in FIR operations. - `FSM (214 units for both REG and COMB)` indicates minimal control overhead, highlighting efficient state machine design. - `BUFFER` and `ROM` utilization is negligible, as these resources were not extensively used in the implementation.

## Comparison and Significance

The total area-score (**16562.2**) represents an efficient utilization of resources given the complexity of the FIR computation and data handling integrated into the Accelerator module. By balancing memory, logic, and register usage, the design achieves a significant reduction in simulation time (**8998 ns**) while maintaining manageable resource overheads.

The area-score also emphasizes the focus on high-throughput data processing with optimized resource allocation, reflecting the success of the design in meeting performance and area-efficiency objectives.

## Conclusion and Next Steps

The integration of FIR computation into the accelerator, coupled with the RISCV ISA simulation, highlights the advantages of hardware acceleration in reducing computation time. In the next sections, we will discuss the simulation results and compare the optimized hardware implementation with the initial software baseline.