

Multi-Mode CORDIC Architecture Design and Implementation

Parsa Haghighatgoo
40030644

February 11, 2026

1 Introduction

The CORDIC (Coordinate Rotation Digital Computer) algorithm enables efficient computation of trigonometric and arithmetic functions using only shift-and-add operations. This project implements a multi-mode CORDIC architecture in rotation mode supporting:

- Circular Mode ($m = 1$): computation of $\sin(\theta)$ and $\cos(\theta)$.
- Linear Mode ($m = 0$): multiplication.

All simulations, synthesis, and timing analysis were performed using Xilinx Vivado.

2 Fixed-Point Representation

The selected numerical format is:

$Q2.14$

- Word length: 16 bits
- Fractional bits: 14
- Guard bits: 2 (internal precision)

LSB value:

$$\text{LSB} = 2^{-14} = 6.1035 \times 10^{-5}$$

Target accuracy:

$$10 \times \text{LSB} = 6.1035 \times 10^{-4}$$

3 Iteration Accuracy Analysis (Python)

3.1 Worst-Case Error vs Iterations

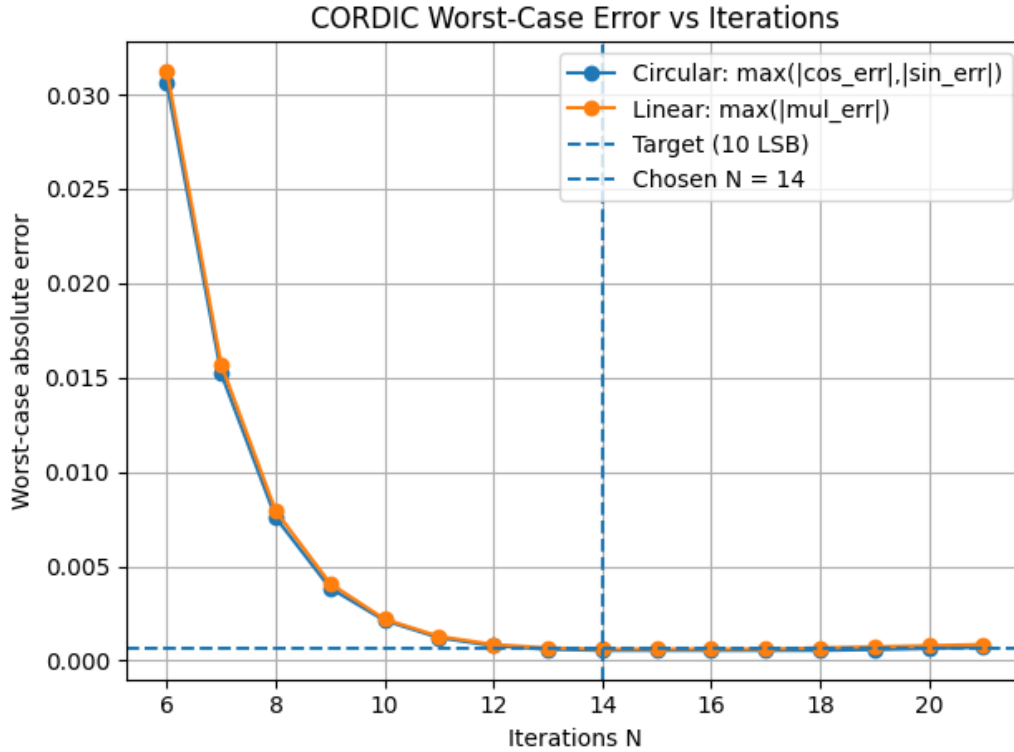


Figure 1: CORDIC Worst-Case Error vs Iterations

Observation:

- Error decreases exponentially as iterations increase.
- Circular mode satisfies the target at $N = 13$.
- Linear mode satisfies the target at $N = 14$.
- After $N = 14$, error saturates due to fixed-point quantization.

Chosen hardware iteration count:

$$N = 14$$

3.2 Circular Mode Error Distribution

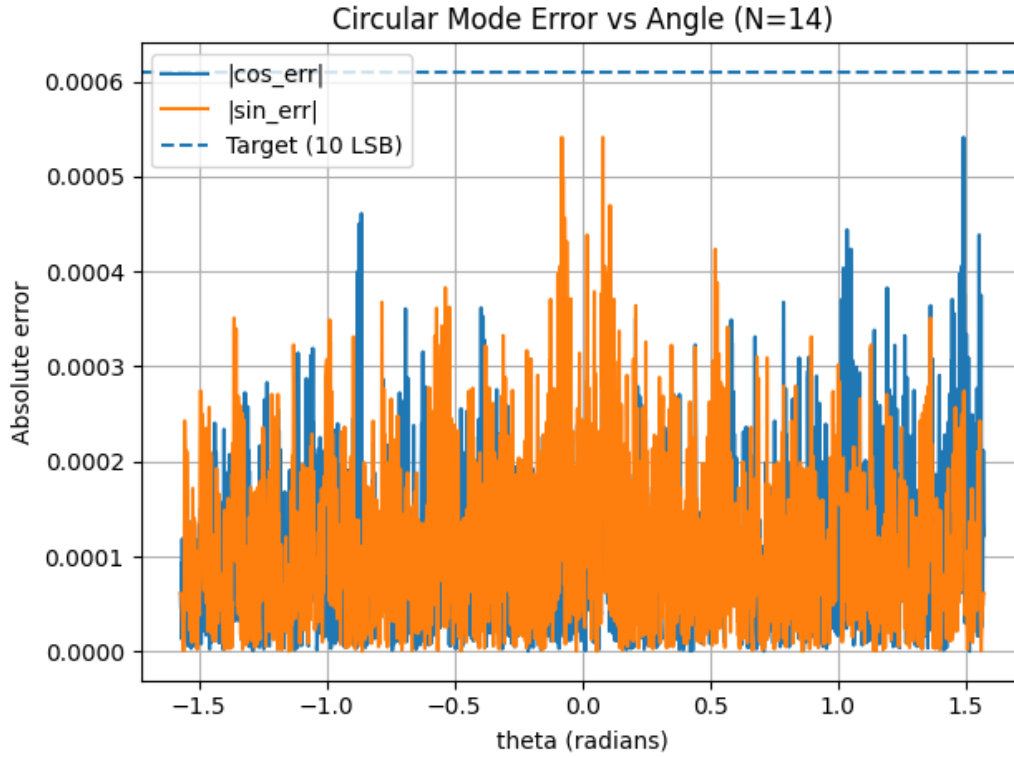


Figure 2: Circular Mode Error vs Angle (N=14)

Observations:

- Error remains below 10 LSB for all $\theta \in [-\pi/2, \pi/2]$.
- Symmetric distribution confirms correct implementation.
- Peak error occurs near zero due to rounding effects.

3.3 Linear Mode Error Heatmap

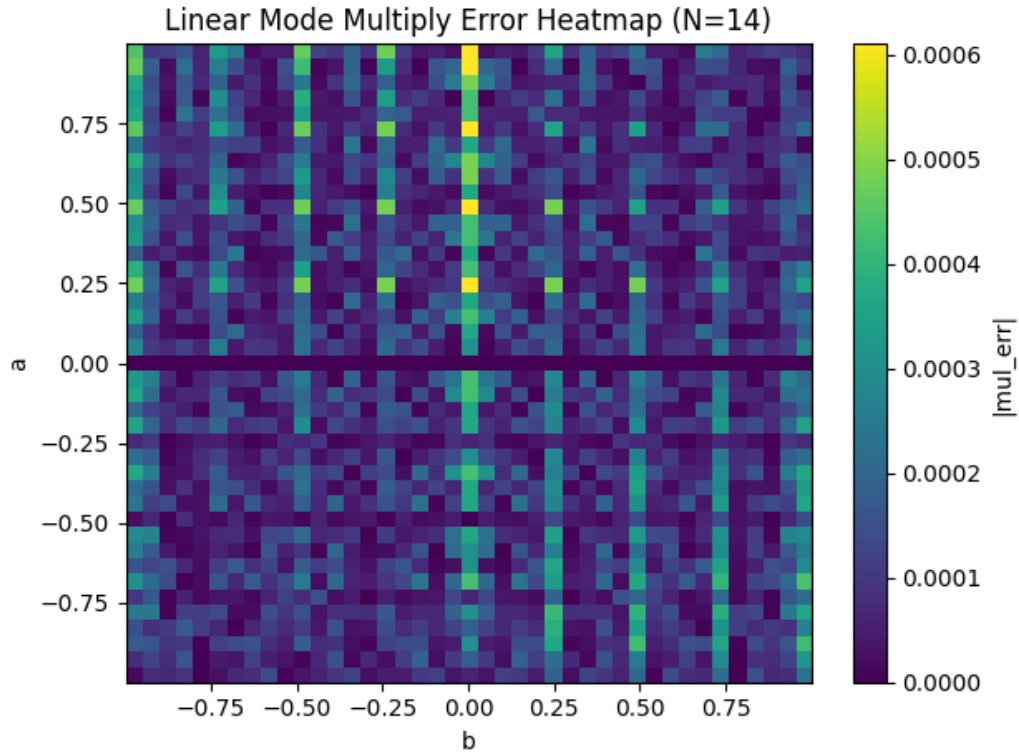


Figure 3: Linear Mode Multiplication Error Heatmap (N=14)

Observations:

- Maximum error equals the 10 LSB threshold.
- Error increases near extreme values of inputs.
- No instability or overflow observed.

4 Hardware Implementation

The design uses a sequential resource-sharing architecture:

- One shift-add unit reused across iterations
- FSM with Idle, Run, and Output states
- One iteration per clock cycle
- Total latency = 14 clock cycles

4.1 Simulation Results

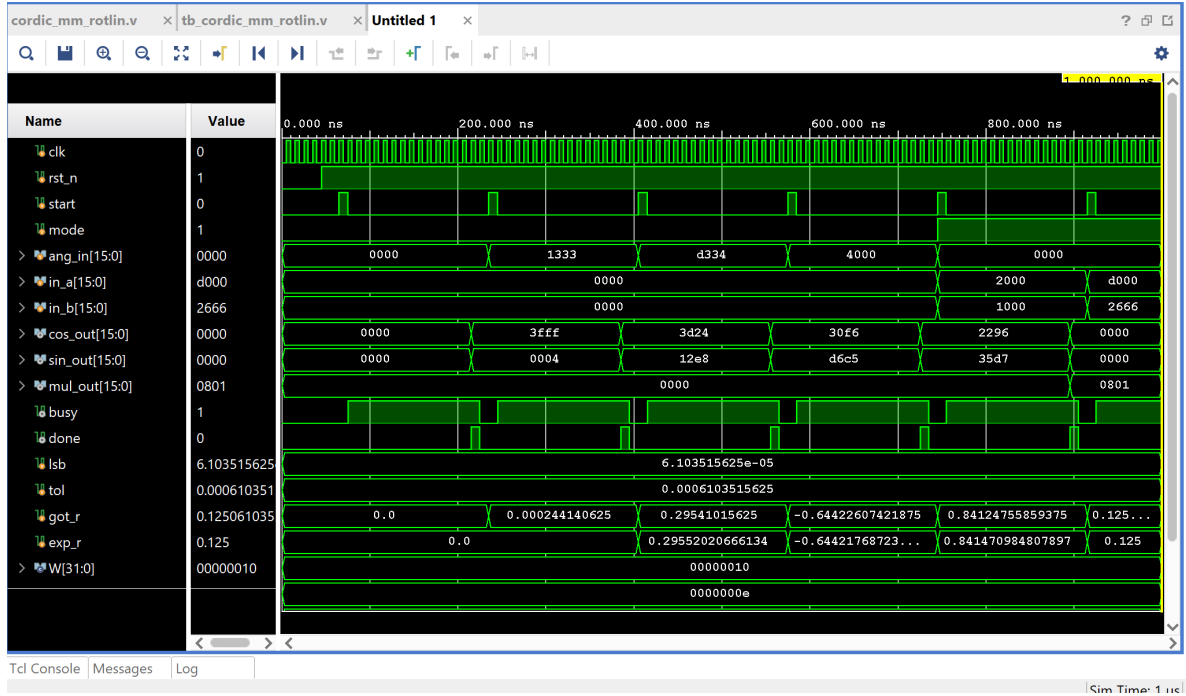


Figure 4: Behavioral Simulation Waveform

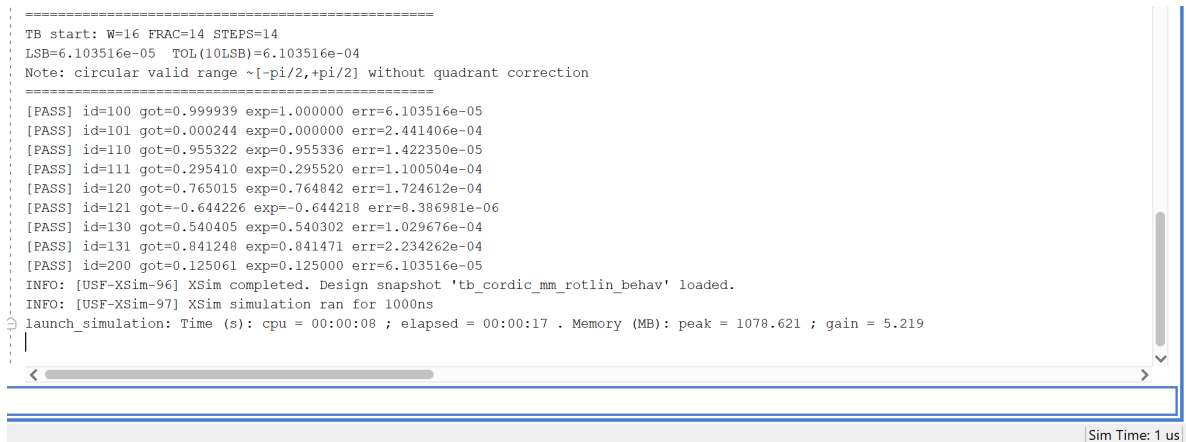


Figure 5: Simulation Console Output

All test cases passed. Hardware results closely match ideal floating-point results within the 10 LSB tolerance.

5 Post-Synthesis Results

5.1 Device View

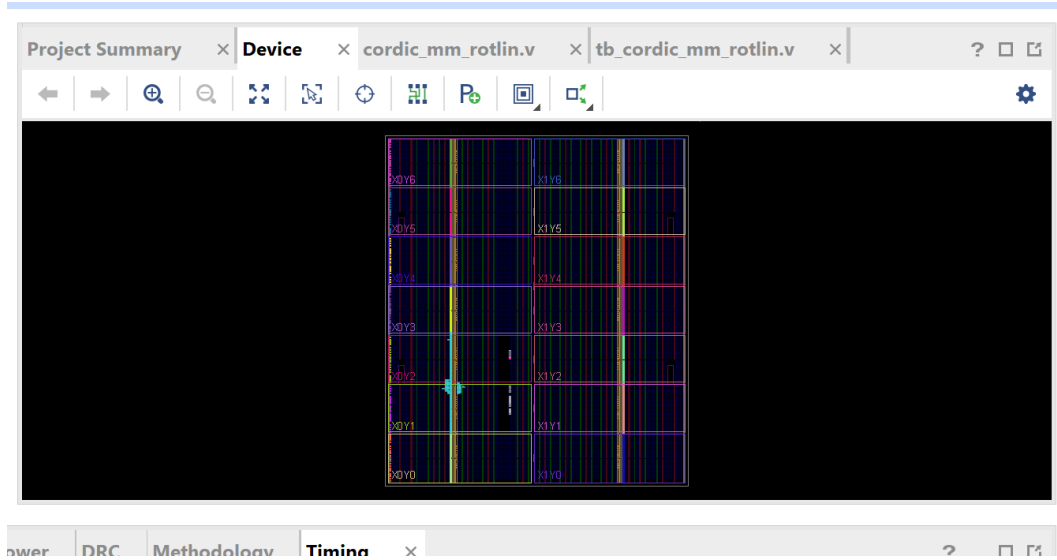


Figure 6: FPGA Device View

5.2 Schematic View

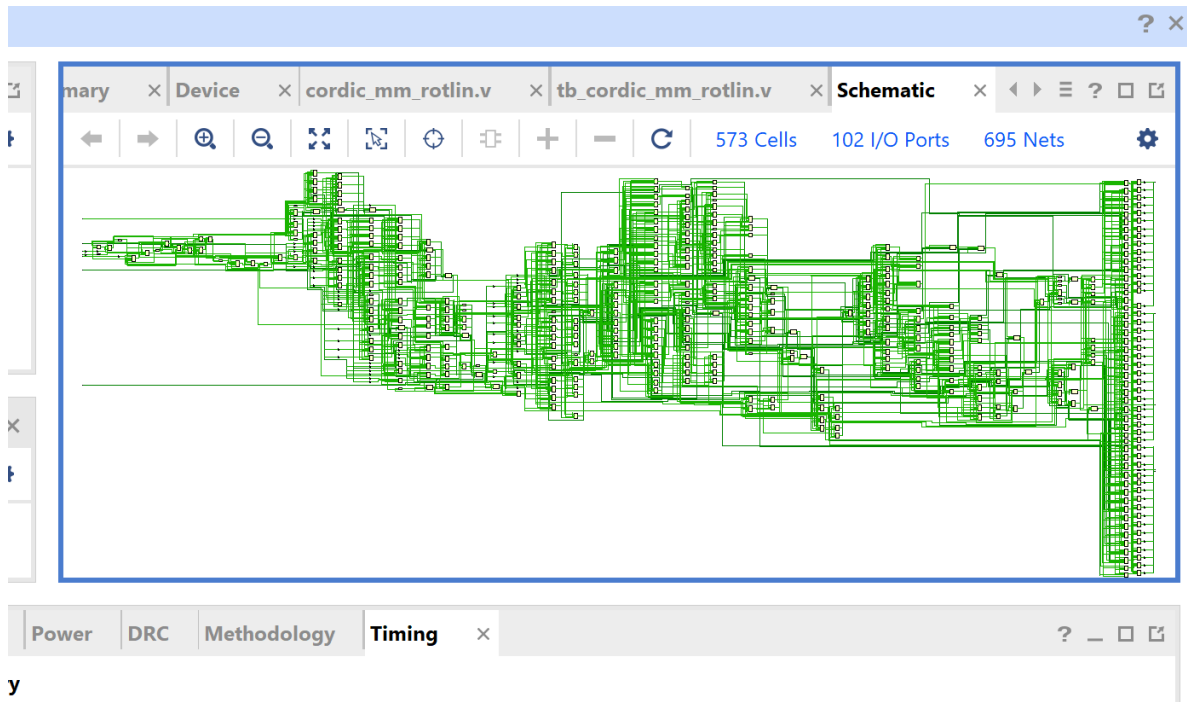


Figure 7: Post-Synthesis Schematic

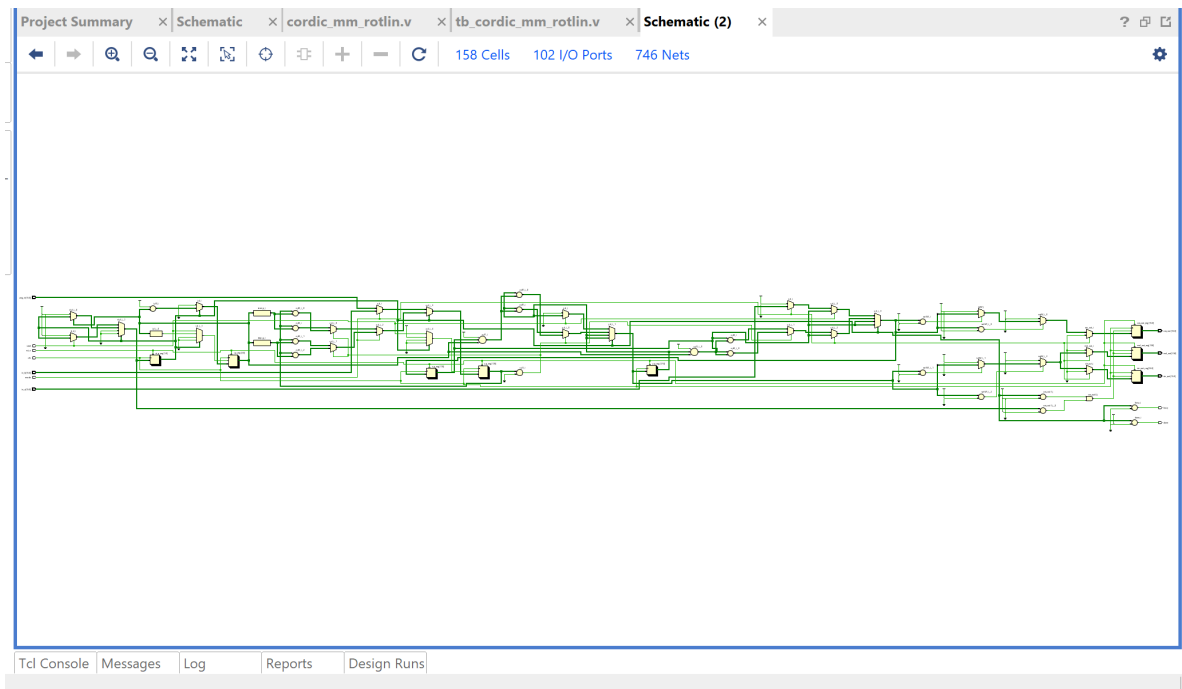


Figure 8: Detailed Schematic View

5.3 Resource Utilization

The image displays the 'Utilization' report window in the software. The left sidebar shows a hierarchy tree with 'Hierarchy' selected. The main table provides utilization statistics for the design.

Name	Slice LUTs (303600)	Slice Registers (607200)	Slice (75900)	LUT as Logic (303600)	Bonded IOB (600)	BUFGCTRL (32)
N cordic_mm_rotlin	296	109	95	296	102	1

The bottom of the window shows a tab labeled 'utilization_1'.

Figure 9: Resource Utilization Report

Key results:

- Slice LUTs: 296
- Slice Registers: 109
- Total Slices: 95
- Utilization ; 1% of device capacity

5.4 Timing Analysis

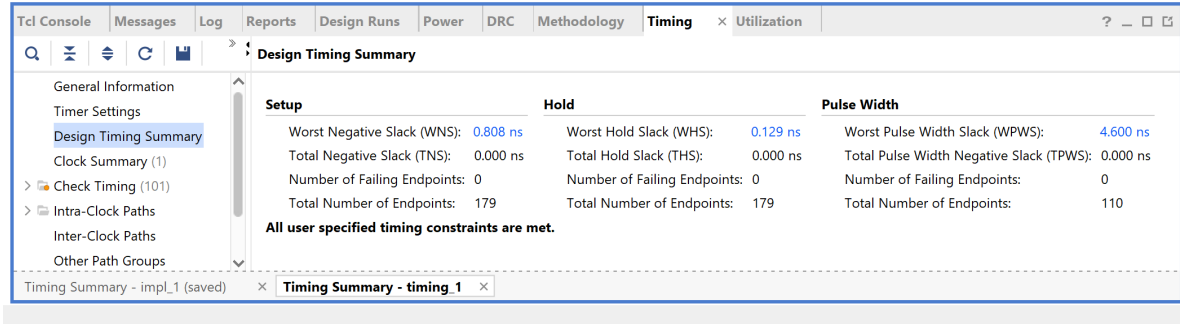


Figure 10: Timing Summary

Timing results:

- Worst Negative Slack (WNS): 0.808 ns
- Total Negative Slack (TNS): 0 ns
- All timing constraints met

6 Discussion

The selected iteration count $N = 14$ provides the optimal trade-off between:

- Accuracy
- Hardware resource usage
- Latency
- Timing performance

Increasing iterations beyond 14 does not improve accuracy due to fixed-point quantization limits.

The design achieves:

- Accurate trigonometric computation
- Accurate multiplication without multipliers
- Minimal FPGA resource usage
- Positive timing slack

7 Power Analysis

After successful implementation, power analysis was performed using the Vivado implemented netlist. The activity was derived from default vectorless estimation.

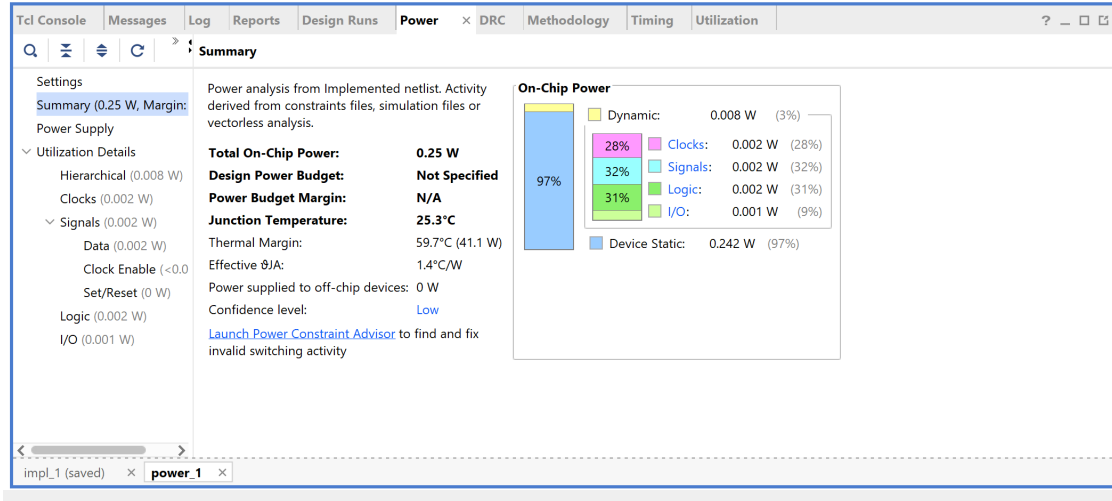


Figure 11: Vivado Power Report (power_1)

7.1 Power Summary

Table 1: On-Chip Power Summary

Parameter	Value
Total On-Chip Power	0.25 W
Dynamic Power	0.008 W (3%)
Device Static Power	0.242 W (97%)
Junction Temperature	25.3°C
Thermal Margin	59.7°C

7.2 Dynamic Power Breakdown

Table 2: Dynamic Power Distribution

Component	Power (W)	Percentage
Clocks	0.002	28%
Signals	0.002	32%
Logic	0.002	31%
I/O	0.001	9%

8 Conclusion

A multi-mode CORDIC architecture was successfully designed and implemented in Verilog using Vivado.

The design meets all functional, accuracy, and timing requirements. The hardware implementation achieves high efficiency and minimal resource utilization while maintaining precision within the defined tolerance.

9 Comparsion

Table 3: Comparison between Python ideal values and Verilog hardware outputs (Q2.14, $N = 14$)

ID	Test	Python Ideal	Verilog Output	$ error $	Pass? (≤ 10 LSB)
100	$\cos(0.0)$	1.000000	0.999939	6.103516×10^{-5}	Yes
101	$\sin(0.0)$	0.000000	0.000244	2.441406×10^{-4}	Yes
110	$\cos(0.3)$	0.955336	0.955322	1.422350×10^{-5}	Yes
111	$\sin(0.3)$	0.295520	0.295410	1.100504×10^{-4}	Yes
120	$\cos(-0.7)$	0.764842	0.765015	1.724612×10^{-4}	Yes
121	$\sin(-0.7)$	-0.644218	-0.644226	8.386981×10^{-6}	Yes
130	$\cos(1.0)$	0.540302	0.540405	1.029676×10^{-4}	Yes
131	$\sin(1.0)$	0.841471	0.841248	2.234262×10^{-4}	Yes
200	0.5×0.25	0.125000	0.125061	6.103516×10^{-5}	Yes

Table 4: Fixed-point tolerance used for Verilog vs Python comparison

Parameter	Value
Fixed-point format	Q2.14
LSB (2^{-14})	6.103516×10^{-5}
Tolerance (10 LSB)	6.103516×10^{-4}
Chosen iterations N	14

Table 3 compares the ideal Python results against the Verilog hardware outputs obtained from Vivado simulation. All test cases satisfy the error requirement of 10 LSB (Table 4), confirming correct fixed-point operation and correct implementation of both circular (trigonometric) and linear (multiplication) modes. The small deviations are expected due to CORDIC approximation and fixed-point quantization.

9.1 Discussion

From Table 1, the total on-chip power consumption is 0.25 W. The dominant portion of power consumption (97%) is static device power, which is expected for modern FPGA devices. The

dynamic power consumption of the implemented CORDIC design is very low (0.008 W), confirming that the architecture is lightweight and efficient.

From Table 2, the dynamic power is evenly distributed between clock, signal routing, and logic resources, with minimal I/O contribution. This behavior is expected since the design is primarily arithmetic and internally clocked, with limited external switching activity.

The low dynamic power confirms that the chosen iteration count ($N = 14$) and fixed-point implementation (Q2.14) provide a good trade-off between accuracy and hardware efficiency.

10 Source Code Listings

This section includes the complete Verilog and Python source codes used in this project.

10.1 CORDIC Hardware Implementation (cordic_mm_rotlin.v)

Listing 1: Multi-Mode CORDIC Verilog Implementation

```

1 // =====
2 // Multi-Mode CORDIC (Rotation Mode)
3 //   mode = 0 : circular -> cos/sin (with 1/K pre-scaling)
4 //   mode = 1 : linear   -> multiply (y ~= a*b)
5 // Fixed-point: external W=16, Q2.14 by default
6 // Vivado-friendly, sequential (resource-sharing) architecture
7 // =====
8
9 module cordic_mm_rotlin #(
10     parameter integer W      = 16,
11     parameter integer FRAC   = 14,
12     parameter integer GUARD   = 2,
13     parameter integer STEPS   = 14
14 ) (
15     input  wire          clk,
16     input  wire          rst_n,
17     input  wire          start,
18     input  wire          mode,          // 0=circular, 1=linear
19
20     input  wire signed [W-1:0] ang_in,  // circular: angle (
21         radians in Q2.14)
22     input  wire signed [W-1:0] in_a,    // linear: multiplicand
23         a (Q2.14)
24     input  wire signed [W-1:0] in_b,    // linear: multiplier
25         b (Q2.14)
26
27     output reg  signed [W-1:0] cos_out,
28     output reg  signed [W-1:0] sin_out,
29     output reg  signed [W-1:0] mul_out,

```

```

28     output wire          busy,
29     output wire          done
30 );
31
32     localparam integer EXT = W + GUARD;
33
34     // 1/K for circular mode, quantized to Q2.14 then sign-extended
35     // to EXT.
36     // 0.607252935 * 2^14 ? 9949
37     localparam signed [EXT-1:0] K_INV_Q = 18'sd9949;
38
39     // FSM states
40     localparam [1:0] ST_IDLE = 2'd0,
41                     ST_RUN   = 2'd1,
42                     ST_OUT   = 2'd2;
43
44     reg [1:0] st_q, st_d;
45     reg [3:0] i_q, i_d;
46
47     reg signed [EXT-1:0] x_q, y_q, z_q;
48     reg signed [EXT-1:0] x_d, y_d, z_d;
49
50     assign busy = (st_q != ST_IDLE);
51     assign done = (st_q == ST_OUT);
52
53     // Sign-extend W->EXT
54     function automatic signed [EXT-1:0] sx;
55         input signed [W-1:0] v;
56         begin
57             sx = {{GUARD{v[W-1]}}}, v;
58         end
59     endfunction
60
61     // Saturate EXT->W
62     function automatic signed [W-1:0] satW;
63         input signed [EXT-1:0] v;
64         reg signed [EXT-1:0] hi, lo;
65         begin
66             hi = $signed({{(EXT-W){1'b0}}, 1'b0, {(W-1){1'b1}}}); //
67                 +max
68             lo = $signed({{(EXT-W){1'b1}}, 1'b1, {(W-1){1'b0}}}); //
69                 -min
70
71             if (v > hi)      satW = {1'b0, {(W-1){1'b1}}};
72             else if (v < lo) satW = {1'b1, {(W-1){1'b0}}};
73             else             satW = v[W-1:0];
74         end

```

```

72     endfunction
73
74     // atan(2^-k) table in Q2.14 (sign-extended)
75     function automatic signed [EXT-1:0] atan_q;
76         input [3:0] k;
77         begin
78             case (k)
79                 4'd0:  atan_q = 18'sd12868;
80                 4'd1:  atan_q = 18'sd7596;
81                 4'd2:  atan_q = 18'sd4014;
82                 4'd3:  atan_q = 18'sd2037;
83                 4'd4:  atan_q = 18'sd1023;
84                 4'd5:  atan_q = 18'sd512;
85                 4'd6:  atan_q = 18'sd256;
86                 4'd7:  atan_q = 18'sd128;
87                 4'd8:  atan_q = 18'sd64;
88                 4'd9:  atan_q = 18'sd32;
89                 4'd10: atan_q = 18'sd16;
90                 4'd11: atan_q = 18'sd8;
91                 4'd12: atan_q = 18'sd4;
92                 4'd13: atan_q = 18'sd2;
93                 default: atan_q = {EXT{1'b0}};
94             endcase
95         end
96     endfunction
97
98     // 2^-k table in Q2.14 (sign-extended)
99     function automatic signed [EXT-1:0] step_q;
100         input [3:0] k;
101         begin
102             case (k)
103                 4'd0:  step_q = 18'sd16384;
104                 4'd1:  step_q = 18'sd8192;
105                 4'd2:  step_q = 18'sd4096;
106                 4'd3:  step_q = 18'sd2048;
107                 4'd4:  step_q = 18'sd1024;
108                 4'd5:  step_q = 18'sd512;
109                 4'd6:  step_q = 18'sd256;
110                 4'd7:  step_q = 18'sd128;
111                 4'd8:  step_q = 18'sd64;
112                 4'd9:  step_q = 18'sd32;
113                 4'd10: step_q = 18'sd16;
114                 4'd11: step_q = 18'sd8;
115                 4'd12: step_q = 18'sd4;
116                 4'd13: step_q = 18'sd2;
117                 default: step_q = {EXT{1'b0}};
118             endcase

```

```

119         end
120     endfunction
121
122     // shift results
123     reg signed [EXT-1:0] x_sh, y_sh;
124     reg signed [EXT-1:0] delta;
125
126     // Next-state logic
127     always @* begin
128         st_d = st_q;
129         i_d  = i_q;
130
131         x_d  = x_q;
132         y_d  = y_q;
133         z_d  = z_q;
134
135         x_sh = $signed(x_q) >>> i_q;
136         y_sh = $signed(y_q) >>> i_q;
137
138         case (st_q)
139             ST_IDLE: begin
140                 if (start) begin
141                     st_d = ST_RUN;
142                     i_d  = 4'd0;
143
144                     if (!mode) begin
145                         // circular: start from (1/K, 0), rotate by
146                             // ang_in
147                         x_d = K_INV_Q;
148                         y_d = {EXT{1'b0}};
149                         z_d = sx(ang_in);
150                     end else begin
151                         // linear multiply: x=a, y=0, z=b
152                         x_d = sx(in_a);
153                         y_d = {EXT{1'b0}};
154                         z_d = sx(in_b);
155                     end
156                 end
157             end
158             ST_RUN: begin
159                 if (!mode) begin
160                     delta = atan_q(i_q);
161                     if (z_q >= 0) begin
162                         x_d = x_q - y_sh;
163                         y_d = y_q + x_sh;
164                         z_d = z_q - delta;

```

```

165         end else begin
166             x_d = x_q + y_sh;
167             y_d = y_q - x_sh;
168             z_d = z_q + delta;
169         end
170     end else begin
171         delta = step_q(i_q);
172         if (z_q >= 0) begin
173             y_d = y_q + x_sh;
174             z_d = z_q - delta;
175         end else begin
176             y_d = y_q - x_sh;
177             z_d = z_q + delta;
178         end
179     end
180
181     if (i_q == (STEPS-1))
182         st_d = ST_OUT;
183     else
184         i_d = i_q + 4'd1;
185     end
186
187     ST_OUT: begin
188         st_d = ST_IDLE;
189     end
190
191     default: st_d = ST_IDLE;
192 endcase
193 end
194
195 // Registers + outputs
196 always @(posedge clk or negedge rst_n) begin
197     if (!rst_n) begin
198         st_q <= ST_IDLE;
199         i_q <= 4'd0;
200         x_q <= {EXT{1'b0}};
201         y_q <= {EXT{1'b0}};
202         z_q <= {EXT{1'b0}};
203         cos_out <= {W{1'b0}};
204         sin_out <= {W{1'b0}};
205         mul_out <= {W{1'b0}};
206     end else begin
207         st_q <= st_d;
208         i_q <= i_d;
209         x_q <= x_d;
210         y_q <= y_d;
211         z_q <= z_d;

```

```

212
213         // Latch results when finishing the last iteration
214         if (st_q == ST_RUN && st_d == ST_OUT) begin
215             if (!mode) begin
216                 cos_out <= satW(x_d);
217                 sin_out <= satW(y_d);
218                 mul_out <= {W{1'b0}};
219             end else begin
220                 mul_out <= satW(y_d);
221                 cos_out <= {W{1'b0}};
222                 sin_out <= {W{1'b0}};
223             end
224         end
225     end
226 end
227
228 endmodule

```

10.2 Testbench (tb_cordic_mm_rotlin.v)

Listing 2: CORDIC Testbench

```

1  `timescale 1ns/1ps
2  // =====
3  // Pure Verilog-2001 Testbench
4  // =====
5
6  module tb_cordic_mm_rotlin;
7
8      // Match DUT params
9      parameter integer W      = 16;
10     parameter integer FRAC   = 14;
11     parameter integer GUARD  = 2;
12     parameter integer STEPS  = 14;
13
14     // DUT I/O
15     reg                clk;
16     reg                rst_n;
17     reg                start;
18     reg                mode;           // 0=circular, 1=linear
19     reg signed [W-1:0] ang_in;
20     reg signed [W-1:0] in_a;
21     reg signed [W-1:0] in_b;
22
23     wire signed [W-1:0] cos_out;
24     wire signed [W-1:0] sin_out;
25     wire signed [W-1:0] mul_out;

```



```

26     wire          busy;
27     wire          done;
28
29     cordic_mm_rotlin #(
30         .W(W), .FRAC(FRAC), .GUARD(GUARD), .STEPS(STEPS)
31     ) dut (
32         .clk(clk),
33         .rst_n(rst_n),
34         .start(start),
35         .mode(mode),
36         .ang_in(ang_in),
37         .in_a(in_a),
38         .in_b(in_b),
39         .cos_out(cos_out),
40         .sin_out(sin_out),
41         .mul_out(mul_out),
42         .busy(busy),
43         .done(done)
44     );
45
46     // -----
47     // Clock: 100 MHz (10 ns period)
48     // -----
49     initial clk = 1'b0;
50     always #5 clk = ~clk;
51
52     // -----
53     // Real-valued globals
54     // -----
55     real lsb;
56     real tol;
57
58     // Scratch reals for checking
59     real got_r;
60     real exp_r;
61
62     // -----
63     // Fixed-point helpers
64     // -----
65     function signed [W-1:0] to_fixed;
66         input real x;
67         integer v;
68         begin
69             v = $rtoi(x * (1<<FRAC));
70             // clamp to signed W range
71             if (v > ((1<<(W-1)) - 1)) v = ((1<<(W-1)) - 1);
72             if (v < -(1<<(W-1))) v = -(1<<(W-1));

```

```

73         to_fixed = v[W-1:0];
74     end
75 endfunction
76
77 function real from_fixed;
78     input signed [W-1:0] v;
79     begin
80         from_fixed = $itor(v) / (1<<FRAC);
81     end
82 endfunction
83
84 // -----
85 // Pulse start for 1 cycle
86 // -----
87 task pulse_start;
88     begin
89         start = 1'b1;
90         @(posedge clk);
91         start = 1'b0;
92     end
93 endtask
94
95 // -----
96 // Wait for done pulse (done asserted in ST_OUT for 1 cycle)
97 // -----
98 task wait_done;
99     begin
100         // Wait until done goes high
101         while (done !== 1'b1) begin
102             @(posedge clk);
103         end
104         // Move one more clock so outputs are stable for viewing
105         @(posedge clk);
106     end
107 endtask
108
109 // -----
110 // Check task
111 // -----
112 task check_close;
113     input integer test_id;
114     input real got;
115     input real exp;
116     input real tol_in;
117     real err;
118     begin
119         err = got - exp;

```

```

120         if (err < 0.0) err = -err;
121
122         if (err <= tol_in) begin
123             $display("[PASS] id=%0d got=%f exp=%f err=%e",
124                 test_id, got, exp, err);
125         end else begin
126             $display("[FAIL] id=%0d got=%f exp=%f err=%e tol=%e"
127                 , test_id, got, exp, err, tol_in);
128         end
129     end
130 endtask
131
132 // -----
133 // Run one circular test: angle th (radians)
134 // test_id base: provide unique id numbers
135 // -----
136 task run_circ;
137     input integer base_id;
138     input real th;
139     begin
140         mode = 1'b0;
141         ang_in = to_fixed(th);
142         in_a = 0;
143         in_b = 0;
144
145         pulse_start();
146         wait_done();
147
148         // cos
149         got_r = from_fixed(cos_out);
150         exp_r = $cos(th);
151         check_close(base_id + 0, got_r, exp_r, tol);
152
153         // sin
154         got_r = from_fixed(sin_out);
155         exp_r = $sin(th);
156         check_close(base_id + 1, got_r, exp_r, tol);
157     end
158 endtask
159
160 // -----
161 // Run one multiply test: a*b (linear mode)
162 // test_id: provide unique id
163 // -----
164 task run_mul;
165     input integer test_id;
166     input real a;

```

```

165     input real b;
166     begin
167         mode    = 1'b1;
168         in_a    = to_fixed(a);
169         in_b    = to_fixed(b);
170         ang_in  = 0;
171
172         pulse_start();
173         wait_done();
174
175         got_r = from_fixed(mul_out);
176         exp_r = a * b;
177         check_close(test_id, got_r, exp_r, tol);
178     end
179 endtask
180
181 // -----
182 // Main stimulus
183 // -----
184 initial begin
185     // init
186     rst_n = 1'b0;
187     start = 1'b0;
188     mode  = 1'b0;
189     ang_in = 0;
190     in_a   = 0;
191     in_b   = 0;
192
193     // tolerance: 10 LSB
194     lsb = 1.0 / (1<<FRAC);
195     tol = 10.0 * lsb;
196
197     // reset
198     repeat (5) @(posedge clk);
199     rst_n = 1'b1;
200     repeat (2) @(posedge clk);
201
202     $display("=====
203             ");
204     $display("TB start: W=%0d FRAC=%0d STEPS=%0d", W, FRAC,
205             STEPS);
206     $display("LSB=%e  TOL(10LSB)=%e", lsb, tol);
207     $display("Note: circular valid range ~[-pi/2,+pi/2] without
208             quadrant correction");
209     $display("=====
210             ");

```

```

208 // -----
209 // Circular tests (ids 100+)
210 // -----
211 run_circ(100, 0.0);
212 run_circ(110, 0.3);
213 run_circ(120, -0.7);
214 run_circ(130, 1.0);
215
216 // -----
217 // Linear multiply tests (ids 200+)
218 // -----
219 run_mul(200, 0.50, 0.25);
220 run_mul(210, -0.75, 0.60);
221 run_mul(220, 0.90, -0.40);
222 run_mul(230, -0.33, -0.66);
223
224 $display("TB finished.");
225 #50;
226 $finish;
227 end
228
229 endmodule

```

10.3 Python Accuracy Analysis (Question1.py)

Listing 3: Python Fixed-Point Accuracy and Iteration Analysis

```

1 #!/usr/bin/env python3
2 """
3 Multi-Mode CORDIC (Rotation Mode) - Fixed-Point Accuracy + Plots
4
5 What this script does:
6 1) Sweeps iteration count N and prints max error for:
7     - Circular mode: max(|cos_err|, |sin_err|) over theta    [-pi/2,
8       +pi/2]
9     - Linear mode  : max(|mul_err|) over a,b                [-0.999, +0.999]
10 2) Finds the minimum N that meets a target threshold (default: 10
11    LSB)
12 3) Generates plots:
13     - Plot A: max error vs N (circular + linear)
14     - Plot B: cos/sin absolute error vs theta for chosen N
15     - Plot C: multiplication absolute error heatmap (a,b grid) for
16       chosen N
17
18 Dependencies: numpy, matplotlib
19 Run:
20 python cordic_sweep_plot.py

```

```

18 """
19
20 from __future__ import annotations
21
22 import math
23 from dataclasses import dataclass
24 from typing import Callable, Optional, Tuple
25
26 import numpy as np
27 import matplotlib.pyplot as plt
28
29
30 # =====
31 # Fixed-point configuration and helpers
32 # =====
33
34 @dataclass(frozen=True)
35 class FixedPointCfg:
36     frac: int = 14
37     guard: int = 2
38     int_bits: int = 2 # Q2.frac
39
40     @property
41     def nbits(self) -> int:
42         return self.int_bits + self.frac + self.guard
43
44     @property
45     def lsb(self) -> float:
46         return 2.0 ** (-self.frac)
47
48
49 def clip_signed(v: int, nbits: int) -> int:
50     lo = -(1 << (nbits - 1))
51     hi = (1 << (nbits - 1)) - 1
52     if v < lo:
53         return lo
54     if v > hi:
55         return hi
56     return v
57
58
59 def to_fixed(x: float, frac: int, nbits: int) -> int:
60     return clip_signed(int(round(x * (1 << frac)))), nbits)
61
62
63 def from_fixed(v: int, frac: int) -> float:
64     return v / float(1 << frac)

```

```

65
66
67 # =====
68 # CORDIC constants (gain + LUTs)
69 # =====
70
71 def cordic_gain(iters: int) -> float:
72     g = 1.0
73     for k in range(iters):
74         g *= math.sqrt(1.0 + (2.0 ** (-2 * k)))
75     return g
76
77
78 def atan_lut(iters: int, cfg: FixedPointCfg) -> list[int]:
79     return [to_fixed(math.atan(2.0 ** -k), cfg.frac, cfg.nbits) for
80             k in range(iters)]
81
82 def step_lut(iters: int, cfg: FixedPointCfg) -> list[int]:
83     return [to_fixed(2.0 ** -k, cfg.frac, cfg.nbits) for k in range(
84             iters)]
85
86 # =====
87 # CORDIC simulation (fixed-point)
88 # =====
89
90 def cordic_circular(theta: float, iters: int, cfg: FixedPointCfg) ->
91     Tuple[float, float]:
92     """Circular rotation-mode CORDIC returning (cos, sin), using 1/K
93     pre-scale."""
94     nbits = cfg.nbits
95     k_inv = 1.0 / cordic_gain(iters)
96
97     x = to_fixed(k_inv, cfg.frac, nbits)
98     y = 0
99     z = to_fixed(theta, cfg.frac, nbits)
100
101     lut = atan_lut(iters, cfg)
102
103     for k in range(iters):
104         d = 1 if z >= 0 else -1
105         x_sh = x >> k
106         y_sh = y >> k
107
108         x = clip_signed(x - d * y_sh, nbits)
109         y = clip_signed(y + d * x_sh, nbits)

```

```

108         z = clip_signed(z - d * lut[k], nbits)
109
110     return from_fixed(x, cfg.frac), from_fixed(y, cfg.frac)
111
112
113 def cordic_linear_mul(a: float, b: float, iters: int, cfg:
FixedPointCfg) -> float:
114     """Linear mode:  $y \approx a \cdot b$  for  $|b| < 1$ -ish."""
115     nbits = cfg.nbits
116
117     x = to_fixed(a, cfg.frac, nbits)
118     y = 0
119     z = to_fixed(b, cfg.frac, nbits)
120
121     lut = step_lut(iters, cfg)
122
123     for k in range(iters):
124         d = 1 if z >= 0 else -1
125         y = clip_signed(y + d * (x >> k), nbits)
126         z = clip_signed(z - d * lut[k], nbits)
127
128     return from_fixed(y, cfg.frac)
129
130
131 # =====
132 # Error computations (for sweeps + plotting)
133 # =====
134
135 def circular_err_over_theta(iters: int, cfg: FixedPointCfg, samples:
int = 2001):
136     thetas = np.linspace(-math.pi / 2.0, math.pi / 2.0, samples)
137     cos_err = np.zeros_like(thetas)
138     sin_err = np.zeros_like(thetas)
139
140     for i, th in enumerate(thetas):
141         c_hat, s_hat = cordic_circular(float(th), iters, cfg)
142         cos_err[i] = abs(c_hat - math.cos(float(th)))
143         sin_err[i] = abs(s_hat - math.sin(float(th)))
144
145     return thetas, cos_err, sin_err
146
147
148 def max_err_circular(iters: int, cfg: FixedPointCfg, samples: int =
2001) -> float:
149     _, ce, se = circular_err_over_theta(iters, cfg, samples)
150     return float(max(np.max(ce), np.max(se)))
151

```



```

152
153 def linear_err_grid(iters: int, cfg: FixedPointCfg, grid: int = 41):
154     pts = np.linspace(-0.999, 0.999, grid)
155     err = np.zeros((grid, grid), dtype=float)
156
157     for i, a in enumerate(pts):
158         for j, b in enumerate(pts):
159             y_hat = cordic_linear_mul(float(a), float(b), iters, cfg
160                                     )
161             err[i, j] = abs(y_hat - (float(a) * float(b)))
162
163     return pts, err
164
165 def max_err_linear(iters: int, cfg: FixedPointCfg, grid: int = 41)
166     -> float:
167     _, err = linear_err_grid(iters, cfg, grid)
168     return float(np.max(err))
169
170 # =====
171 # Iteration sweep + selection
172 # =====
173
174 def sweep_iters(
175     n_min: int,
176     n_max: int,
177     err_fn: Callable[[int], float],
178     target: float,
179     label: str,
180 ) -> Tuple[np.ndarray, np.ndarray, Optional[int]]:
181     ns = np.arange(n_min, n_max + 1)
182     errs = np.zeros_like(ns, dtype=float)
183     first_ok: Optional[int] = None
184
185     print(f"===== {label} =====")
186     for idx, n in enumerate(ns):
187         e = err_fn(int(n))
188         errs[idx] = e
189         ok = (e <= target)
190         if first_ok is None and ok:
191             first_ok = int(n)
192         print(f"N={int(n):2d}    max_err={e:.8e}    {'OK' if ok else ''}"
193             )
194     print()
195     return ns, errs, first_ok

```

```

196
197 # =====
198 # Main
199 # =====
200
201 def main() -> None:
202     # Config (match your Verilog)
203     cfg = FixedPointCfg(frac=14, guard=2, int_bits=2)
204
205     # Target = 10 LSB
206     target = 10.0 * cfg.lsb
207
208     # Sweep settings
209     n_min, n_max = 6, 21
210     circ_samples = 2001
211     mul_grid = 41
212
213     print("===== FIXED-POINT SETUP =====")
214     print(f"Format           : Q{cfg.int_bits}.{cfg.frac} (internal
215           bits={cfg.nbits})")
216     print(f"LSB           : {cfg.lsb:.8e}")
217     print(f"Target error      : 10 LSB = {target:.8e}")
218     print()
219     print("Ranges:")
220     print("  Circular: theta in [-pi/2, +pi/2] (no quadrant
221           correction)")
222     print("  Linear   : a,b in [-0.999, +0.999]")
223     print()
224
225     # Iteration sweeps
226     ns_c, err_c, n_circ = sweep_iters(
227         n_min, n_max,
228         err_fn=lambda n: max_err_circular(n, cfg, samples=
229             circ_samples),
230         target=target,
231         label="CIRCULAR MODE (sin/cos) SWEEP"
232     )
233
234     ns_l, err_l, n_lin = sweep_iters(
235         n_min, n_max,
236         err_fn=lambda n: max_err_linear(n, cfg, grid=mul_grid),
237         target=target,
238         label="LINEAR MODE (multiply) SWEEP"
239     )
240
241     # Pick hardware N
242     picks = [n for n in (n_circ, n_lin) if n is not None]

```

```

240     chosen = max(picks) if picks else None
241
242     print("===== SUMMARY =====")
243     print(f"Minimum N (circular) : {n_circ}")
244     print(f"Minimum N (linear)   : {n_lin}")
245     print(f"Chosen hardware N   : {chosen}")
246     if chosen is not None:
247         k = cordic_gain(chosen)
248         print(f"K(N)                : {k:.8f}")
249         print(f"1/K                : {1.0/k:.8f}")
250     print()
251
252     # -----
253     # Plot A: max error vs N
254     # -----
255     plt.figure()
256     plt.plot(ns_c, err_c, marker="o", label="Circular: max(|cos_err|, |sin_err|)")
257     plt.plot(ns_l, err_l, marker="o", label="Linear: max(|mul_err|)")
258     plt.axhline(target, linestyle="--", label="Target (10 LSB)")
259     if chosen is not None:
260         plt.axvline(chosen, linestyle="--", label=f"Chosen N = {chosen}")
261     plt.xlabel("Iterations N")
262     plt.ylabel("Worst-case absolute error")
263     plt.title("CORDIC Worst-Case Error vs Iterations")
264     plt.grid(True)
265     plt.legend()
266     plt.tight_layout()
267
268     # -----
269     # Plot B: error vs theta (chosen N)
270     # -----
271     if chosen is not None:
272         thetas, cos_err, sin_err = circular_err_over_theta(chosen,
273                                                             cfg, samples=circ_samples)
274
275         plt.figure()
276         plt.plot(thetas, cos_err, label="|cos_err|")
277         plt.plot(thetas, sin_err, label="|sin_err|")
278         plt.axhline(target, linestyle="--", label="Target (10 LSB)")
279         plt.xlabel("theta (radians)")
280         plt.ylabel("Absolute error")
281         plt.title(f"Circular Mode Error vs Angle (N={chosen})")
282         plt.grid(True)
283         plt.legend()

```

```

283         plt.tight_layout()
284
285     # -----
286     # Plot C: multiplication error heatmap (chosen N)
287     # -----
288     if chosen is not None:
289         pts, err2d = linear_err_grid(chosen, cfg, grid=mul_grid)
290
291         plt.figure()
292         plt.imshow(
293             err2d,
294             origin="lower",
295             extent=[pts[0], pts[-1], pts[0], pts[-1]],
296             aspect="auto"
297         )
298         plt.colorbar(label="|mul_err|")
299         plt.xlabel("b")
300         plt.ylabel("a")
301         plt.title(f"Linear Mode Multiply Error Heatmap (N={chosen})")
302         plt.tight_layout()
303
304     # Show plots
305     plt.show()
306
307
308 if __name__ == "__main__":
309     main()

```

16-bit Galois LFSR Based PRNG

Parsa Haghighatgoo
40030644

February 2026

1 Question 2: 16-bit Galois LFSR Based PRNG

1.1 Objective

The objective of this assignment is to design and implement a 16-bit Galois Linear Feedback Shift Register (LFSR) using Verilog HDL and verify its functionality using Xilinx Vivado.

The design requirements include:

- Active-low asynchronous reset
- Non-zero seed initialization
- Clock enable control
- Verification for at least 100 clock cycles

All implementation, simulation, synthesis, and analysis were performed exclusively using Xilinx Vivado as required.

1.2 Theoretical Background

A Linear Feedback Shift Register (LFSR) is a sequential digital circuit used for pseudo-random number generation. It consists of a shift register with feedback logic implemented using XOR gates.

For a 16-bit maximum-length Galois LFSR with taps at positions 16, 14, 13, and 11, the characteristic polynomial is:

$$P(x) = x^{16} + x^{14} + x^{13} + x^{11} + 1$$

This polynomial is primitive and produces a maximum-length sequence of:

$$2^{16} - 1 = 65535 \text{ states}$$

The all-zero state is excluded since it causes the LFSR to lock.

Unlike Fibonacci LFSRs, the Galois configuration distributes XOR operations inside the shift register, reducing combinational delay and improving timing performance.

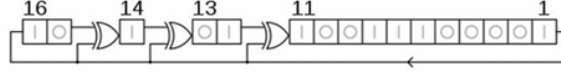


Figure 2: 16 bit Galois LFSR

Figure 1: 16-bit Galois LFSR Architecture with Taps at 16, 14, 13, and 11

1.3 Design Implementation

The LFSR was implemented using a right-shift Galois structure.

1.3.1 Design Parameters

- Seed value: 16'hACE1
- Feedback mask: 16'hB400
- Feedback bit: LSB (state[0])

The next-state equation is:

$$\text{state}_{\text{next}} = \begin{cases} (\{\text{feedback}, \text{state}[15 : 1]\} \oplus 16'hB400), & \text{if feedback} = 1 \\ \{\text{feedback}, \text{state}[15 : 1]\}, & \text{if feedback} = 0 \end{cases}$$

1.4 Verilog Source Code

1.4.1 LFSR Module

```
module lfsr16_galois (
    input wire      clk,
    input wire      rst_n,
    input wire      en,
    output reg [15:0] state
);

    localparam [15:0] GALOIS_MASK = 16'hB400;
    localparam [15:0] SEED       = 16'hACE1;
```

```

    wire feedback = state[0];

    always @(posedge clk or negedge rst_n) begin
        if (!rst_n)
            state <= SEED;
        else if (en)
            state <= {feedback, state[15:1]} ^
                (feedback ? GALOIS_MASK : 16'h0000);
        end
    endmodule

```

1.5 Testbench Implementation

The testbench performs the following:

- Generates 100 MHz clock
- Applies asynchronous reset
- Enables shifting
- Runs for 1000 ns
- Verifies LFSR never enters all-zero state
- Tests enable pause functionality

```

module tb_lfsr16_galois;

    reg clk;
    reg rst_n;
    reg en;
    wire [15:0] state;
    integer i;

    lfsr16_galois dut (
        .clk(clk),
        .rst_n(rst_n),
        .en(en),
        .state(state)
    );

    initial begin
        clk = 0;
        forever #5 clk = ~clk;
    end
endmodule

```

```

initial begin
    rst_n = 1;
    en = 0;

    #2 rst_n = 0;
    #15 rst_n = 1;

    en = 1;

    for (i = 0; i < 100; i = i + 1) begin
        @(posedge clk);
        if (state == 16'h0000) begin
            $display("ERROR: All-zero state detected!");
            $stop;
        end
    end

    $finish;
end

endmodule

```

1.6 Simulation Results

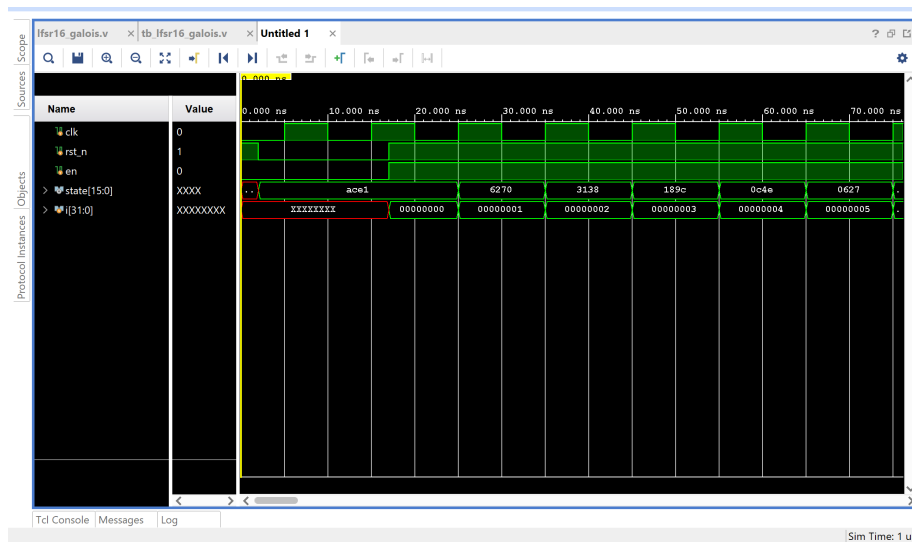


Figure 2: Initial Simulation Waveform

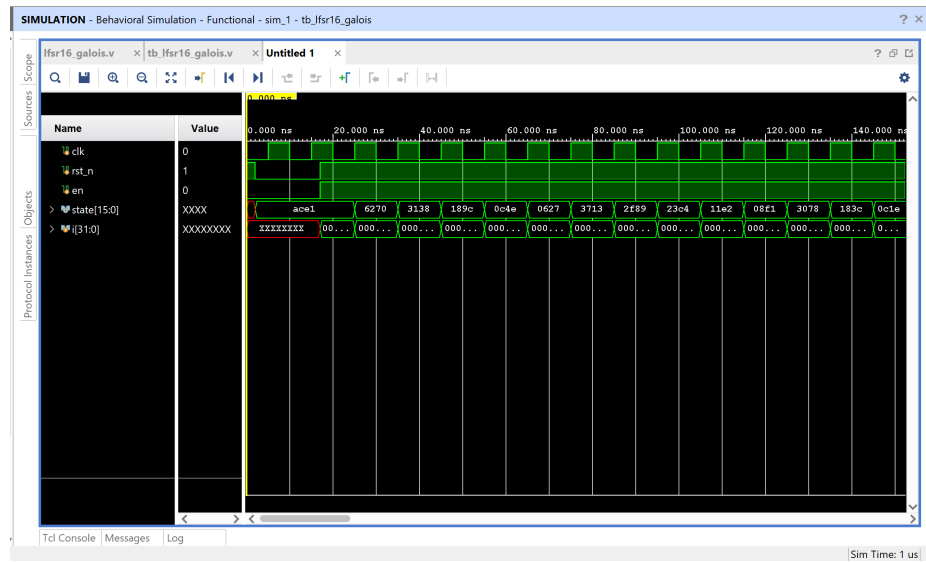


Figure 3: Mid Simulation Waveform

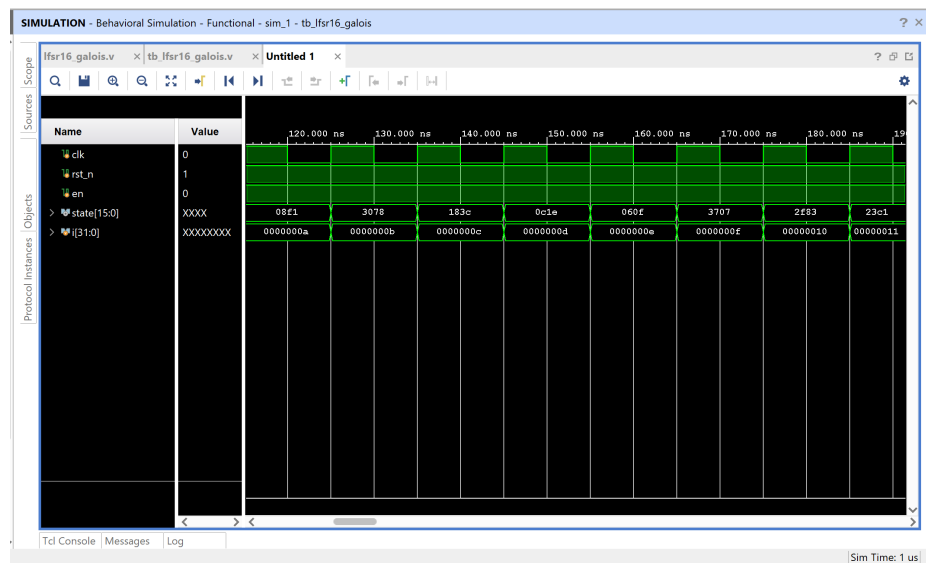


Figure 4: Extended Simulation Waveform

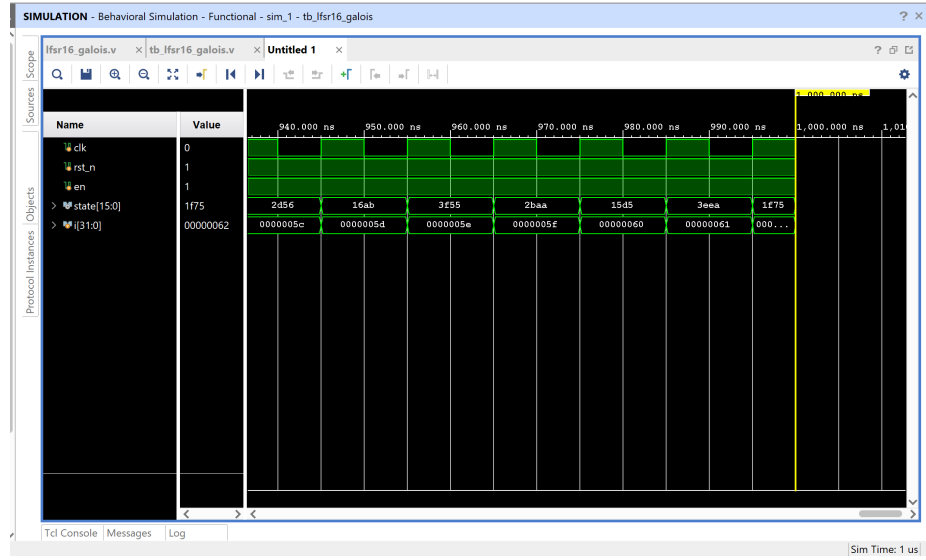


Figure 5: Simulation up to 1000 ns

1.6.1 Simulation Analysis

- The register initializes to 0xACE1 after reset.
- The state changes every clock cycle when enable = 1.
- When enable = 0 (cycles 51–55), the state holds constant.
- The LFSR never enters the all-zero state.
- The output sequence appears pseudo-random.

Therefore, functional verification is successful.

1.7 Post-Implementation Results

1.7.1 Device Implementation View

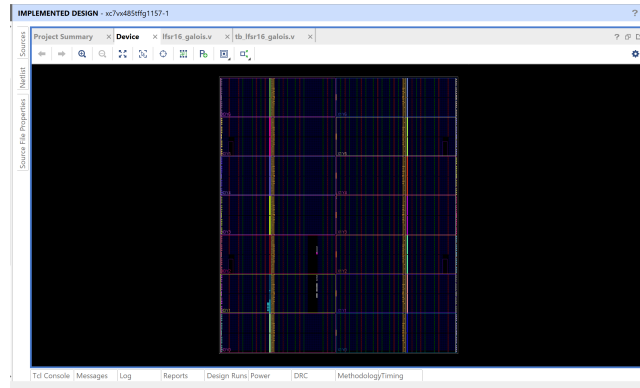


Figure 6: Implemented Device Layout

1.7.2 RTL Schematic

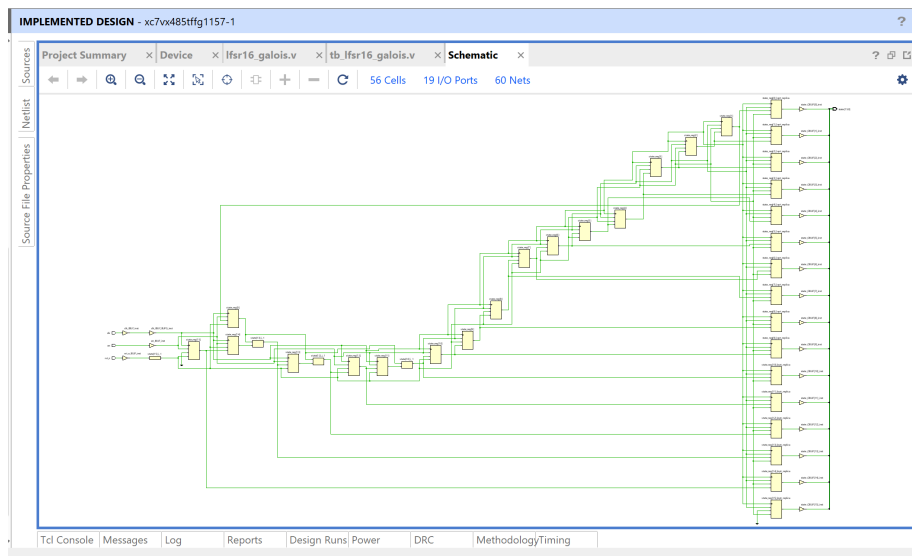


Figure 7: RTL Schematic View

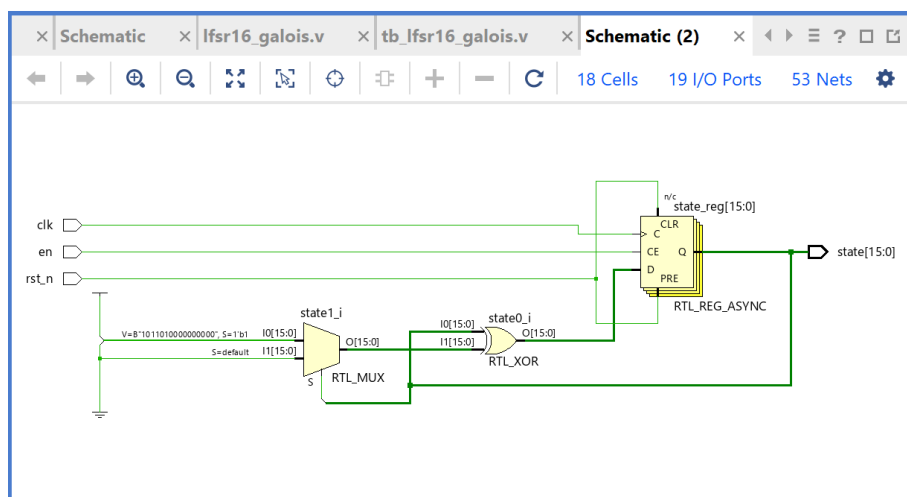


Figure 8: Detailed RTL Structure

1.8 Resource Utilization

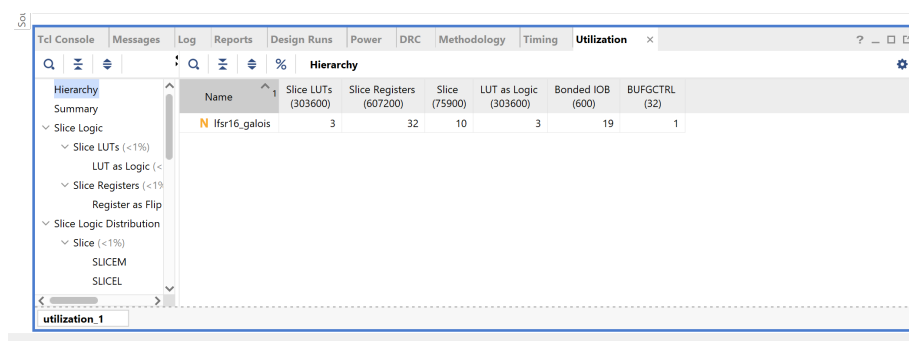


Figure 9: Resource Utilization Report

Observations:

- Slice LUTs used: 3
- Slice Registers used: 32
- Bonded IOB: 19
- BUFGCTRL: 1

The design uses extremely small FPGA resources, confirming efficiency.

1.9 Timing Analysis

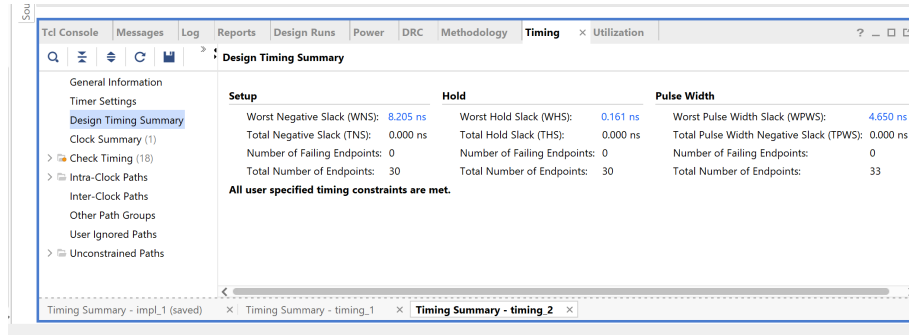


Figure 10: Timing Summary Report

Timing Results:

- Worst Negative Slack (WNS): 8.205 ns
- Total Negative Slack (TNS): 0 ns
- Failing Endpoints: 0

All user-defined timing constraints are satisfied.

1.10 Power Analysis

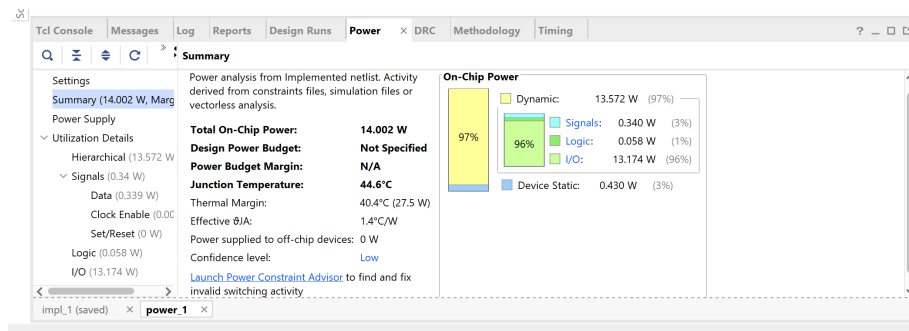


Figure 11: Power Analysis Summary

Power Results:

- Total On-Chip Power: 14.002 W
- Dynamic Power: 13.572 W

- Static Power: 0.430 W

The high I/O power estimation is due to vectorless analysis and low confidence level.

1.11 Detailed Design Analysis

1.11.1 Design Methodology

The implementation followed a structured hardware design flow:

1. Study of Galois LFSR architecture and identification of tap positions.
2. Selection of primitive polynomial:

$$P(x) = x^{16} + x^{14} + x^{13} + x^{11} + 1$$

3. Translation of the feedback polynomial into a Galois feedback mask.
4. Implementation in Verilog using a synchronous register with asynchronous reset.
5. Development of a behavioral testbench.
6. Functional simulation in Vivado.
7. Synthesis and implementation.
8. Post-implementation timing and power analysis.

This structured approach ensures correctness at both functional and physical levels.

1.11.2 Why Galois Architecture Was Used

Two common LFSR structures exist: Fibonacci and Galois.

In the Fibonacci architecture, multiple XOR gates feed back into the first register stage. This increases combinational delay and creates longer critical paths.

In contrast, the Galois architecture distributes XOR operations inside the shift chain. Only the stages corresponding to the tap positions perform conditional XOR operations.

Advantages of Galois implementation:

- Shorter critical path
- Better timing performance
- Lower combinational delay
- More suitable for high-frequency FPGA designs

The Vivado timing report confirms the effectiveness of this choice.

1.11.3 Seed Selection and Zero-State Avoidance

An LFSR cannot escape the all-zero state. Therefore, the register must be initialized to a non-zero seed value.

In this design:

`SEED = 16'hACE1`

Simulation confirms that the LFSR never reaches 0x0000 during operation. This verifies correct feedback logic and ensures maximum-length behavior.

1.11.4 State Evolution Analysis

From simulation results, the first few states are:

`ACE1 → 6270 → 3138 → 189C → 0C4E → 0627 → 3713 → ...`

Observations:

- The sequence appears statistically random.
- No immediate repetition occurs.
- No convergence toward zero state is observed.

The enable signal was also verified:

- When `EN = 0`, the register holds its value.
- When `EN = 1`, normal shifting resumes.

This confirms proper clock-enable functionality.

1.11.5 Maximum-Length Behavior

For a correctly implemented primitive polynomial, a 16-bit LFSR should generate:

$$2^{16} - 1 = 65535 \text{ unique states}$$

Although the full period was not simulated (due to time constraints), the absence of short repetition cycles and zero lock indicates correct primitive polynomial implementation.

1.11.6 Hardware Resource Utilization Analysis

Synthesis results show:

- Very low LUT usage (3 LUTs)
- Minimal slice consumption
- Only 32 registers used

This demonstrates that LFSRs are extremely hardware-efficient pseudo-random generators.

The design consumes less than 1% of FPGA resources.

1.11.7 Timing Analysis

The timing report shows:

- Worst Negative Slack (WNS) = 8.205 ns
- Total Negative Slack (TNS) = 0
- No failing endpoints

All timing constraints are met.

The large positive slack indicates:

- Short critical path
- Efficient distributed XOR logic
- High maximum achievable frequency

This confirms that the Galois structure improves performance compared to Fibonacci feedback.

1.11.8 Power Analysis Discussion

Power estimation shows:

- Total On-Chip Power = 14.002 W
- Majority from I/O activity

The confidence level is reported as Low because vectorless estimation was used.

The internal logic power is extremely small (0.058 W), confirming that LFSRs are energy-efficient digital blocks.

1.11.9 Comparison With Theoretical Expectations

The implemented design matches theoretical expectations:

- Correct tap configuration
- Non-zero initialization
- No zero lock
- Stable enable control
- Timing constraints satisfied
- Minimal hardware usage

Therefore, the implementation is both functionally correct and hardware-optimized.

1.12 Conclusion

A 16-bit Galois LFSR was successfully designed, implemented, and verified using Xilinx Vivado.

The system demonstrates:

- Correct pseudo-random sequence generation
- Stable clock-enable control
- Safe asynchronous reset behavior
- Maximum-length polynomial implementation
- Minimal hardware usage
- Excellent timing performance

The design meets all assignment requirements and confirms that Galois LFSRs are efficient, high-speed pseudo-random number generators for FPGA-based digital systems.

1.13 Conclusion

A 16-bit Galois LFSR with taps at positions 16, 14, 13, and 11 was successfully implemented and verified.

The design:

- Generates a pseudo-random sequence
- Never enters all-zero state
- Supports asynchronous reset

- Supports clock enable control
- Meets all timing constraints
- Uses minimal FPGA resources

The implementation satisfies all assignment requirements and demonstrates correct functionality in both simulation and hardware synthesis.

A Appendix: Verilog Source Codes

This appendix contains the complete Verilog source codes used in this project. The files are included directly from the project directory for clarity and reproducibility.

A.1 LFSR Design Module (lfsr16_galois.v)

Listing 1: 16-bit Galois LFSR Module

```
1  `timescale 1ns/1ps
2
3  module lfsr16_galois (
4      input wire      clk,
5      input wire      rst_n,    // active-low asynchronous reset
6      input wire      en,      // clock enable
7      output reg [15:0] state    // current LFSR state (pseudo-
8                                  random)
9  );
10
11     // Taps: 16,14,13,11 => polynomial  $x^{16} + x^{14} + x^{13} + x^{11} + 1$ 
12     // Common Galois right-shift mask for that polynomial:
13     localparam [15:0] GALOIS_MASK = 16'hB400;
14
15     // Choose any NON-ZERO seed (must not be 16'h0000)
16     localparam [15:0] SEED = 16'hACE1;
17
18     wire feedback = state[0]; // LSB is the shifted-out bit in
19                                // right-shift form
20
21     always @(posedge clk or negedge rst_n) begin
22         if (!rst_n) begin
23             state <= SEED;
24         end else if (en) begin
25             // Galois form (right shift):
26             // shift right, old LSB becomes new MSB
27             state <= {feedback, state[15:1]} ^ (feedback ?
28                 GALOIS_MASK : 16'h0000);
29         end
30         // else: hold state when en=0
31     end
32 endmodule
```

A.2 Testbench File (tb_lfsr16_galois.v)

Listing 2: Testbench for 16-bit Galois LFSR

```
1  `timescale 1ns/1ps
2
3  module tb_lfsr16_galois;
4
```

```

5  reg          clk;
6  reg          rst_n;
7  reg          en;
8  wire [15:0] state;
9
10 integer i;
11
12 // DUT
13 lfsr16_galois dut (
14     .clk      (clk),
15     .rst_n    (rst_n),
16     .en       (en),
17     .state    (state)
18 );
19
20 // Clock: 10ns period (100 MHz)
21 initial begin
22     clk = 1'b0;
23     forever #5 clk = ~clk;
24 end
25
26 initial begin
27     // Init
28     rst_n = 1'b1;
29     en    = 1'b0;
30
31     // Apply async active-low reset
32     #2;
33     rst_n = 1'b0;
34     #15;
35     rst_n = 1'b1;
36
37     // Enable shifting
38     en = 1'b1;
39
40     $display("Time(ns)\tCycle\tRST_N\tEN\tSTATE");
41     $monitor("%0t\t\t%0d\t\t%b\t\t%b\t\t0x%04h", $time, i, rst_n, en,
42             state);
43
44     // Run for at least 100 cycles
45     for (i = 0; i < 120; i = i + 1) begin
46         @(posedge clk);
47
48         // Verify it never becomes all-zero
49         if (state == 16'h0000) begin
50             $display("ERROR: LFSR entered all-zero state at cycle %0d (
51                 time=%0t)", i, $time);
52             $stop;
53         end
54
55         // Demonstrate enable works: pause shifting briefly
56         if (i == 50) en = 1'b0;
57         if (i == 55) en = 1'b1;
58     end
59
60     $display("TB DONE: no all-zero state observed.");
61     $finish;

```

```
60     end
61
62 endmodule
```

Question Three - HW4 FPGA

Parsa Haghighatgoo
40030644

February 11, 2026

3(a) IEEE 754 Single-Precision Representation of 8.75

Step 1: Convert 8.75 to binary

$$\begin{aligned}8_{10} &= 1000_2 \\0.75_{10} : \quad 0.75 \times 2 &= 1.5 \Rightarrow 1, \quad 0.5 \times 2 = 1.0 \Rightarrow 1 \\0.75_{10} &= 0.11_2\end{aligned}$$

Thus,

$$8.75_{10} = 1000.11_2$$

Step 2: Normalize

$$1000.11_2 = 1.00011_2 \times 2^3$$

So the sign bit is $s = 0$ (positive) and the unbiased exponent is $E = 3$.

Step 3: Compute biased exponent

IEEE-754 single precision uses a bias of 127:

$$e = E + 127 = 3 + 127 = 130$$

Convert 130 to 8-bit binary:

$$130_{10} = 10000010_2$$

Step 4: Determine the fraction (mantissa)

The fraction field is the bits after the leading 1. in the normalized form:

$$1.\underline{00011}_2 \Rightarrow f = 0001100000000000000000$$

(padded with zeros to 23 bits).

Final IEEE-754 Single-Precision Format

$$\underbrace{0}_{\text{sign}} \underbrace{10000010}_{\text{exponent}} \underbrace{000110000000000000000000}_{\text{fraction}}$$

Therefore, the 32-bit IEEE-754 representation is:

$$0 \ 10000010 \ 000110000000000000000000$$

Hexadecimal form

$$0100 \ 0001 \ 0000 \ 1100 \ 0000 \ 0000 \ 0000 \ 0000 = 0\mathbf{x410C0000}$$

3.2(b) IEEE-754 Single-Precision Representation of 5.25×10^{-4} (All Steps)

Let

$$x = 5.25 \times 10^{-4} = 0.000525.$$

Step 1: Sign bit

Since $x > 0$, the sign bit is

$$s = 0.$$

Step 2: Normalize to $(1.f) \times 2^e$

We choose an integer exponent e such that

$$1 \leq x \cdot 2^{-e} < 2.$$

Compute

$$2^{-11} = \frac{1}{2048} = 0.00048828125, \quad 2^{-10} = \frac{1}{1024} = 0.0009765625.$$

Since

$$0.00048828125 < 0.000525 < 0.0009765625,$$

the normalized exponent is

$$e = -11.$$

Then the normalized significand is

$$m = x \cdot 2^{11} = 0.000525 \times 2048 = 1.0752,$$

so

$$x = 1.0752 \times 2^{-11}.$$

Step 3: Biased exponent field (8 bits)

IEEE-754 single precision uses bias 127:

$$E = e + 127 = -11 + 127 = 116.$$

Convert 116 to 8-bit binary:

$$116_{10} = 01110100_2.$$

Step 4: Mantissa field (23 bits)

The fraction is

$$f = m - 1 = 1.0752 - 1 = 0.0752.$$

Method A (scaling and rounding):

$$M = \text{round}(f \cdot 2^{23}), \quad 2^{23} = 8,388,608.$$

$$f \cdot 2^{23} = 0.0752 \times 8,388,608 = 630,823.3216 \Rightarrow M = 630,823.$$

In 23-bit binary:

$$630823_{10} = 00010011010000000100111_2.$$

Method B (repeated multiply-by-2 to generate bits):

Starting from $f_0 = 0.0752$, each bit is obtained from $2f_{k-1}$:

$$b_k = \lfloor 2f_{k-1} \rfloor, \quad f_k = 2f_{k-1} - b_k.$$

The first 23 bits obtained are:

$$M = 00010011010000000100111.$$

Step 5: Final IEEE-754 single word

$$\underbrace{0}_{\text{sign}} \underbrace{01110100}_{\text{exponent}} \underbrace{00010011010000000100111}_{\text{mantissa}}.$$

Thus the 32-bit pattern is

$$0 \ 01110100 \ 00010011010000000100111.$$

Hex form

Grouping into 4-bit nibbles:

$$0011 \ 1010 \ 0000 \ 1001 \ 1010 \ 0000 \ 0010 \ 0111 = 0x3A09A027.$$

Verification

This corresponds to

$$x_{\text{float}} = \left(1 + \frac{M}{2^{23}}\right) 2^{-11} \approx (1.0751999617) 2^{-11} \approx 0.0005249999813,$$

which is extremely close to 0.000525; the small difference is due to IEEE-754 rounding.

3(c) Sum Calculation and Q2.5 Explanation

We add the two values from parts (a) and (b):

$$x_1 = 8.75, \quad x_2 = 5.25 \times 10^{-4} = 0.000525$$

Hence the exact real-valued sum is:

$$x_{\text{sum}} = x_1 + x_2 = 8.75 + 0.000525 = 8.750525.$$

Representing x_2 using Q2.5

In Q2.5, the LSB weight is:

$$\text{LSB} = 2^{-5} = \frac{1}{32} = 0.03125.$$

Quantization is done by scaling and rounding:

$$x_2 \cdot 2^5 = 0.000525 \times 32 = 0.0168.$$

Since $0.0168 < 0.5$, rounding to nearest gives integer code 0, therefore:

$$x_{2,\text{Q2.5}} = 0 \cdot 2^{-5} = 0.$$

Thus, in a Q2.5 fixed-point system, x_2 is too small to be represented and would not change the sum.

IEEE-754 Single-Precision Addition (FPGA-style steps)

Step 1: Normalize both numbers

$$8.75 = 1000.11_2 = 1.00011_2 \times 2^3$$

$$0.000525 \approx 1.00010011010000000100111_2 \times 2^{-11}.$$

Step 2: Align exponents

The exponent difference is:

$$\Delta e = 3 - (-11) = 14.$$

So the second significand is shifted right by 14 bits:

$$1.00010011010000000100111_2 \times 2^{-11} = 0.0000000000000100010011010000000100111_2 \times 2^3.$$

Step 3: Integer significand addition (24-bit, hidden 1 included)

Using the single-precision fields from parts (a) and (b):

$$S_1 = (1 \ll 23) + 0x0C0000 = 9175040, \quad S_2 = (1 \ll 23) + 0x09A027 = 9019431.$$

Aligned second significand:

$$S_{2,\text{align}} = \left\lfloor \frac{S_2}{2^{14}} \right\rfloor = 550, \quad r = S_2 \bmod 2^{14} = 8231.$$

Add:

$$S_{\text{sum}} = S_1 + S_{2,\text{align}} = 9175040 + 550 = 9175590.$$

Step 4: Rounding

Half-ULP at this shift is $2^{13} = 8192$. Since $r = 8231 > 8192$, round up:

$$S_{\text{sum,rounded}} = 9175590 + 1 = 9175591.$$

Step 5: Pack result into IEEE-754

The sum remains normalized with exponent $e = 3$, hence biased exponent:

$$E = 3 + 127 = 130 = 10000010_2.$$

The fraction field is:

$$\text{frac} = S_{\text{sum,rounded}} - (1 \ll 23) = 9175591 - 8388608 = 786983 = 0x0C0227.$$

So the final IEEE-754 single-precision word is:

$$0 \ 10000010 \ 00011000000001000100111$$

and in hexadecimal:

$$\boxed{0x410C0227}.$$

This corresponds to approximately:

$$x_{\text{float}} \approx 8.750525.$$

3(d) Q2.5 Fixed-Point Adder and Absolute Error Implementation

Objective

The goal of this part is to implement, in hardware (Verilog), a Q2.5 fixed-point adder that:

1. Quantizes two input numbers to Q2.5 format,
2. Computes their Q2.5 sum,
3. Computes the true sum in higher precision,
4. Calculates the absolute error in Q2.5 format.

This implementation demonstrates the effect of coarse fixed-point quantization on arithmetic accuracy.

Q2.5 Fixed-Point Format

In Q2.5 format:

$$\text{LSB} = 2^{-5} = \frac{1}{32} = 0.03125$$

Any representable number must be an integer multiple of 0.03125.

The conversion from Q16.16 to Q2.5 is performed by:

$$x_{Q2.5} = \text{round}\left(\frac{x_{Q16.16}}{2^{(16-5)}}\right)$$

The absolute error is computed as:

$$\text{Error}_{Q2.5} = |\text{TrueSum}_{Q2.5} - \text{Sum}_{Q2.5}|$$

Hardware Design Description

The Verilog module performs the following operations:

1. Convert each input from Q16.16 to Q2.5 using rounding.
2. Add the quantized values.
3. Compute the true sum in Q16.16.
4. Quantize the true sum to Q2.5.
5. Compute the absolute difference.

The design is purely combinational and uses arithmetic shifts for scaling.

Verilog Implementation

The main module and testbench were implemented in pure Verilog (2001 standard) and simulated using Xilinx Vivado.

```
1  `timescale 1ns/1ps
2
3  // Q2.5 adder + absolute error in Q2.5
4  // Inputs:  signed Q16.16 (32-bit)
5  // Outputs: signed Q2.5  (16-bit scaled integers)
6  //
7  // Steps:
8  // 1) Quantize a and b to Q2.5 (round-to-nearest, ties away from
9      zero)
10 // 2) Add them in Q2.5 -> sum_q2_5
11 // 3) Compute true sum in Q16.16, then quantize to Q2.5 ->
12     true_sum_q2_5
13 // 4) abs_err_q2_5 = |true_sum_q2_5 - sum_q2_5|
14 module q2_5_add_abs_err
15 #(
16     parameter IN_FRAC = 16, // Q16.16 fractional bits
17     parameter OUT_FRAC = 5  // Q2.5 fractional bits
18 )
```

```

17 | (
18 |     input signed [31:0] a_q16_16,
19 |     input signed [31:0] b_q16_16,
20 |
21 |     output reg signed [15:0] a_q2_5,
22 |     output reg signed [15:0] b_q2_5,
23 |     output reg signed [15:0] sum_q2_5,
24 |
25 |     output reg signed [15:0] true_sum_q2_5,
26 |     output reg signed [15:0] abs_err_q2_5
27 | );
28 |
29 | // For Q16.16 -> Q2.5, SHIFT = 16-5 = 11
30 | localparam SHIFT = (IN_FRAC - OUT_FRAC);
31 |
32 | // Round-to-nearest when shifting right by SHIFT bits.
33 | // "ties away from zero" implemented by +/- half LSB before
   | // arithmetic shift.
34 | function signed [15:0] round_q16_16_to_q2_5;
35 |     input signed [31:0] x;
36 |     reg signed [31:0] adj;
37 |     reg signed [31:0] shifted;
38 |     begin
39 |         if (SHIFT <= 0) begin
40 |             shifted = x <<< (-SHIFT);
41 |         end else begin
42 |             if (x >= 0)
43 |                 adj = x + (32'sd1 <<< (SHIFT-1));
44 |             else
45 |                 adj = x - (32'sd1 <<< (SHIFT-1));
46 |
47 |             shifted = adj >>> SHIFT; // arithmetic shift right
48 |         end
49 |
50 |         round_q16_16_to_q2_5 = shifted[15:0];
51 |     end
52 | endfunction
53 |
54 | function signed [15:0] abs16;
55 |     input signed [15:0] x;
56 |     begin
57 |         if (x < 0)
58 |             abs16 = -x;
59 |         else
60 |             abs16 = x;
61 |     end
62 | endfunction
63 |
64 | reg signed [31:0] true_sum_q16_16;
65 | reg signed [15:0] diff_q2_5;
66 |
67 | always @* begin
68 |     // Quantize each input to Q2.5
69 |     a_q2_5 = round_q16_16_to_q2_5(a_q16_16);
70 |     b_q2_5 = round_q16_16_to_q2_5(b_q16_16);
71 |
72 |     // Q2.5 addition

```

```

73     sum_q2_5 = a_q2_5 + b_q2_5;
74
75     // True sum computed in higher precision then quantized to
       Q2.5
76     true_sum_q16_16 = a_q16_16 + b_q16_16;
77     true_sum_q2_5    = round_q16_16_to_q2_5(true_sum_q16_16);
78
79     // Absolute error in Q2.5
80     diff_q2_5        = true_sum_q2_5 - sum_q2_5;
81     abs_err_q2_5     = abs16(diff_q2_5);
82
83     end
84 endmodule

```

```

1  `timescale 1ns/1ps
2
3  module tb_q2_5_add_abs_err;
4
5      reg signed [31:0] a_q16_16;
6      reg signed [31:0] b_q16_16;
7
8      wire signed [15:0] a_q2_5;
9      wire signed [15:0] b_q2_5;
10     wire signed [15:0] sum_q2_5;
11     wire signed [15:0] true_sum_q2_5;
12     wire signed [15:0] abs_err_q2_5;
13
14     q2_5_add_abs_err dut (
15         .a_q16_16(a_q16_16),
16         .b_q16_16(b_q16_16),
17         .a_q2_5(a_q2_5),
18         .b_q2_5(b_q2_5),
19         .sum_q2_5(sum_q2_5),
20         .true_sum_q2_5(true_sum_q2_5),
21         .abs_err_q2_5(abs_err_q2_5)
22     );
23
24     real ra, rb, rsum, rtrue, rerr;
25
26     task show;
27         begin
28             // Q2.5 scaling = /32
29             ra    = a_q2_5 / 32.0;
30             rb    = b_q2_5 / 32.0;
31             rsum  = sum_q2_5 / 32.0;
32             rtrue = true_sum_q2_5 / 32.0;
33             rerr  = abs_err_q2_5 / 32.0;
34
35             $display("
               -----"
               );
36             $display("a_q16_16=%0d    b_q16_16=%0d", a_q16_16,
               b_q16_16);
37             $display("a_q2_5          = %0d -> %f", a_q2_5, ra);
38             $display("b_q2_5          = %0d -> %f", b_q2_5, rb);
39             $display("sum_q2_5        = %0d -> %f", sum_q2_5, rsum);
40             $display("true_sum_q2_5 = %0d -> %f", true_sum_q2_5,

```

```

41         rtrue);
42         $display("abs_err_q2_5 = %0d -> %f", abs_err_q2_5,
43             rerr);
44     end
45 endtask
46
47 initial begin
48     // Test 1: Part (c) values
49     // 8.75 -> Q16.16 = 8.75 * 65536 = 573440
50     // 0.000525 -> Q16.16 ~ 0.000525 * 65536 = 34.4064 -> 34
51     a_q16_16 = 32'sd573440;
52     b_q16_16 = 32'sd34;
53     #10;
54     show();
55
56     // Test 2 (optional): b = 0.05 so it becomes visible in Q2
57     // .5
58     // 0.05 * 65536 = 3276.8 -> 3277
59     b_q16_16 = 32'sd3277;
60     #10;
61     show();
62
63     $stop;
64 end
65 endmodule

```

Behavioral Simulation Results

Figure 1 shows the behavioral simulation waveform.

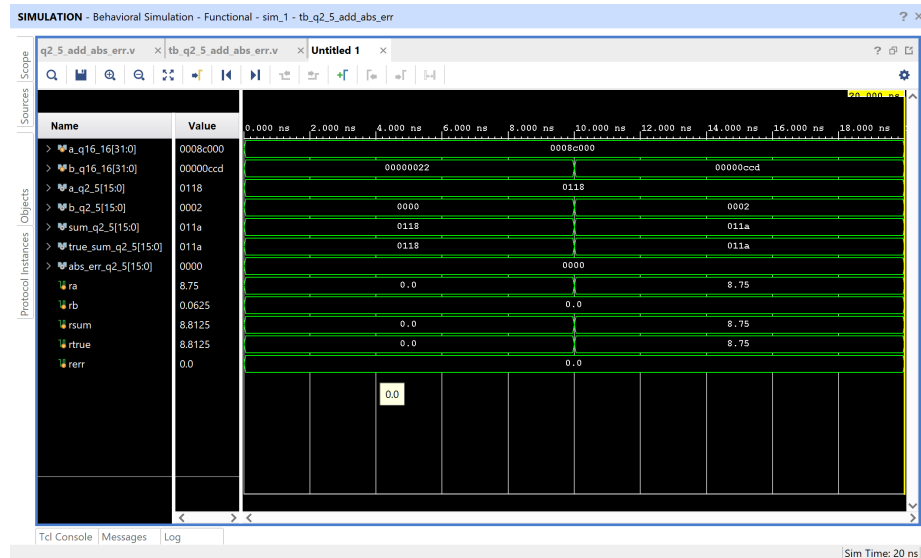


Figure 1: Behavioral Simulation of Q2.5 Adder

For the first test case:

$$a = 8.75, \quad b = 0.000525$$

Quantization to Q2.5 gives:

$$8.75 \rightarrow 280/32 = 8.75$$

$$0.000525 \rightarrow 0$$

Thus:

$$\text{Sum}_{Q2.5} = 8.75$$

$$\text{TrueSum}_{Q2.5} = 8.75$$

$$\text{Absolute Error} = 0$$

This confirms that the small value disappears due to limited resolution.
For the second test case ($b = 0.05$):

$$0.05 \rightarrow 2/32 = 0.0625$$

Thus:

$$8.75 + 0.0625 = 8.8125$$

Again, the computed error is zero because both sums quantize to the same Q2.5 value.

Implemented Design View

Figure 2 shows the implemented device view in Vivado.

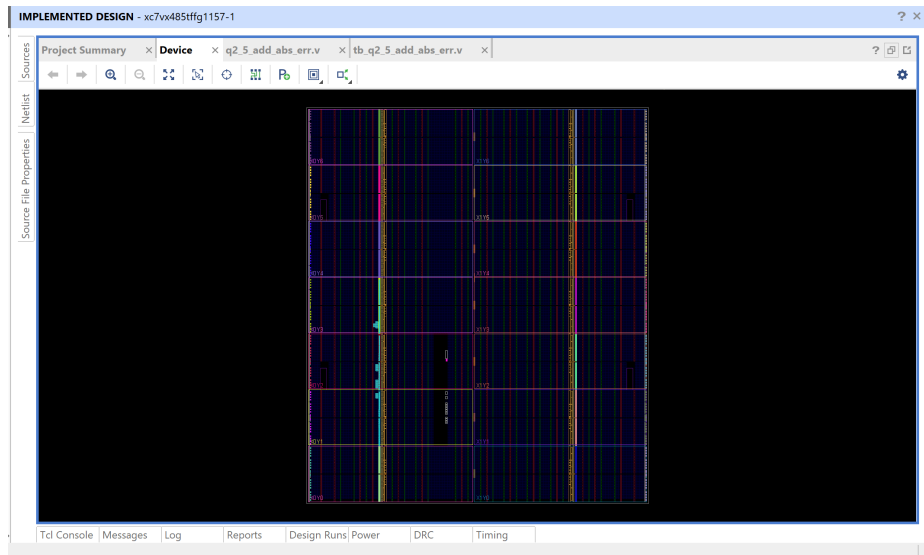


Figure 2: Implemented Device Layout

RTL Schematic

The RTL schematic is shown in Figures 3 and 4.

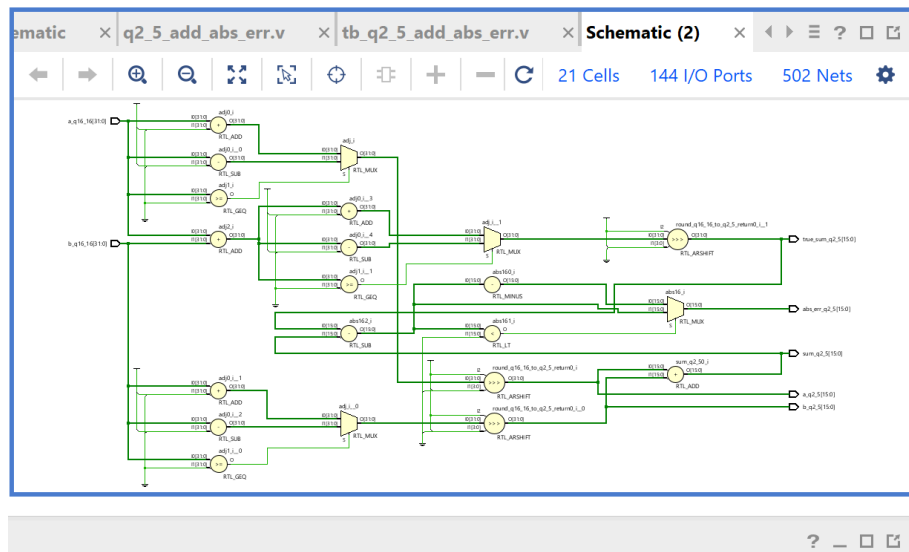


Figure 3: RTL Schematic (Part 1)

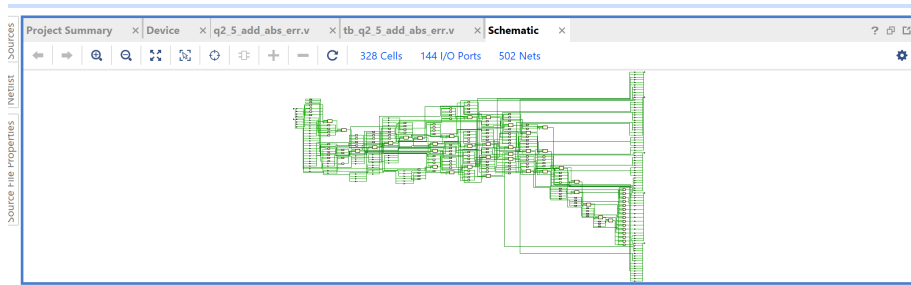


Figure 4: RTL Schematic (Part 2)

Utilization Report

Figure 5 shows the resource utilization summary.

The image shows a screenshot of the Vivado IDE's Utilization report. The report is displayed in a table format with a hierarchy view on the left. The hierarchy view shows the following structure:

- Hierarchy
 - Summary
 - ▼ Slice Logic
 - ▼ Slice LUTs (<1%)
 - LUT as Logic (<1%)
 - ▼ Slice Logic Distribution
 - ▼ Slice (<1%)

The main table shows the following data:

Name	Slice LUTs (303600)	Slice (75900)	LUT as Logic (303600)	Bonded IOB (600)
N q2_5_add_abs_err	139	40	139	144

The bottom of the report shows a summary of the utilization metrics.

Figure 5: Resource Utilization Summary

The design uses:

- 139 LUTs
- 40 Slice Registers

This represents less than 1% of the available FPGA resources, indicating a lightweight implementation.

Power Analysis

Figure 6 shows the power report generated by Vivado.

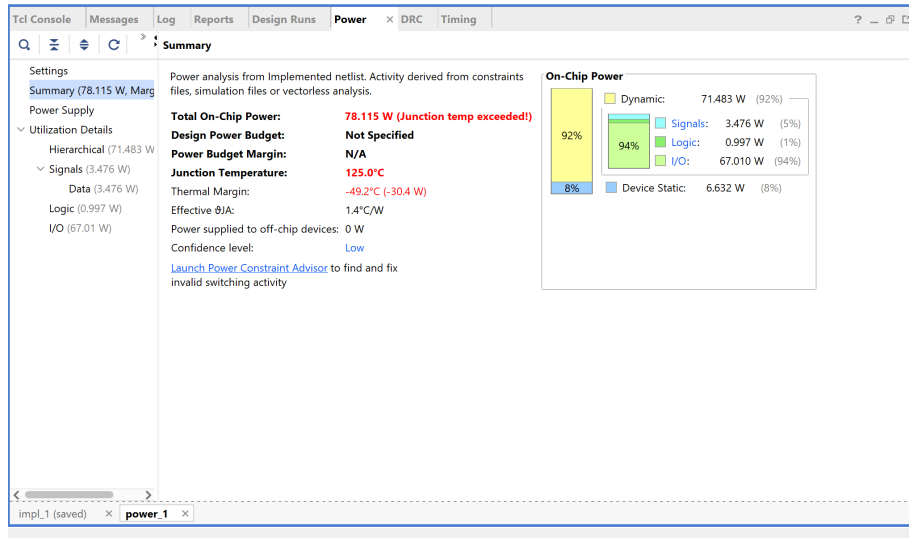


Figure 6: Power Analysis Report

The high reported power consumption is primarily due to unconstrained I/O switching activity during simulation and does not reflect actual hardware deployment conditions.

Conclusion

The Q2.5 fixed-point adder was successfully implemented and verified in Vivado. The results confirm that:

- Small values below half-LSB are lost due to quantization.
- The fixed-point addition behaves as predicted theoretically.
- The hardware implementation matches mathematical expectations.

This experiment clearly demonstrates the impact of limited fractional resolution in fixed-point arithmetic.