

Exploring Int-8 Optimization Methods for Multiply-Accumulation in Neural Networks

By Jan Silva and Oleh Kondratyuk
June 2021

Timing Analysis of Multiplier

Files: *int8.sv*

IP Modules used: *clk_wiz_0*, *xbip_dsp48_macro_0*

Using one instance of the DSP 48 Macro and a PLL generated by a clock wizard IP module, the highest clock frequency for one multiplication is 645 MHz. Two test weights and one test input were hard-coded.

```
// in_i * (w_k + w_j)
// A and D are inputs to pre-adder
xbip_dsp48_macro_0 dsp_slice (
    .CLK(clk),           // input wire CLK
    .A({weight_j, {17{1'b0}}}), // input wire [24 : 0] A - weight
    .B({{10{1'b0}}, input_i}), // input wire [17 : 0] B - input
    .C(accum),           // input wire [47 : 0] C - 3rd input -
                        // used as another addition here for accumulation
    .D({{10{1'b0}}, weight_k}), // input wire [17 : 0] D - weight
    .P(P)                // output wire [47 : 0] P
);
```

Figure 1: DSP Macro Instantiation

Highest clock frequency = 645 Hz

Counter pipeline stages = 3

Output pipeline stages = 2

Fine-Tuning the Packed-Weight Multiplication

Files: *macc.sv*

IP Modules: *xbip_dsp48_macro_0*

Note: This report mainly focuses on signed numbers.

The figure below shows the implementation of packed-weight multiplication using the DSP48 macro's pre-adder, multiplier, and accumulator.

Using the method above for different bit-widths and worst-case values, we found the estimated number of padding for worst-case packed multiplication. The table below shows the calculated amount of padding bits for each input bit-width. See Appendix section II for more details on how these values were calculated.

Table 1: Padding Analysis

Width	Range (worst-case values)	Padding (bits)
Int-9	-256, -256	-1
Int-8	-128, -128	0
Int-7	-64, -64	3
Int-6	-32, -32	6
Int-5	-16, -16	9

Four Methods of Correction and Accumulation

Files: *macc.sv*, *chain_adder_tree.sv*, *chain.sv*

IP Modules: *xbip_dsp48_macro_0*

During the padding analysis, we found that there are slight quirks when certain combinations of negative and positive two's-complement numbers are multiplied together. For example, when the lower-bit term wk is negative and the upper-bit term and input wj and i are positive, the resulting product term ji is one bit lower than expected.

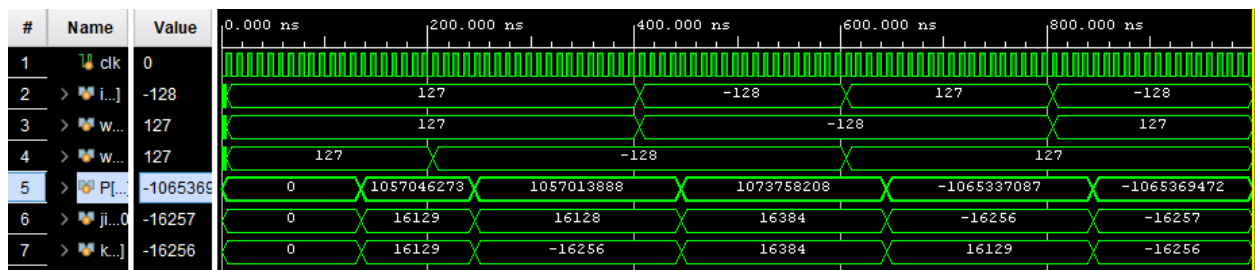


Figure 3: $127 \times (127 + 127)$ results in $ji = ki = 16129$, but $127 \times (127 + -128)$ results in $ji = 16128$ instead of $ji = 16129$.

To correct these quirks, we explored four different methods of multiply-accumulation. First, we examined the difference between using an internal correction and an external adder tree, and an external correction and internal adder tree. We also tested methods

with both correction and no correction to investigate the effects of corrections on resource utilization and timing. If the effects of the corrections are close to negligible, it may be worth it to manually manipulate the weights to avoid such quirks.

- Internal Correction, External Adder Tree
 - The correction for the one-bit off phenomenon is applied through a simple conditional statement in the C port of the DSP. The C port can no longer be used to accumulate results, and an external adder tree can be implemented to accumulate multiplications. See chain.sv for this implementation.
- External Correction, Internal Adder Tree
 - A separate adder tree, external to the DSP slices was needed to accumulate the errors that occur in the process of chaining DSP slices together. Whether a chained pair of DSP slices required a correction (a +1 or -1 to the result) was determined by some logic taking in the sign bit of the 1) Input, 2) w_k 3) preceding DSP result and 4) overall DSP result. Some of these depend on each other, so some pipelining was required to properly apply the correction logic. The final sum of the adder tree is added to the result of the overall adder tree to come to a final, correct, result.
- No correction, External Adder Tree
- No correction, Internal Adder Tree

We tested these cases with 5, 8-bit inputs first, then ran synthesis and implementation with 1024 inputs. This required a new top-level module called test_use_case that read 1024 input values and weights from memory files.

Resource Utilization Analysis

Files: macc.sv, chain_adder_tree.sv, chain.sv, test_use_case.sv

IP Modules: xbip_dsp48_macro_0

Note: The results below are slightly outdated. As of 6/8/2021, we implemented a new binary adder tree which raised the LUT usage for both of the modules significantly.

Utilization - Total LUTs by Module

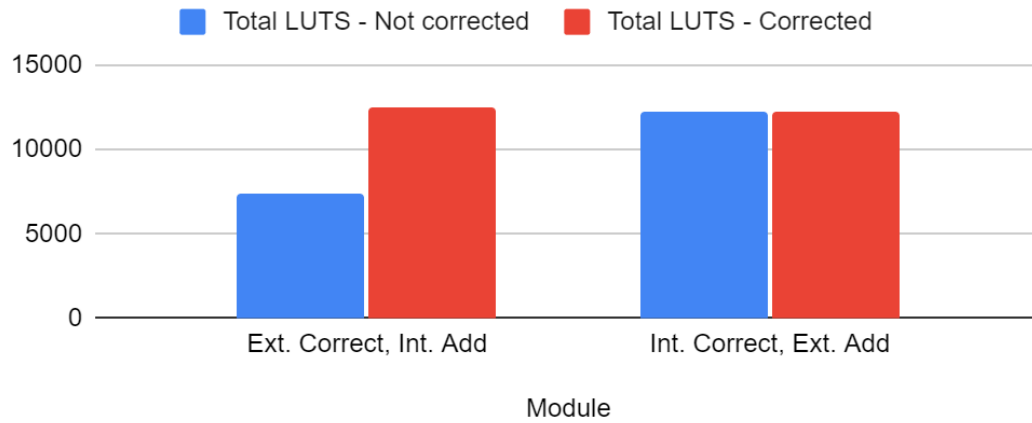


Figure 4: Total LUTs by Module

Utilization with Corrections

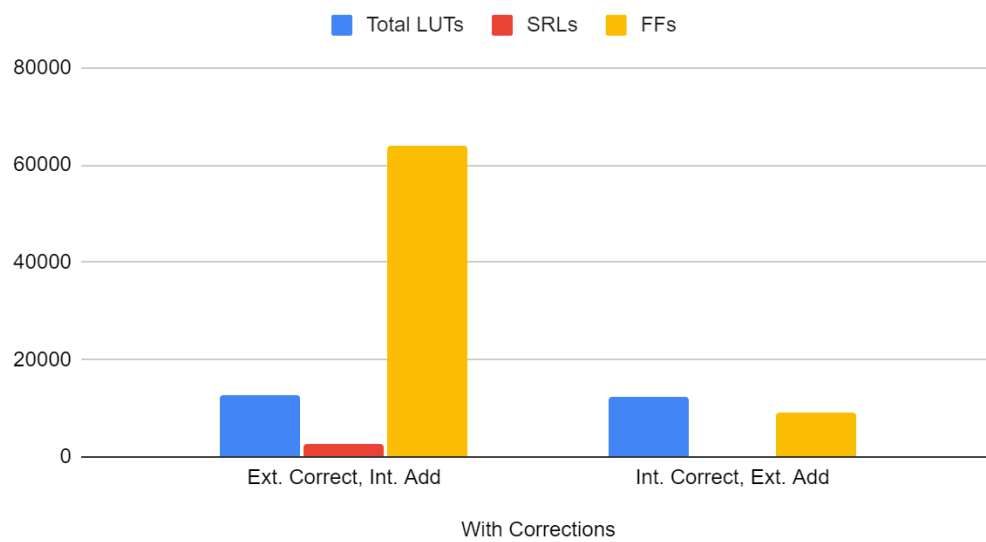


Figure 5: Utilization of each method with corrections

Utilization without Corrections

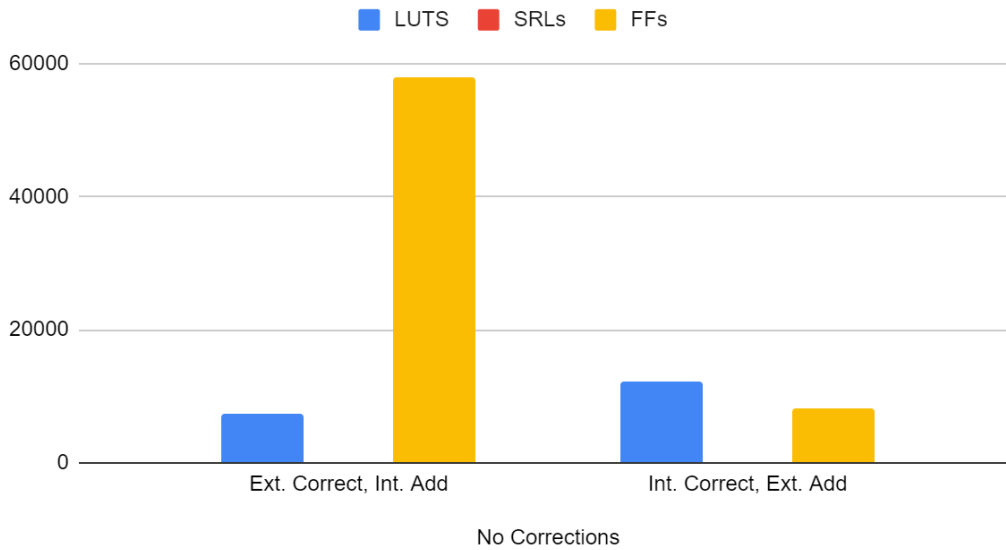


Figure 6: Utilization of each method without corrections

Timing Analysis

Files: *macc.sv*, *chain_adder_tree.sv*, *chain.sv*, *test_use_case.sv*, *timed_case.sv*

IP Modules: *xbip_dsp48_macro_0*, *clk_wiz_0*

For the timing analysis, we created a new top-level module called *timed_case* with a PLL to generate a 400 MHz clock. This module contains an instance of *test_use_case*, which can be modified to implement each of the four methods.

400 MHz

File	Method	Synthesis	Implementation
chain_adder_tree	External correction, internal addition	WNS = 0.778 ns WHS = 0.043 ns WPWS = 0.607 ns	WNS = -0.113 ns WHS = 0.043 ns WPWS = 0.607 ns (Failed timing)
chain	Internal correction, external addition	WNS = 0.855 ns WHS = 0.049 ns WPWS = 0.899 ns	WNS = 0.097 WHS = 0.023 WPWS = 0.849

WNS = Worst Negative Slack, WHS = Worst Hold Slack, WPWS = Worst Pulse Width Slack

Conclusion

From the above resource utilization reports and timing analyses, the internally-corrected multiply-accumulate implementation can handle faster clock speeds. While the total LUT usage of the uncorrected internal adder method is much lower, the method still uses a much higher number of flip-flops. The corrected versions of both the internal adder and external adder methods have similar LUT usage, but the latter uses much less FFs. Thus, the internally corrected method with the external adder tree is more resource-efficient, and along with the potential to go faster than 400MHz, this method seems to be the better one out of the four.

References

https://www.xilinx.com/support/documentation/white_papers/wp486-deep-learning-int8.pdf

Appendices

- I. [Multiplier Timing Spreadsheet](#)
- II. Screenshots of Simulations for analyzing Int-5 to Int-9

A. Calculated binary values

1. $d'127 \times d'127 =$
00000000000000000000000000000000 0111111000000001

2. $D'127 \times -d'128 =$
00000000000000000000000000000000 111111100000000

3. $-d'128 \times d'127 =$
00000000000000000000000000000000 111111100000000
(-111111100000000)

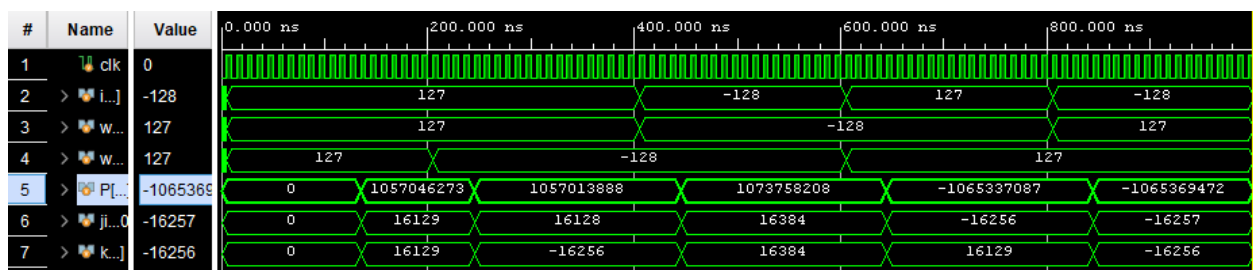
4. $-d'128 \times -d'128 =$
00000000000000000000000000000000 0100000000000000

B. Int-9

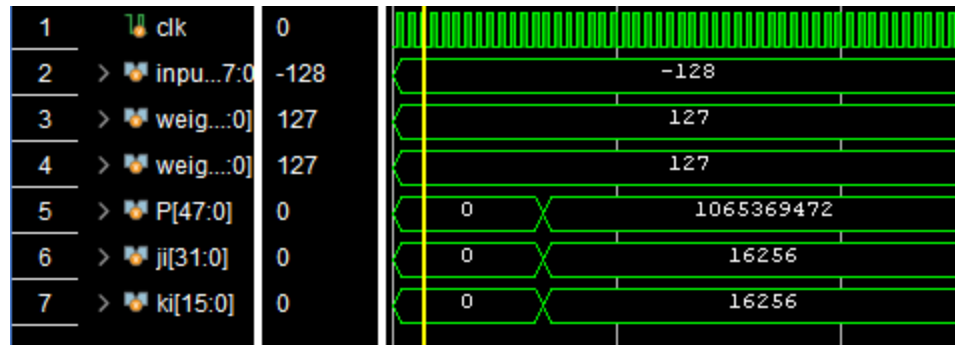
1. 0000000000000000001111111000000010111111000000001
11111110000000001

2.

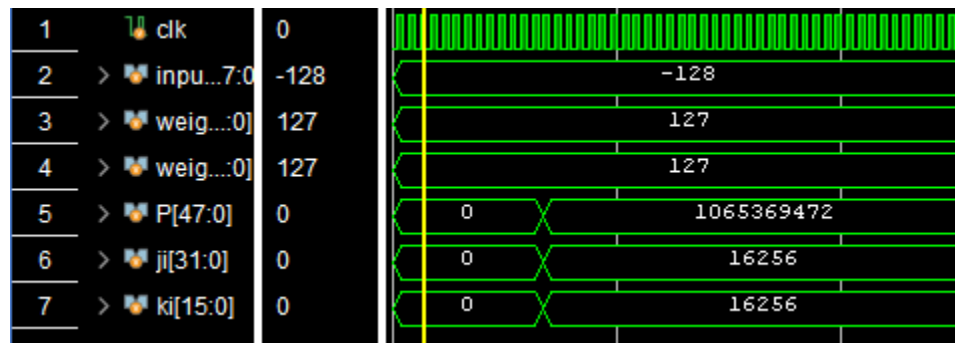
C. Int-8



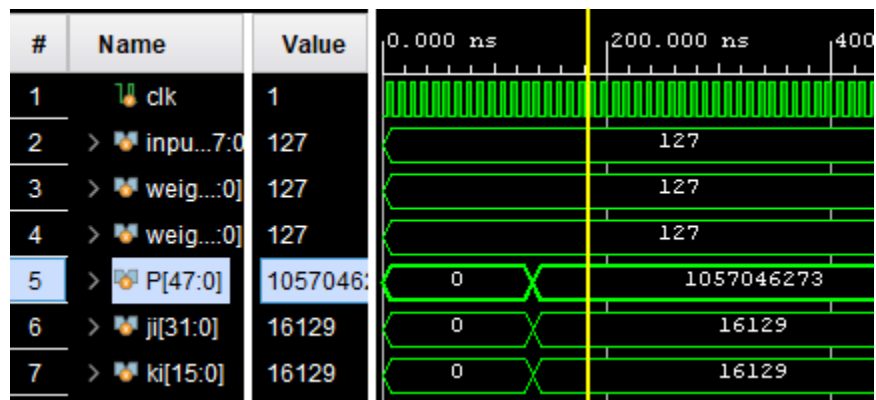
1. $-128 \times (127 + 127)$



2. $127 \times (-128 + -128)$

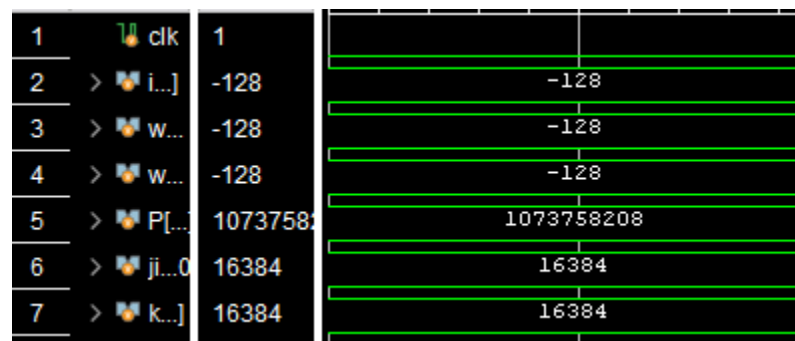


3. $127 \times (127 + 127)$



4. $-128 \times (-128 + -128) =$

1111111111111111 10000000000000 0100000000000000



5. $-128 \times (127 + -128) =$

000000000000000000 11111110000000 0100000000000000

6. $-128 \times (-128 + 127) =$

1111111111111111111100000000000000 00 11111110000000

7. $127 \times (127 + -128)$

1	clk	1	
2	> input_i[7:0]	127	127
3	> weight_j[7:0]	127	127
4	> weight_k[7:0]	-128	-128
5	> P[47:0]	105701388	1057013888
6	> ji[31:0]	16128	16128
7	> ki[15:0]	-16256	-16256

8. $127 \times (-128 + 127)$

#	Name	Value	162,144,990 ps	162,144,995 ps
1	clk	1		
2	> input_i[7:0]	-128	-128	-128
3	> weight_j[7:0]	-128	-128	-128
4	> weight_k[7:0]	127	127	127
5	> P[47:0]	107372556	1073725568	1073725568
6	> ji[31:0]	16383	16383	16383
7	> ki[15:0]	-16256	-16256	-16256

D. Int-7

1. $(-64 \times (-64 + -64))$

2. 00000000000000000000

01 0000 0000 0000 000 01 0000 0000 0000

assign ji = P[47:17]; assign ki = P[13:0];

#	Name	Value	0.000 ns	200.000 ns	400.000 ns	600.000 ns	800.000 ns
1	clk	0					
2	> i...	63	63	63	-64	-64	63
3	> w...	0	0	63	-64	-64	63
4	> w...	63	63	63	-64	-64	63
5	> P[...]	0	0	3969	520228737	520220736	536875008
6	> ji...	0	0	3969	3968	4096	2147479616
7	> k...	0	0	3969	-4032	4096	3969

E. Int-6

1. $(-32 \times (-32 + -32))$

2	> i...]	-32	-32
3	> w...	-32	-32
4	> w...	-32	-32
5	> P[...	26843648	268436480
6	> ji...0	1024	1024
7	> k...]	1024	1024

- 2.
3. 000000000000000000 010000000000 000000
010000000000
4. 6 bit
5. Assign ji = P[47:18], assign ki = P[11:0]

F. Int-5

1. 00000000000000000000 0100000000 0000000000
010000000000
2. 9 bit

G. Int-4

1. 7x(7+7)
2. 0000000000000000000000000000110001 0000 0000 0000 0011 0001
3. 12-bit

III. Timing Analysis Summaries at 400 MHz

A. Chain_adder_tree

1. Synthesis

Setup	Hold	Pulse Width
Worst Negative Slack (WNS): 0.778 ns	Worst Hold Slack (WHS): 0.043 ns	Worst Pulse Width Slack (WPWS): 0.607 ns
Total Negative Slack (TNS): 0.000 ns	Total Hold Slack (THS): 0.000 ns	Total Pulse Width Negative Slack (TPWS): 0.000 ns
Number of Failing Endpoints: 0	Number of Failing Endpoints: 0	Number of Failing Endpoints: 0
Total Number of Endpoints: 238702	Total Number of Endpoints: 238702	Total Number of Endpoints: 178292

All user specified timing constraints are met.

2. Implementation

Design Timing Summary

Setup	Hold	Pulse Width
Worst Negative Slack (WNS): -0.113 ns	Worst Hold Slack (WHS): 0.043 ns	Worst Pulse Width Slack (WPWS): 0.607 ns
Total Negative Slack (TNS): -0.113 ns	Total Hold Slack (THS): 0.000 ns	Total Pulse Width Negative Slack (TPWS): 0.000 ns
Number of Failing Endpoints: 1	Number of Failing Endpoints: 0	Number of Failing Endpoints: 0
Total Number of Endpoints: 110168	Total Number of Endpoints: 110168	Total Number of Endpoints: 86622

Timing constraints are not met.

B. Chain

1. Synthesis

Design Timing Summary

Setup	Hold	Pulse Width
Worst Negative Slack (WNS): 0.855 ns	Worst Hold Slack (WHS): 0.049 ns	Worst Pulse Width Slack (WPWS): 0.899 ns
Total Negative Slack (TNS): 0.000 ns	Total Hold Slack (THS): 0.000 ns	Total Pulse Width Negative Slack (TPWS): 0.000 ns
Number of Failing Endpoints: 0	Number of Failing Endpoints: 0	Number of Failing Endpoints: 0
Total Number of Endpoints: 232465	Total Number of Endpoints: 232465	Total Number of Endpoints: 191819

All user specified timing constraints are met.

2. Implementation

Setup	Hold	Pulse Width
Worst Negative Slack (WNS): 0.097 ns	Worst Hold Slack (WHS): 0.023 ns	Worst Pulse Width Slack (WPWS): 0.849 ns
Total Negative Slack (TNS): 0.000 ns	Total Hold Slack (THS): 0.000 ns	Total Pulse Width Negative Slack (TPWS): 0.000 ns
Number of Failing Endpoints: 0	Number of Failing Endpoints: 0	Number of Failing Endpoints: 0
Total Number of Endpoints: 49057	Total Number of Endpoints: 49057	Total Number of Endpoints: 50087
















All user specified timing constraints are met.

IV. Registers as a Flip Flop

A. Chained_adder_tree

test_use_case	69937
chains(chained_adder_tree)	63655
Leaf Cells (5959)	5959
maccs[0].even_macc(chained_macc__xdcDup__1)	104
dsp_slice(xbip_dsp48_macro_0_HD1024)	104
U0(xbip_dsp48_macro_0_xbip_dsp48_macro_v3_1)	104
i_synth(xbip_dsp48_macro_0_xbip_dsp48_macro_v3_1)	104
i_synth_option.i_synth_model(xbip_dsp48_macro_0_xbip_dsp48_macro_v3_1)	104
i_c4(xbip_dsp48_macro_0_xbip_pipe_v3_1)	48
i_has_c.i_c3(xbip_dsp48_macro_0_xbip_pipe_v3_1)	48
i_op3(xbip_dsp48_macro_0_xbip_pipe_v3_1)	4
i_op4(xbip_dsp48_macro_0_xbip_pipe_v3_1)	4

B. Chain

				Register as Flip Flop
Name	Used			
✓ N timed_case	49057			
✓  test (test_use_case)	49057			
✓  dut (chain)	49057			
 Leaf Cells (3986)	39869			
✓  eachSlice[0].slic	10			
✓  dsp_slice (xb	10			
>  U0 (xbip_c	10			
>  eachSlice[1].slic	10			
>  eachSlice[5].slic	10			
>  eachSlice[6].slic	10			
>  eachSlice[7].slic	10			
>  eachSlice[9].slic	10			

V. Resource Utilization Summaries

A. External adder tree with internal correction (chain)

Resource	Utilization	Available	Utilization %
LUT	83222	433200	19.21
FF	186880	866400	21.57
DSP	1280	3600	35.56
IO	30785	850	3621.76

B. Internal adder tree with external correction (chained_adder_tree)

Resource	Utilization	Available	Utilization %
LUT	81914	433200	18.91
FF	186880	866400	21.57
DSP	1280	3600	35.56
IO	30773	850	3620.35

C. Internal adder tree, no correction (chain)

Resource	Utilization	Available	Utilization %
LUT	64000	433200	14.77
FF	186880	866400	21.57
DSP	1280	3600	35.56
IO	30785	850	3621.76

VI. Early work on observing overflow

01111111 x 01111111

00000000000000000000000000000000 0011111100000001 term1
00000000000000000000000000000000 0111111000000010 term2
00000000000000000000000000000000 1011110100000011 term3
00000000000000000000000000000000 1111110000000100 term4
00000000000000000000000000000001 0011101100000101 term 5 << overflow

10000000 x 10000000

00000000000000000000000000000000 0000000000000000 term 0
00000000000000000000000000000000 0100000000000000 term 1
00000000000000000000000000000000 1000000000000000 term 2
00000000000000000000000000000000 1100000000000000 term 3
00000000000000000000000000000001 0000000000000000 term 4 << overflow

10000000 x 01111111

00000000000000000000000000000000 0000000000000000 term 0
00000000000000000000000000000000 0011111110000000 term 1
00000000000000000000000000000000 0111111100000000 term 2
00000000000000000000000000000000 1011111010000000 term 3
00000000000000000000000000000000 1111111000000000 term 4
00000000000000000000000000000001 0011110110000000 term 5 << overflow

Int 8 opt with 01111111

0000000000000000 0000000000000000 0 0000000000000000 term0
0000000000000000 0011111100000001 0 0011111100000001 term1
0000000000000000 0111111000000010 0 0111111000000010 term2
0000000000000000 1011110100000011 0 1011110100000011 term3
0000000000000000 1111110000000100 0 1111110000000100 term4
0000000000000001 0011101100000101 1 0011101100000101 term5 << overflow

11111111 x 11111111

00000000000000000000000000000000 0 0000000000000000 term0
00000000000000000000000000000000 0 1111111000000001 term1
00000000000000000000000000000000 1 1111110000000010 term2 << overflow