

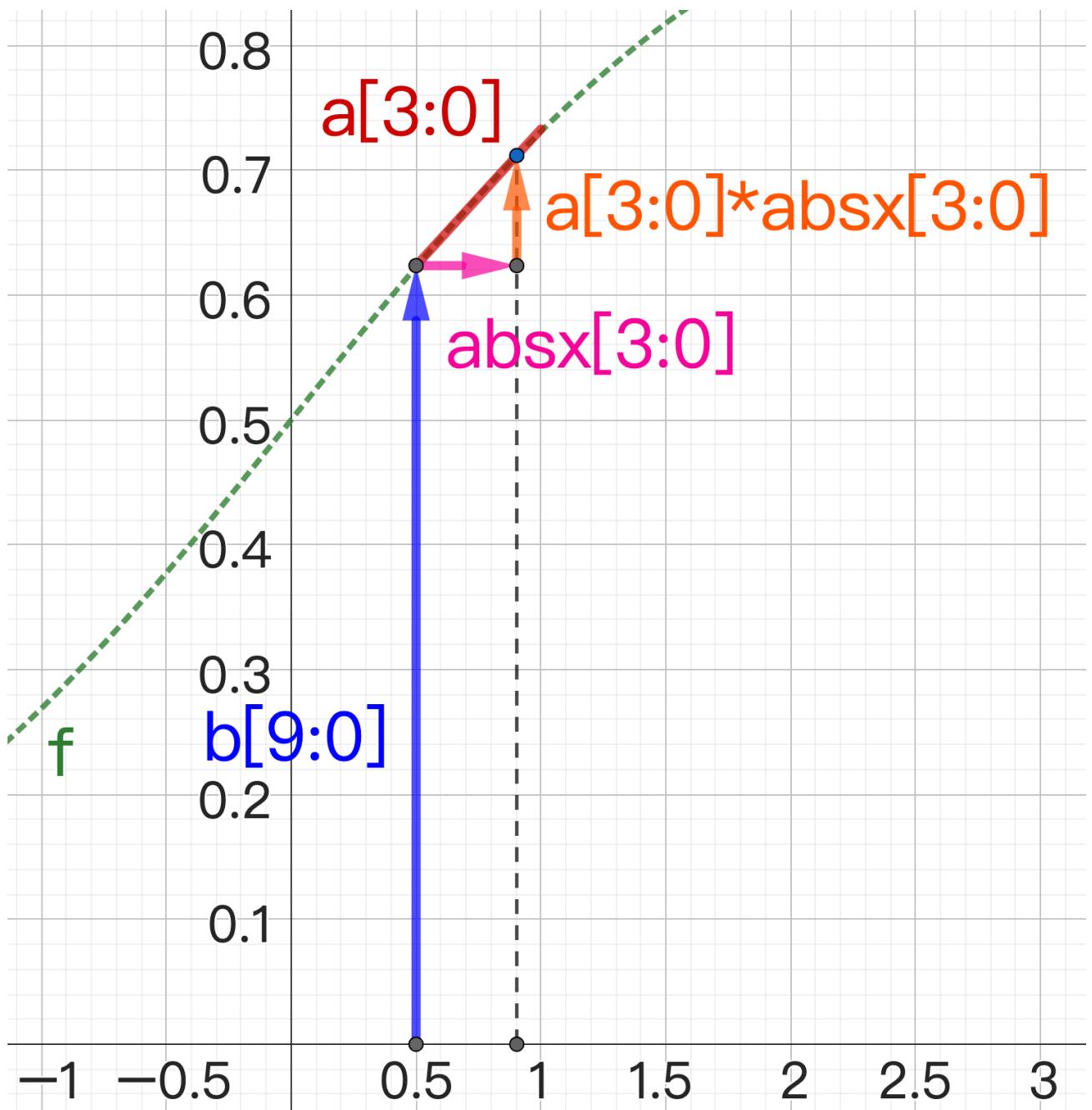
IC Design HW4

B10901176 蔡弘祥

1. Solution Idea

Since sigmoid function has a special property, which is $f(x)+f(-x) = 1$, we only need to handle the non-negative part of the function.

1. Getting the 8-bit absolute value abs_x of the input x .
 - i. The absolute value of most of the negative cases can be expressed in 8 bits.
 - ii. The only case that cannot be expressed is 1000_0000, which is -4 in decimal. In this case, we will directly output the golden value. This will consume lots of resources in the circuit.
2. Deriving the 4-bit slope and the 10-bit intercept by range of abs_x
 - i. Splitting [0,4] into 8 equivalent intervals, which are [0,0.5), [0.5,1), ..., [3.5,4). In this case, we only need the three bits, $\text{abs_x}[6:4]$, to determine the interval that abs_x lies on.
 - ii. Each interval corresponds to a specific slope and intercept, by an 8 to 1 mux bus, we can determine the value of the 4-bit slope, $a[3:0]$, and the 10-bit intercept, $b[9:0]$.
 - iii. We can observe that the first 4 bits, $\text{abs_x}[7:4]$, are the same in every interval no matter abs_x is, so when doing multiplication, we only need to multiply the last 4 bits, $\text{abs_x}[3:0]$, which are 2^{-2} to 2^{-5} .
 - iv. The 4-bit slope are 2^{-3} to 2^{-6} . This is enough precise. More bits will only increase little precision, but cost much more time on calculating.
 - v. The 10-bit intercept are 2^0 to 2^{-11} . Since 2^0 will always be 0 and 2^{-1} will always be 1 when considering the non-negative input of the sigmoid function, so we don't have to consider it when deciding the intercept.
 - vi. Checkout the graph below
3. Calculating the $f(\text{abs_x})$
 - i. The first step is to calculate $a[3:0] * \text{abs_x}[3:0]$. The result is $\text{mul}[7:0]$ with 8 bits, which are 2^{-4} to 2^{-11} .
 - ii. The next step is to calculate $\text{mul}[7:0] + b[9:0]$. The result is $\text{funcOut}[9:0]$ with 10 bits, which are 2^{-2} to 2^{-11} . You might wonder why the result is 10 bits (2^{-2} to 2^{-11}) while not 11 bits (2^{-1} to 2^{-11}). The general case of an 8-bit number plus a 10-bit number is a 11-bit number. But in our case, we don't have to consider about 2^{-1} since 2^{-1} is always 1 in every case. Thus, we can use less computational resources to get the same result.
4. If $x \geq 0$, $y = f(\text{abs_x})$. If $x < 0$, $y = 1 - f(\text{abs_x})$.
 - i. We will first check the sign of the input, and then determine the output.
 - ii. The positive cases are easy to handle, just set $y[13:4]$ to the value we just calculate. Then set $y[15:14]$ to 01, $y[3:0]$ to 0011 (set $y[3:0]$ to 0011 will be the least error case).
 - iii. In general, calculating $1 - f(\text{abs_x})$ is hard. But if we just invert every bit of $f(\text{abs_x})$, we will also get the similar result since the error can be ignored (the error is 2^{-11}). After inverting, just set $y[13:4]$ to the value. Then set $y[15:14]$ to 00, $y[3:0]$ to 1011 (set $y[3:0]$ to 1011 will be the least error case).
 - iv. The last consideration is the case of 1000_000, we will directly output the golden value.



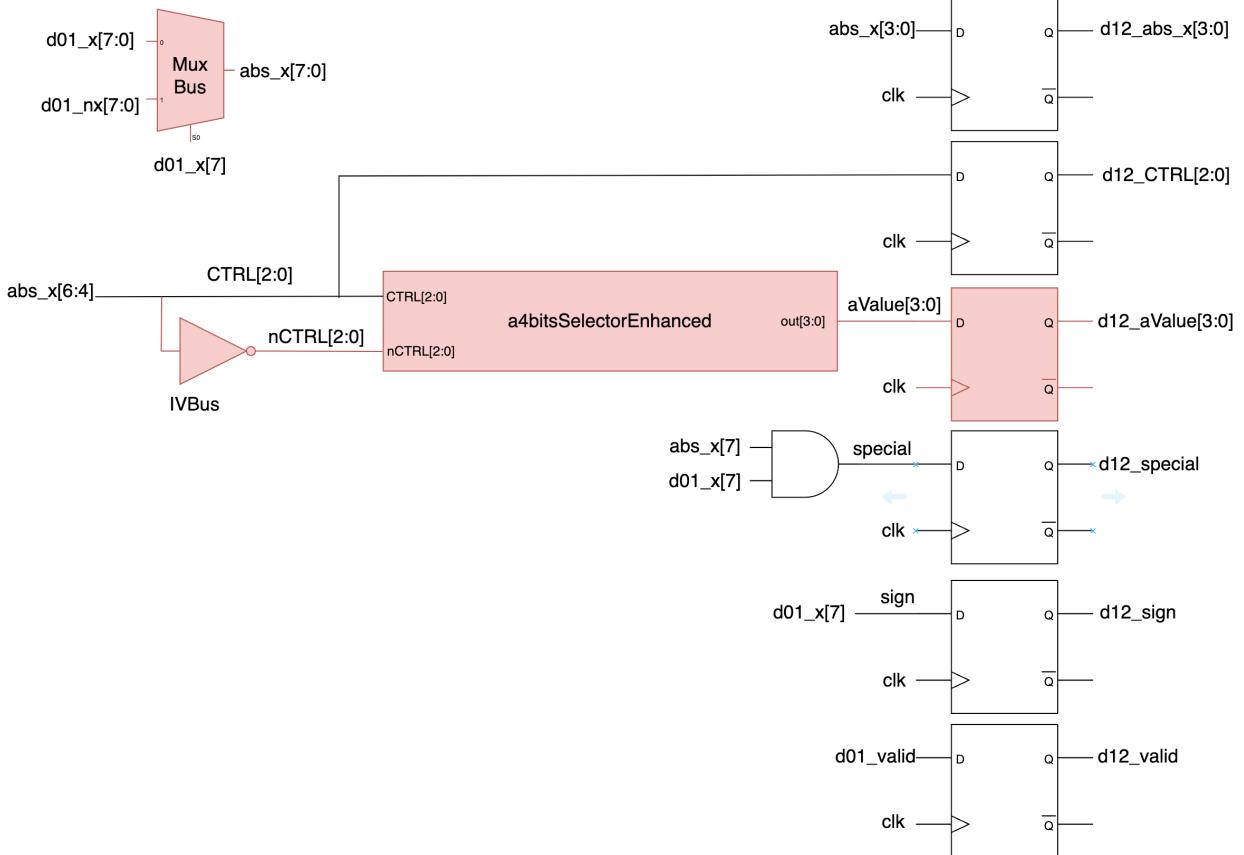
2. Pipeline Stages Overview

Stage	Input Signals	Output Signals	Explanation
0	$i_x[7:0]$ $i_{\text{in_valid}}$	$d01_x[7:0]$ $d01_{\text{nx}}[7:0]$ $d01_{\text{valid}}$	2's complement handling
1	$d01_x[7:0]$ $d01_{\text{nx}}[7:0]$ $d01_{\text{valid}}$	$d12_{\text{abs_x}}[3:0]$ $d12_{\text{CTRL}}[2:0]$ $d12_{\text{aValue}}[3:0]$ $d12_{\text{special}}$ $d12_{\text{sign}}$ $d12_{\text{valid}}$	Absolute value of x deriving Slope deriving

2	d12_abs_x[3:0] d12_CTRL[2:0] d12_aValue[3:0] d12_special d12_sign d12_valid	d23_add01[5:0] d23_add23[5:0] d23_CTRL[2:0] d23_special d23_sign d23_valid	Multiplication Stage 1
3	d23_add01[5:0] d23_add23[5:0] d23_CTRL[2:0] d23_special d23_sign d23_valid	d34_mul[7:0] d34_bValue[9:0] d34_special d34_sign d34_valid	Multiplication Stage 2 Intercept deriving
4	d34_mul[7:0] d34_bValue[9:0] d34_special d34_sign d34_valid	d45_funcOut[9:0] d34_special d34_sign d34_valid	Addition Stage
5	d45_funcOut[9:0] d34_special d34_sign d34_valid	o_y[15:0] o_out_valid	Output handling

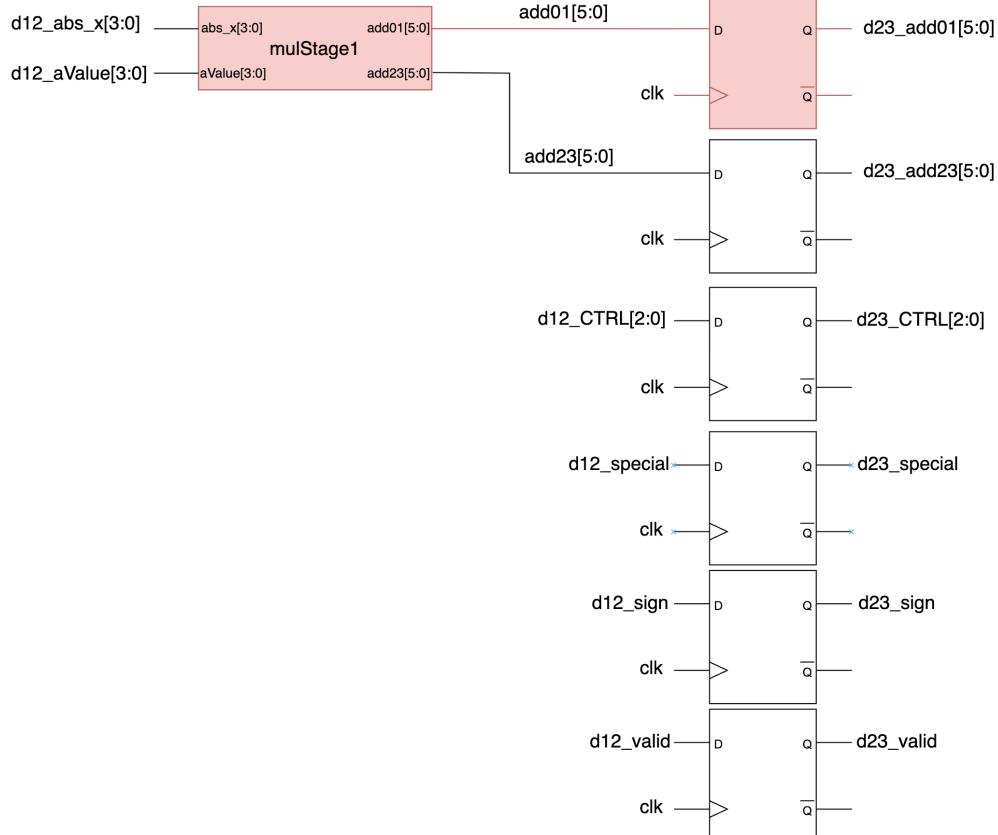
3. Each Stage Circuit Diagram and Combinational Part Explanation

Stage	Circuit Diagram and Combinational Part Explanation (Red Parts Are Critical Paths)		
0	<p>The diagram shows the combinational logic for Stage 0. An input $i_x[7:0]$ is inverted to $nIn[7:0]$, which is then processed by an IVBus block. The output of the IVBus is connected to the $in[7:0]$ port of an addOne block. The addOne block also receives i_in_valid as its clock signal (clk). Its output $out[7:0]$ is connected to the D inputs of three D flip-flops. The Q outputs of these flip-flops are $d01_nx[7:0]$, $d01_x[7:0]$, and $d01_valid$. The clk signals for the flip-flops are also i_in_valid.</p>		
	<p>2's complement handling</p> <p>$d01_nx[7:0]$ is the 2's complement of $i_x[7:0]$. We will regard 1000_0000 as a special case since we cannot express -128's 2's complement in 8 bits. In this special case, we will directly output its golden value.</p>		



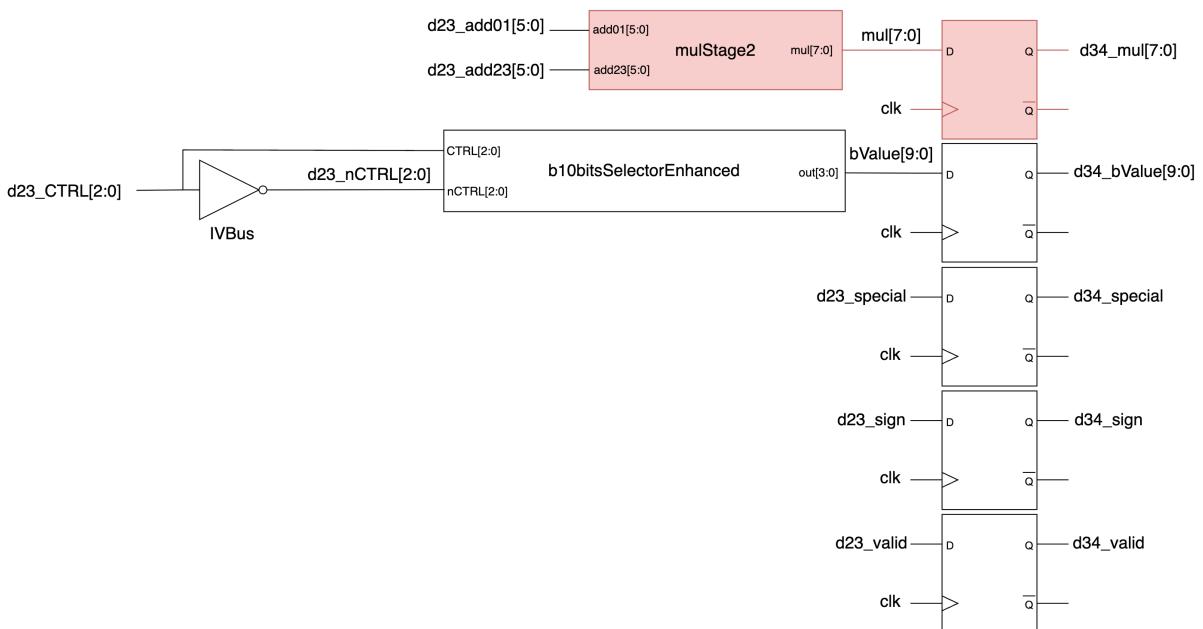
1

Absolute value of x deriving	We can determine whether the input value is negative or not by the most significant bit of the input. To get the absolute value of the input, we only need to use a 2 to 1 mux bus.
Slope deriving	a4bitsSelectorEnhanced is a module that can determine the value of a[3:0]. It used to be an 8-1 mux named a4bitsSelector. Since the inputs are all constants, we can simplify it to lessen the number of transistors and lower the latency.
CTRL	CTRL signal needs to propagate until the intercept is determined.
Special	The special case will be determined here. If it is true, then it means the input is 1000_0000.
Sign	Sign signal needs to propagate to the last stage.



Multiplication Stage 1

Handle the first stage of multiplication. A 4-bit number times a 4-bit number can be regarded as adding 4 numbers, which is three additions. We will handle two additions in advanced and handle the last addition in the next stage.

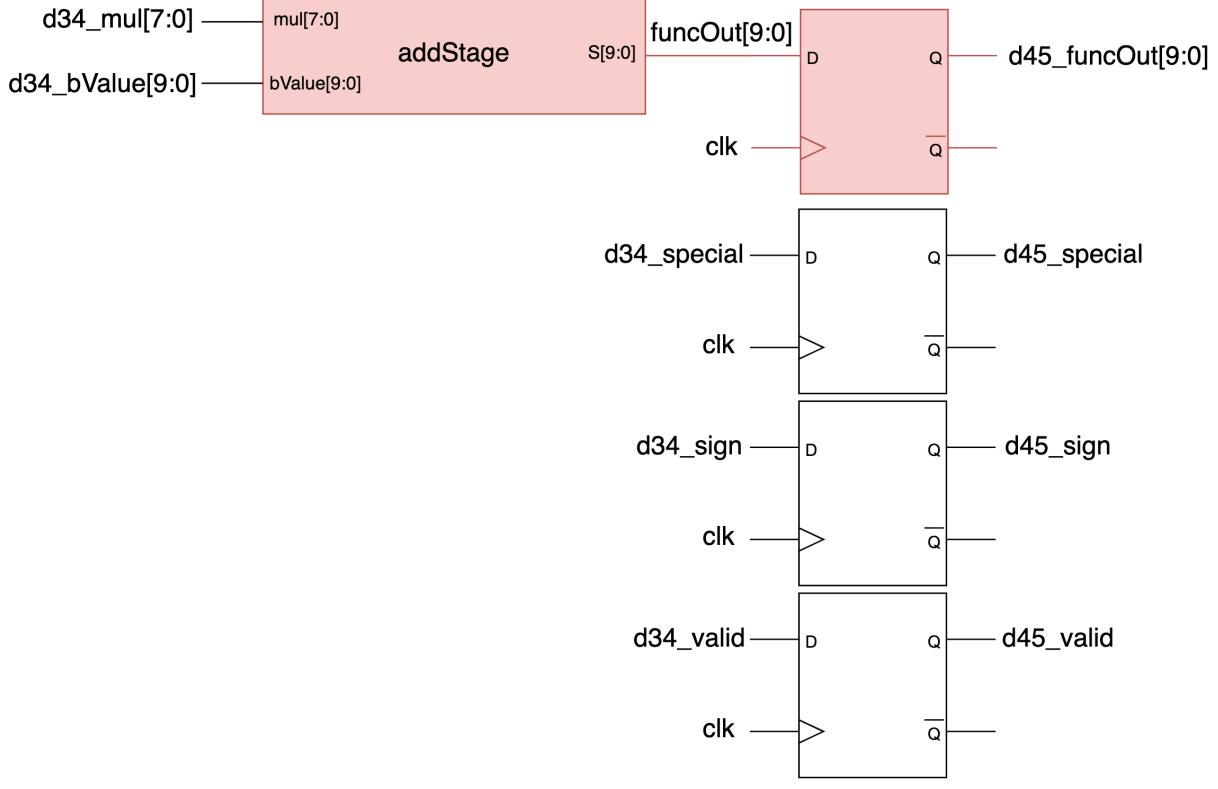


Multiplication Stage 2

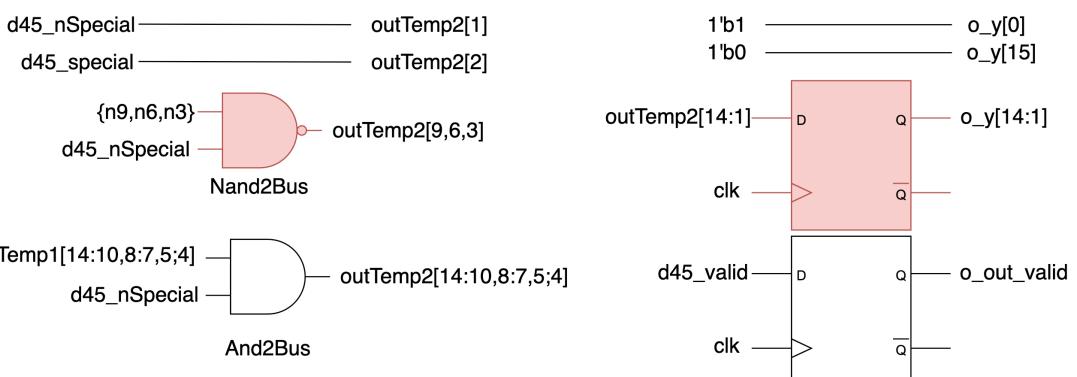
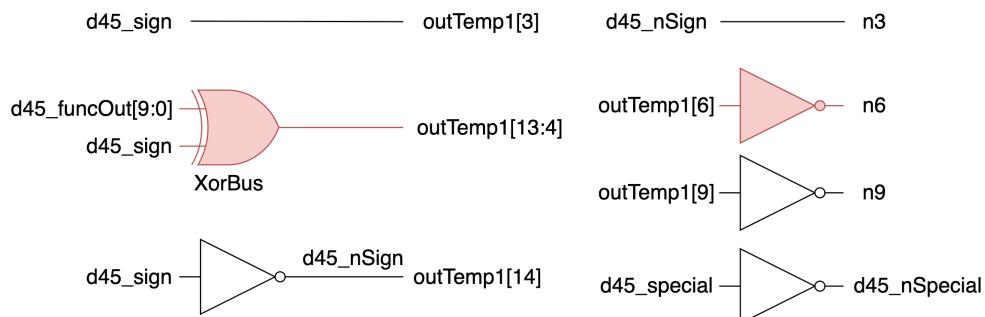
Handle the second stage of multiplication.

Intercept deriving

b10bitsSelectorEnhanced is a module that can determine the value of b[9:0]. It used to be an 8-1 mux named b10bitsSelector. Since the inputs are all constants, we can simplify it to lessen the number of transistors and lower the latency.



Addition Stage | Handle the addition of `mul[7:0]` and `b[9:0]`.



`outTemp1` | `outTemp1` is the output of the input case except (1000_0000).

`outTemp2` | `outTemp2` is the output considering the output of 1000_0000, which is 0000_0010_0100_1101. It used to be a 2 to 1 mux, but we can simplify into some simple logic gates.

Output handling | Send `outTemp2` and `d45_valid` into flip-flops, it will output the result in the next cycle.

4. Self-defined Sub Modules Overview

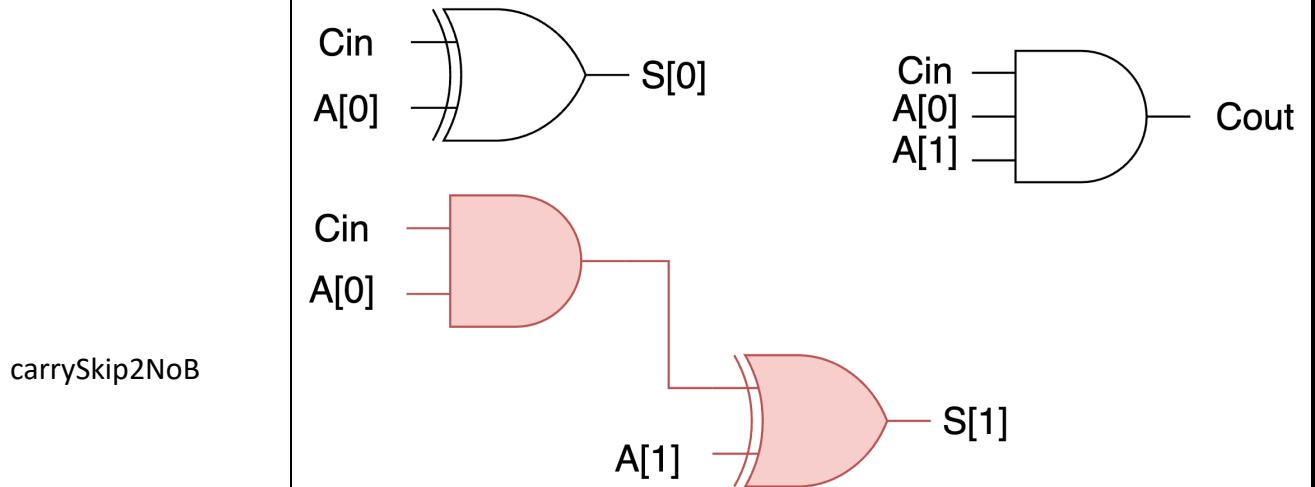
Sub Module Name	Input signals	Output signals	Explanation
carrySkip2	A[1:0] B[1:0] Cin	S[1:0] Cout	Carry skip adder with 2-bit inputs
carrySkip2NoCin	A[1:0] B[1:0]	S[1:0] Cout	Carry skip adder with 2-bit inputs No Cin signal
carrySkip2NoB	A[1:0] Cin	S[1:0] Cout	Carry skip adder with 2-bit inputs No B[1:0] signals
carrySkip2NoBCin1	A[1:0]	S[1:0] Cout	Carry skip adder with 2-bit inputs Cin=1 No B[1:0] signals
carrySkip2NoBNoCout	A[1:0] Cin	S[1:0]	Carry skip adder with 2-bit inputs No B[1:0] signals No Cout signals
carrySkip4	A[3:0] B[3:0] Cin	S[3:0] Cout	Cascading 2 carrySkip2s
carrySkip4NoCin	A[3:0] B[3:0]	S[3:0] Cout	Cascading a carrySkip2NoCin and a carrySkip2
carrySkip8NoCin	A[7:0] B[7:0]	S[4:0] Cout	Cascading a carrySkip2NoCin and 3 carrySkip2s
addOne	in[7:0]	out[7:0]	Input a 8-bit number and add 1. Do not need to consider the case of overflowing since it won't happen.
a4bitsSelectorEnhanced	CTRL[2:0] nCTRL[2:0]	out[3:0]	Simplify the original 8 to 1 mux bus (a4bitsSelector) into simple logic gates
b10bitsSelectorEnhanced	CTRL[2:0] nCTRL[2:0]	out[9:0]	Simplify the original 8 to 1 mux bus (b10bitsSelector) into simple logic gates
mulStage1	aValue[5:0] abs_x[3:0]	add01[5:0] add23[5:0]	The first stage of multiplication, which is two additions. add01[5:0] : 2^{-6} to 2^{-11} add01[5:0] : 2^{-4} to 2^{-9}
mulStage2	add01[5:0] add23[5:0]	mul[7:0]	The second stage of multiplication, which is a addition. mul[7:0] : 2^{-4} to 2^{-11}
addStage	mul[7:0] bValue[9:0]	S[9:0]	The stage of addition S[9:0] : 2^{-2} to 2^{-11}

5. Slope and Intercept

x	a (2^{-3} to 2^{-6})	b (2^{-2} to 2^{-11})
[0,0.5)	1111	00_0000_0011
[0.5,1)	1110	00_1111_1101
[1,1.5)	1011	01_1101_1110
[1.5,2)	1000	10_1000_1111
[2,2.5)	0110	11_0000_1011
[2.5,3)	0100	11_0110_0100
[3,3.5)	0010	11_1010_0010
[3.5,4)	0001	11_1100_1000

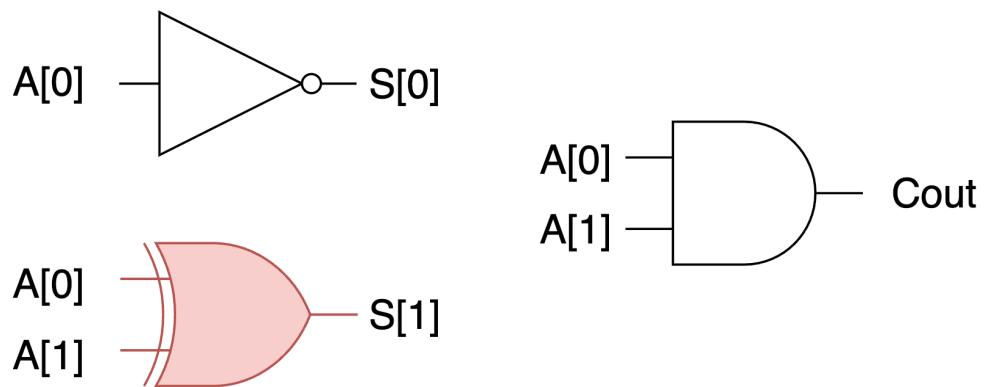
6. Self-defined Sub Module Circuit Diagrams

Sub Module Name	Circuit Diagram (Red Parts Are Critical Paths)
carrySkip2	<p>Simplified from CHAP8 slide</p>
carrySkip2NoCin	<p>Simplified from carrySkip2</p>



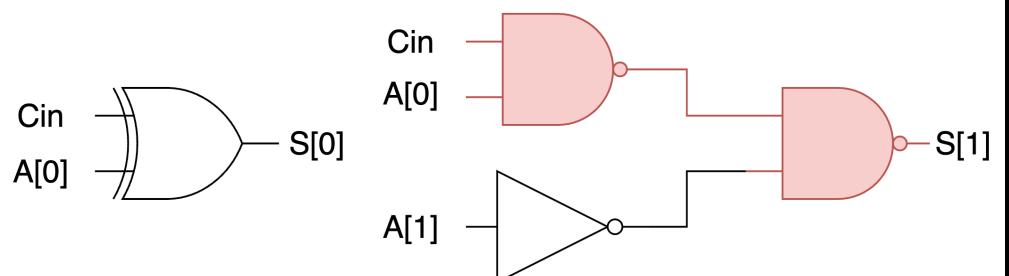
$\{\text{Cout}, S[1:0]\}$	$\text{Cin}=0$	$\text{Cin}=1$	$S[0]$	$\text{Cin} \wedge A[0]$
$A[1:0]=00$	000	001		
$A[1:0]=01$	001	010		
$A[1:0]=11$	011	100	$S[1]$	$(\text{Cin} \wedge A[0]) \wedge A[1]$
$A[1:0]=10$	010	011	Cout	$\text{Cin} \wedge A[0] \wedge A[1]$

carrySkip2NoBCin1



$\{\text{Cout}, S[1:0]\}$	$A[1]=0$	$A[1]=1$	$S[0]$	$A[0]'$
$A[0]=0$	001	011	$S[1]$	$A[0] \wedge A[1]$
$A[0]=1$	010	100	Cout	$A[0] \wedge A[1]$

carrySkip2NoBNoCout



$S[1:0]$	$\text{Cin}=0$	$\text{Cin}=1$	$S[0]$	$\text{Cin} \wedge A[0]$
$A[1:0]=00$	00	01		
$A[1:0]=01$	01	10		
$A[1:0]=11$	11	xx	$S[1]$	$\text{Cin} \wedge A[0] + A[1]$
$A[1:0]=10$	10	11		



carrySkip4



Cascading 2 carrySkip2s

carrySkip4NoCin



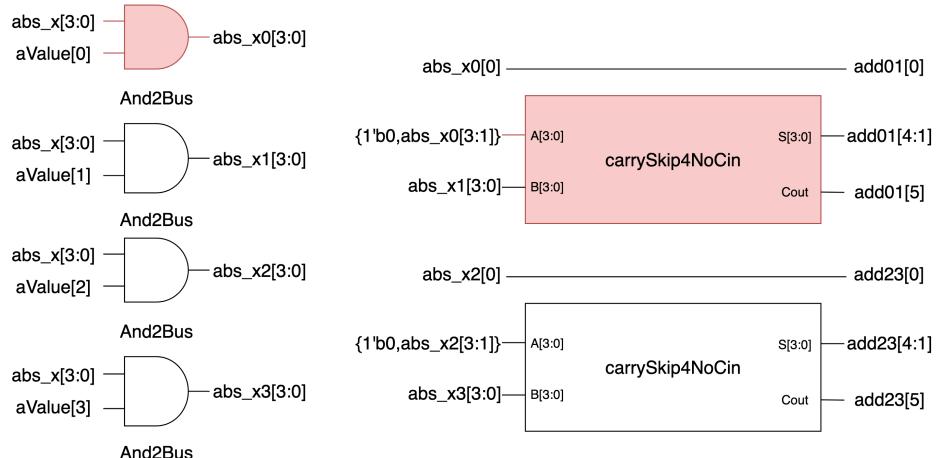
Cascading a carrySkip2NoCin and a carrySkip2

	<p>carrySkip4NoCin</p>								
carrySkip8NoCin	<p>carrySkip8NoCin</p> <p>carrySkip4</p>								
	<p>Cascading a carrySkip4NoCin and a carrySkip4</p>								
addOne	<p>addOne</p> <p>carrySkip2NoBCin1</p> <p>carrySkip2NoB</p> <p>carrySkip2NoB</p> <p>carrySkip2NoBNoCout</p>								
	<p>Cascading a carrySkip2NoBCin1, two carrySkip2NoBs, and a carrySkip2NoBNoCout</p>								
a4bitsSelectorEnhanced	<table border="1"> <tbody> <tr> <td>out[0]</td><td>$\text{CTRL}[2]' \text{CTRL}[0]' + \text{CTRL}[2]\text{CTRL}[1]\text{CTRL}[0]$ This is the critical path.</td></tr> <tr> <td>out[1]</td><td>$\text{CTRL}[2]' \text{CTRL}[1]' + \text{CTRL}[0]' \text{CTRL}[1]$</td></tr> <tr> <td>out[2]</td><td>$\text{CTRL}[1]$</td></tr> <tr> <td>out[3]</td><td>$\text{CTRL}[2]$</td></tr> </tbody> </table> <p>The left diagram is the original 8 to 1 mux bus. The right table is the simplified logics. Use nand gates to realize.</p>	out[0]	$\text{CTRL}[2]' \text{CTRL}[0]' + \text{CTRL}[2]\text{CTRL}[1]\text{CTRL}[0]$ This is the critical path.	out[1]	$\text{CTRL}[2]' \text{CTRL}[1]' + \text{CTRL}[0]' \text{CTRL}[1]$	out[2]	$\text{CTRL}[1]$	out[3]	$\text{CTRL}[2]$
out[0]	$\text{CTRL}[2]' \text{CTRL}[0]' + \text{CTRL}[2]\text{CTRL}[1]\text{CTRL}[0]$ This is the critical path.								
out[1]	$\text{CTRL}[2]' \text{CTRL}[1]' + \text{CTRL}[0]' \text{CTRL}[1]$								
out[2]	$\text{CTRL}[1]$								
out[3]	$\text{CTRL}[2]$								

b10bitsSelectorEnhanced		out[0]	$CTRL[1]'CTRL[0]' + CTRL[2]'CTRL[0]$
		out[1]	$CTRL[0]' + CTRL[2]'CTRL[1]$
		out[2]	$CTRL[1]'CTRL[0] + CTRL[2]'CTRL[1]$
		out[3]	$CTRL[2]'CTRL[1] + CTRL[2]'CTRL[0] + CTRL[1]'CTRL[0] + CTRL[2]CTRL[1]'CTRL[0]$ This is the critical path.
		out[4]	$CTRL[2]'CTRL[1]'CTRL[0] + CTRL[2]'CTRL[1]CTRL[0]$
		out[5]	$CTRL[1]'CTRL[0] + CTRL[2]CTRL[1]CTRL[0]$
		out[6]	$CTRL[1]'CTRL[0] + CTRL[2]CTRL[0] + CTRL[2]'CTRL[1]CTRL[0]$
		out[7]	$CTRL[1] + CTRL[2]'CTRL[0]$
		out[8]	$CTRL[2] + CTRL[1]'CTRL[0]$
		out[9]	$CTRL[2] + CTRL[1]'CTRL[0]$

The left diagram is the original 8 to 1 mux bus.
The right table is the simplified logics. Use nand gates to realize.

mulStage1



Go through an And2Bus to determine whether to add or not.

There are two additions.

$$\{1'b0, abs_x0[3:1]\} + abs_x1[3:0]$$

$$\{1'b0, abs_x2[3:1]\} + abs_x3[3:0]$$

Thus we need 2 carrySkip4NoCins.

	<p>add01[1:0]</p> <p>add01[5:2]</p> <p>add23[3:0]</p> <p>mul[1:0]</p> <p>mul[5:2]</p> <p>Cout</p> <p>S[3:0]</p> <p>carrySkip4NoCin</p> <p>A[3:0]</p> <p>B[3:0]</p> <p>Cin</p> <p>carrySkip2NoBNoCout</p> <p>s[1:0]</p> <p>out[7:6]</p> <p>add23[5:4]</p> <p>A[1:0]</p>
mulStage2	<p>There is only one addition $\{2'b00, \text{add01}[5:0]\} + \{\text{add23}[5:0], 2'b00\}$</p> <p>But we can split into two</p> <p>$\text{add01}[5:2] + \text{add23}[3:0]$</p> <p>$\text{add23}[5:4] + \text{carry}$</p>
addStage	<p>mul[7:0]</p> <p>bValue[7:0]</p> <p>S[7:0]</p> <p>carrySkip8NoCin</p> <p>A[7:0]</p> <p>B[7:0]</p> <p>Cout</p> <p>S[7:0]</p> <p>Cin</p> <p>carrySkip2NoBNoCout</p> <p>s[1:0]</p> <p>S[9:8]</p> <p>bValue[9:8]</p> <p>A[1:0]</p>

7. Discussion

1. Result

```
=====
Simulation finished
=====

=====
Summary
=====

Clock cycle: 3.1 ns
Number of transistors: 3916
Total execution cycle: 261
Approximation Error Score: 1500.0
Performance Score: 3168435.6
=====
```

2. How to determine the splitting points ?

For convenience, splitting the non-negative sides into 2,4,8,16... equivalent intervals are easier. Since the error of 4 intervals is too large and splitting into 16 intervals costs too many transistors, splitting into 8 intervals is the most efficient and economical.

3. How to determine the number of bits of the slope and the intercept of every interval ?

Since we need to do multiplication on the value of the slope, it is obviously to choose less bits. After plotting the sigmoid function on GeoGebra, I found out that the slope of every interval is less than 0.25. Then I started testing the number of bits and it turned out that choosing 4 bits, 2^3 to 2^6 , is the most economical solution. Thus, we can also determine the intercept should be 2^0 to 2^{-11} . While the first two bits of the intercept will always be 01, so choosing 10 bits is enough. More bits will only cost more resources but increase little precision.

4. How to determine the value of the slope and the intercept of every interval ?

After deciding the manner of splitting and the number of bits, we can now determine the value of the slope and the intercept of every interval. I first used GeoGebra to approximate some values. The second step is to write some codes to determine them. I've wrote JavaScript code to calculate the error.

Below are some parts of my code.

```
3 const aComponent = [
4   [3, 4, 5, 6], // 0 0.5
5   [3, 4, 5], // 0.5 1
6   [3, 5, 6], // 1 1.5
7   [3], // 1.5 2
8   [4, 5], // 2 2.5
9   [4], // 2.5 3
10  [5], // 3 3.5
11  [6] // 3.5 4
12 ];
13
14 const bComponent = [
15   [1, 9, 10], // 0 0.5
16   [1, 7, 8, 9, 11], // 0.5 1
17   [1, 5, 6, 7, 8, 9, 10], // 1 1.5
18   [1, 3, 8, 9, 10, 11], // 1.5 2
19   [1, 3, 4, 8, 10, 11], // 2 2.5
20   [1, 2, 6, 9], // 2.5 3
21   [1, 2, 4, 5, 6, 10], // 3 3.5
22   [1, 2, 3, 5, 7, 8] // 3.5 4
23 ];
24
25 const a = aComponent.map(each => {
26   // combine the components to get the value
27   return each.reduce((acc, current) => acc + Math.pow(2, -current), 0)
28 });
29
30 const b = bComponent.map(each => {
31   // combine the components to get the value
32   return each.reduce((acc, current) => acc + Math.pow(2, -current), 0)
33 });

64 // Convert binary to decimal and write to output file
65 const outputPath = '../jsResults/outPractice5.dat';
66 var mse = 0;
67 const outputData = inputData.map((line, index) => {
68   const decimalX = binaryToDecimal(line.trim()) / 32;
69   const absX = Math.abs(decimalX)
70   const decimalY = goDataDec[index]
71   let approximateY = a[Math.floor(absX * 2)] * absX + b[Math.floor(absX * 2)];
72   if (absX == 4) approximateY = a[7] * absX + b[7];
73   if (line.trim()[0] === '!') approximateY = 1 - Math.pow(2, -11) - approximateY + Math.pow(2, -12) + Math.pow(2, -14) + Math.pow(2, -15);
74   else approximateY += Math.pow(2, -14) + Math.pow(2, -15);
75   const error = approximateY - decimalY;
76   mse += error * error;
77   // console.log(`Math.floor(decimalX * 2) + 8], b[Math.floor(decimalX * 2) + 8])`)
78   return (`InputDec : ${decimalX}toString()
79   + `n`inputBin : ${line}
80   + `n`goldenDec : ${goDataBin[index]}
81   + `n`goldenBin : ${decimalY}toString()
82   + `n`approximateDec : ${approximateY}toString()
83   + `n`approximateBin : ${padBinaryTo16Bits((approximateY * Math.pow(2, 15)).toString(2))
84   + `n`approximateHex : ${(approximateY * Math.pow(2, 15)).toString(16)}
85   + `n`error : ${error}toString()
86   + `n`n`);
87 });
88 
```

5. How to choose adders ?

There are many kinds of adder we can choose. Carry-skip adder performs well while not costing too many transistors. We've learned a 4-bit based carry-skip adder in the class. After googling, I found out that using 2-bit based carry-skip-adder is better though it cost a bit more transistors. Since I at most need a 10-bit adder and most of them are 8-bit or 4-bit adders, using 2-bit based adders is wise. By wiki, when $n=8$, $m=\sqrt{8/2}=2$.

I used the 4-bit based carry-skip adder in my circuit. Although it saved some transistors, but it cost 0.4 ns more than the 2-bit based carry-skip adder.

Fixed size block-carry-skip adders [edit]

Fixed size block-carry-skip adders split the n bits of the input bits into blocks of m bits each, resulting in $k = \frac{n}{m}$ blocks. The critical path consists of the ripple path and the skip element of the first block, the skip paths that are enclosed between the first and the last block, and finally the ripple-path of the last block.

$$T_{FCSA}(n) = T_{CRA_{[0:n_{out}]}}(m) + T_{CSK} + (k-2) \cdot T_{CSK} + T_{CRA}(m) = 3D + m \cdot 2D + (k-1) \cdot 2D + (m+2)2D = (2m+k) \cdot 2D + 5D$$

The optimal block size for a given adder width n is derived by equating to 0

$$\begin{aligned} \frac{dT_{FCSA}(n)}{dm} &= 0 \\ 2D \cdot \left(2 - n \cdot \frac{1}{m^2}\right) &= 0 \\ \Rightarrow m_{1,2} &= \pm \sqrt{\frac{n}{2}} \end{aligned}$$

Only positive block sizes are realizable

$$\Rightarrow m = \sqrt{\frac{n}{2}}$$

Reference: https://en.wikipedia.org/wiki/Carry-skip_adder

6. How to lower transistors between different cycles ?

Encoding is important when sending signals between cycles. By encoding, we can save a lot of DFFs (DFFs cost a lot of transistors). Meanwhile, we only need to send useful signals. If the signal is not useful, just don't send it.

For example, I determined the intercept in stage 3 while not stage 1. Since the bandwidth of the intercept is 10 while the CTRL signals is only 3, we can instead send CTRL to stage 3 to conserve some DFFs. Lastly, just decode CTRL into intercept in stage 3.

7. How to accelerate multiplication ?

In our case, multiplication is 4-bit * 4-bit. We can divide the task into three additions : $[(A+B)+(C+D)]$, where A,B,C,D are the 4-bit numbers that multiplicand times each bit of the multiplier and shifts. First, do two additions parallelly : $(A+B)$ and $(C+D)$. Second, add two results together, which is handling the middle addition of $[(A+B)+(C+D)]$. Although it takes more transistors to realize, we can save much time on it and thus decrease the cycle time. Note that

A : 2^{-8} to 2^{-11} , B : 2^{-7} to 2^{-10} , C : 2^{-6} to 2^{-9} , D : 2^{-5} to 2^{-8}

$A+B$: 2^{-6} to 2^{-11} , $C+D$: 2^{-4} to 2^{-9} , final result : 2^{-4} to 2^{-11}

8. How to determine pipeline stages ?

There is no easy way to determine pipelining. I first did the work without pipeline, and then found out that I could not lower the cycle time much (9.5 ns). Then I started splitting tasks into different stages. After testing different number of stages, I found out that 6-stage case performs well. More stages only cost more transistors, but decrease little cycle time. Although the number of transistors increased, but I could further lower the cycle time into (3.1 ns).