

Computer Arithmetics: Floating Point Arithmetic Circuits



Dr. Ahmed Elhossini
Fachgebiet: Architektur eingebetteter Systeme
Institut für Technische Informatik und Mikroelektronik
Fak. IV – Elektrotechnik und Informatik

Previous Lecture

Floating Point Numbers

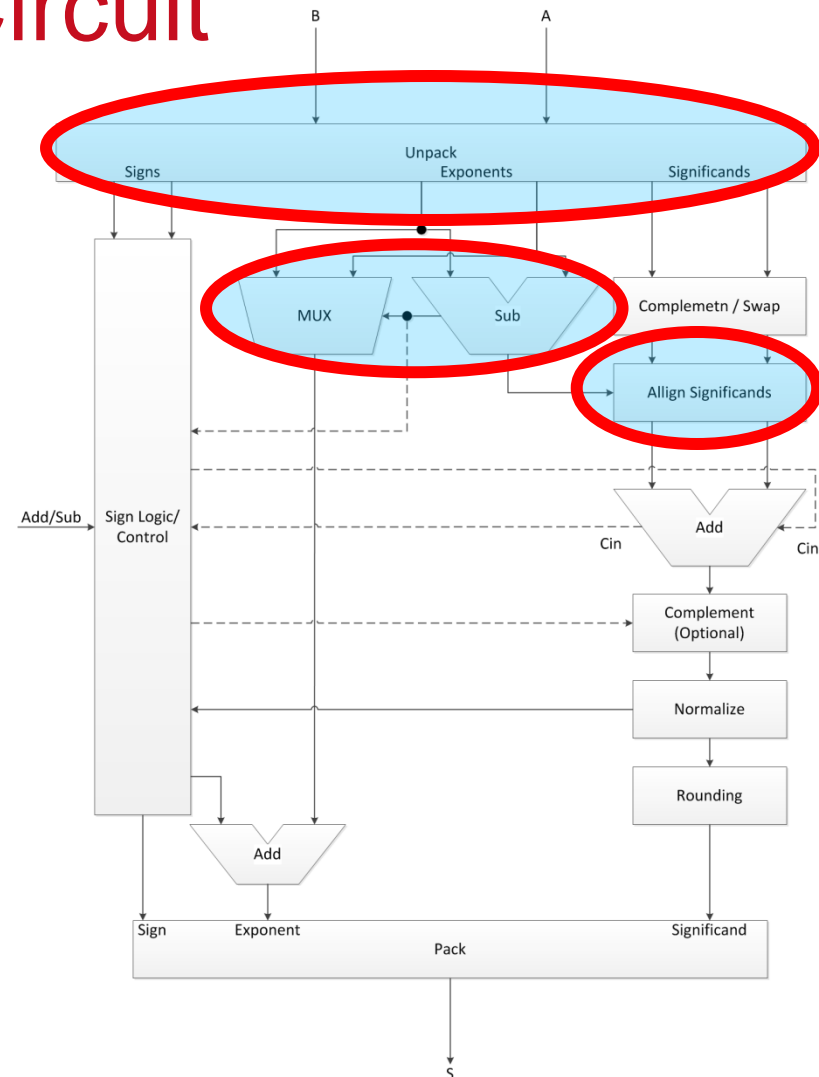
- Floating Point Numbers
 - Finite numbers
 - Floating point number format
 - IEEE floating point standard
 - Single / Double formats
 - Exponent Encoding
 - Special Operands and Subnormal's
- Arithmetic operations
 - Requirements
 - Addition/Subtraction
 - Multiplication/Division
 - Exception

Floating Point Circuits

- Adder Circuit
 - Controlled Shifter – Significant alignment
 - Normalization
 - Rounding
 - Sign Control and Exception Generation
- Multiplier Circuit
- Divider Circuit
- Multiply Accumulate Circuit

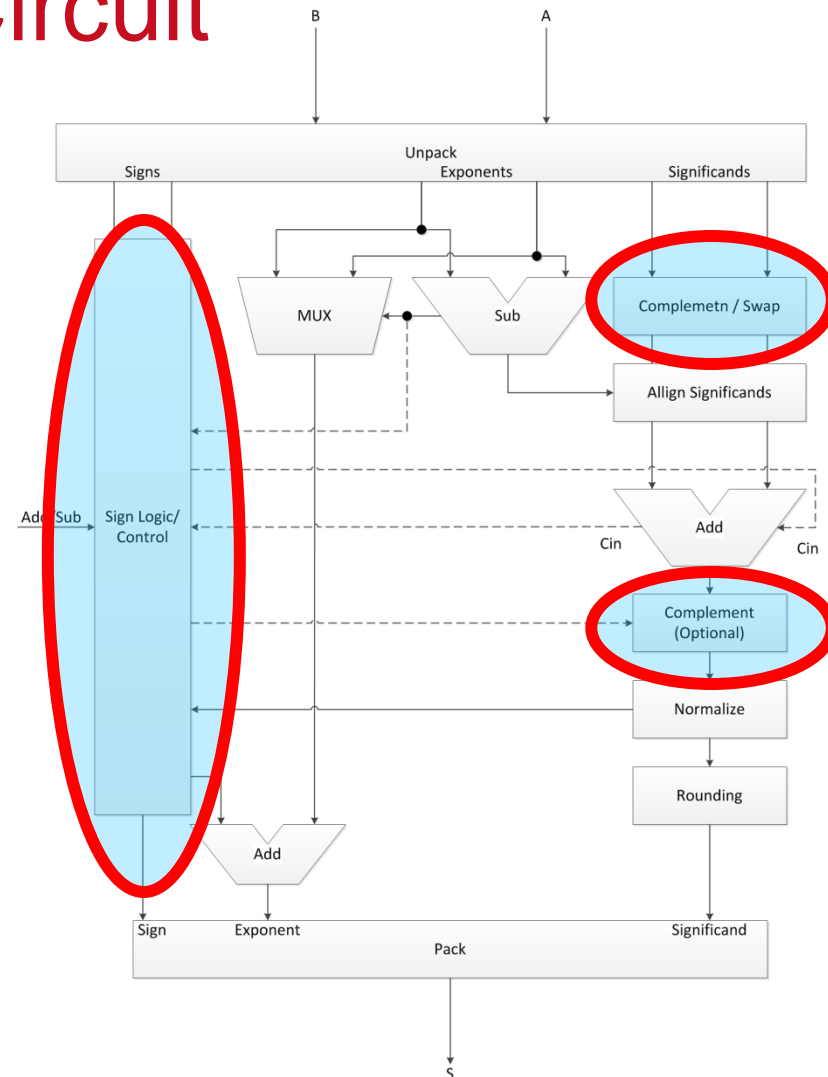
Floating Point Adder Circuit

- Separate the exponent and significant components from the inputs
- Subtract Exponents, depending on the result select the maximum value
- Use the difference between the two exponents to align the significant by shifting



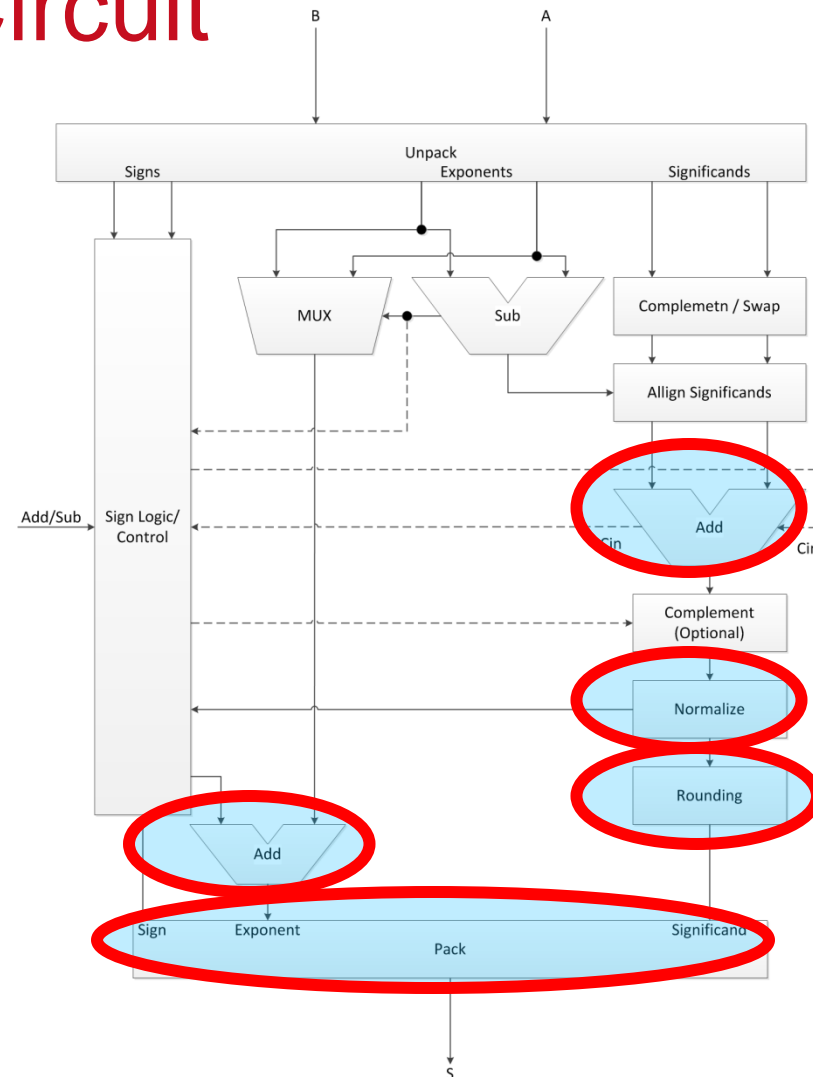
Floating Point Adder Circuit

- The sign logic is used to control the adder/subtractor
- Depending on the difference of the exponents, the significant can be swapped or complemented
- Depending on the sign of both operands, complementing the adder output may be required



Floating Point Adder Circuit

- Add (subtract significant)
- Normalize the addition result (after complement)
- Round the output to fit the significant
- Adjust the Exponent depending on the normalization result
- Pack the output parts into the output

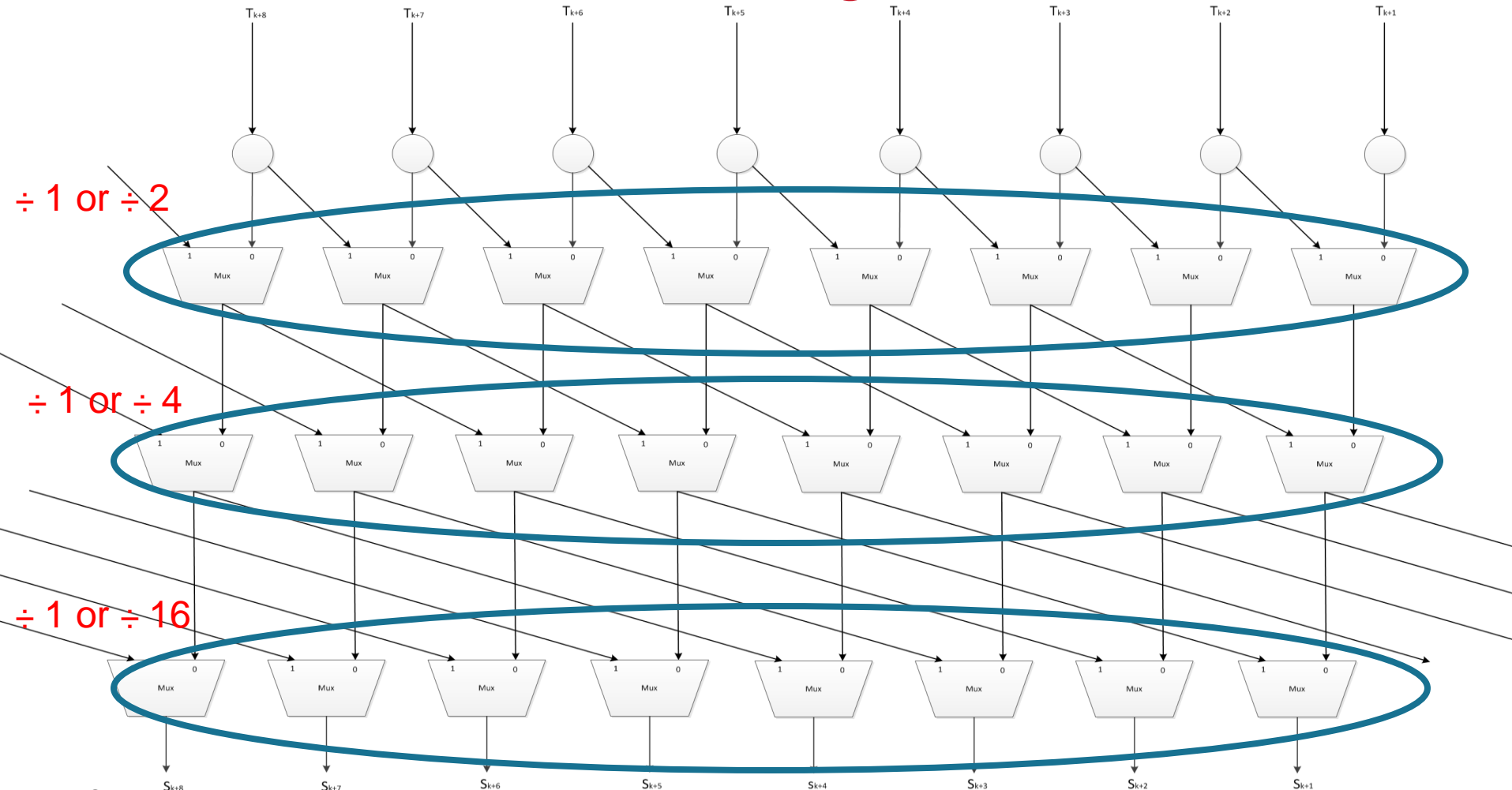




Significant Alignment – Barrel Shifter

- Significant Alignment is performed using Barrel Shifter
- Barrel shifter is able to perform controlled shifting (left/right) depending on a specific input
- Multiplexers are usually used to implement such a shifter
- Multiplexer stage is added for each control bit
- Total number of multiplexers needed = number of bits in the input vector * number of bits in the control word

Controlled Shifter – Right Barrel Shifter

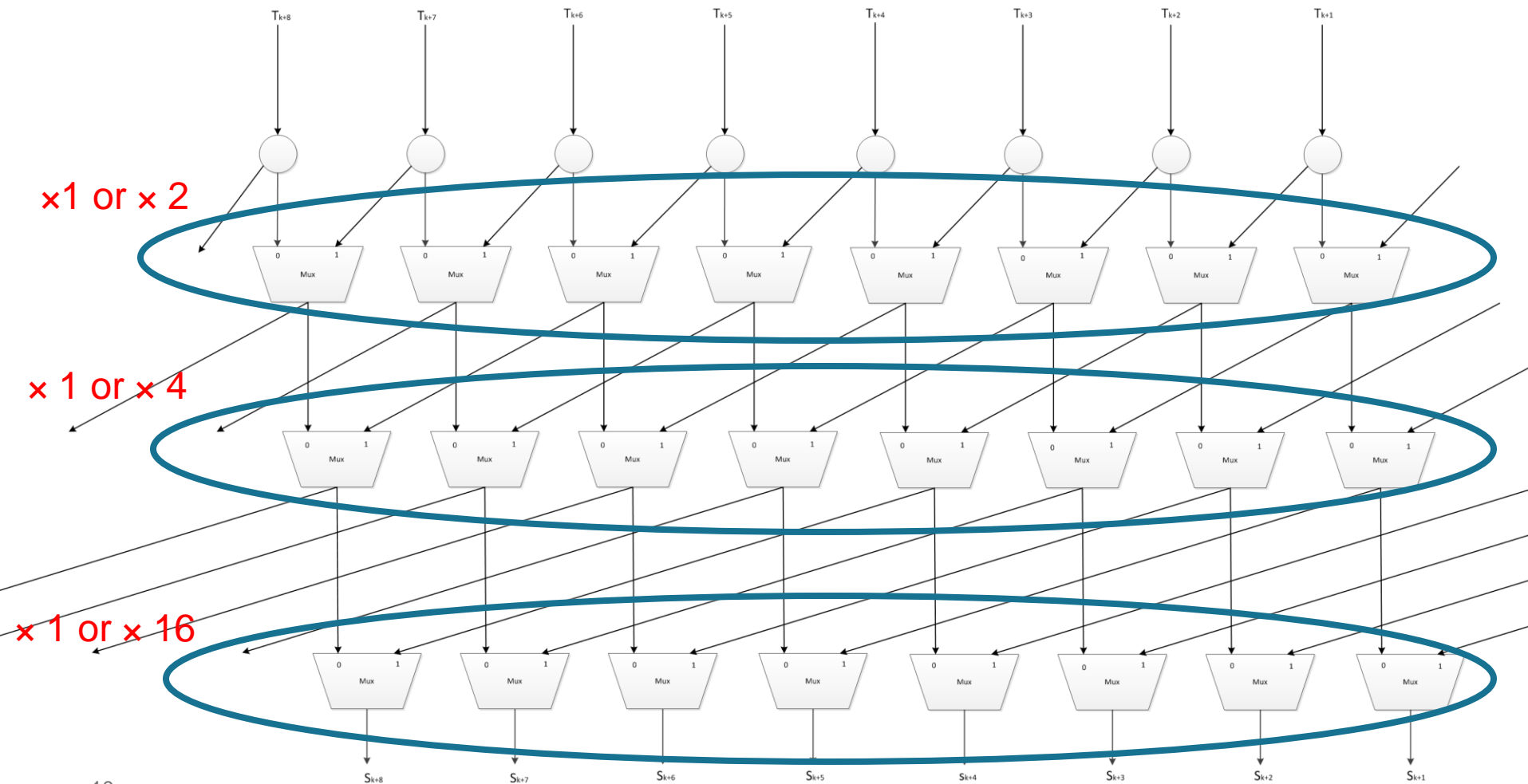


Controller Shifter – Right Barrel Shifter

S_2	S_1	S_0	Operation	Value	Shift Right Value
0	0	0	1 / 1×1×1	1 / 1	0
0	0	1	1 / 1×1×2	1 / 2	1
0	1	0	1 / 1×4×1	1 / 4	2
0	1	1	1 / 1×4×2	1 / 8	3
1	0	0	1 / 16×1×1	1 / 16	4
1	0	1	1 / 16×1×2	1 / 32	5
1	1	0	1 / 16×4×1	1 / 64	6
1	1	1	1 / 16×4×2	1 / 128	7

- Add more stages to allow more shifts
- Can be used in both significant alignment and normalization
- Similar structure can be used to make left shift

Controlled Shifter – Left Barrel Shifter



Controller Shifter – Left Barrel Shifter

S_2	S_1	S_0	Operation	Value	Shift Left Value
0	0	0	1×1×1	1	0
0	0	1	1×1×2	2	1
0	1	0	1×4×1	4	2
0	1	1	1×4×2	8	3
1	0	0	16×1×1	16	4
1	0	1	16×1×2	32	5
1	1	0	16×4×1	64	6
1	1	1	16×4×2	128	7

- Add more stages to allow more shifts
- Can be used in both significant alignment and normalization

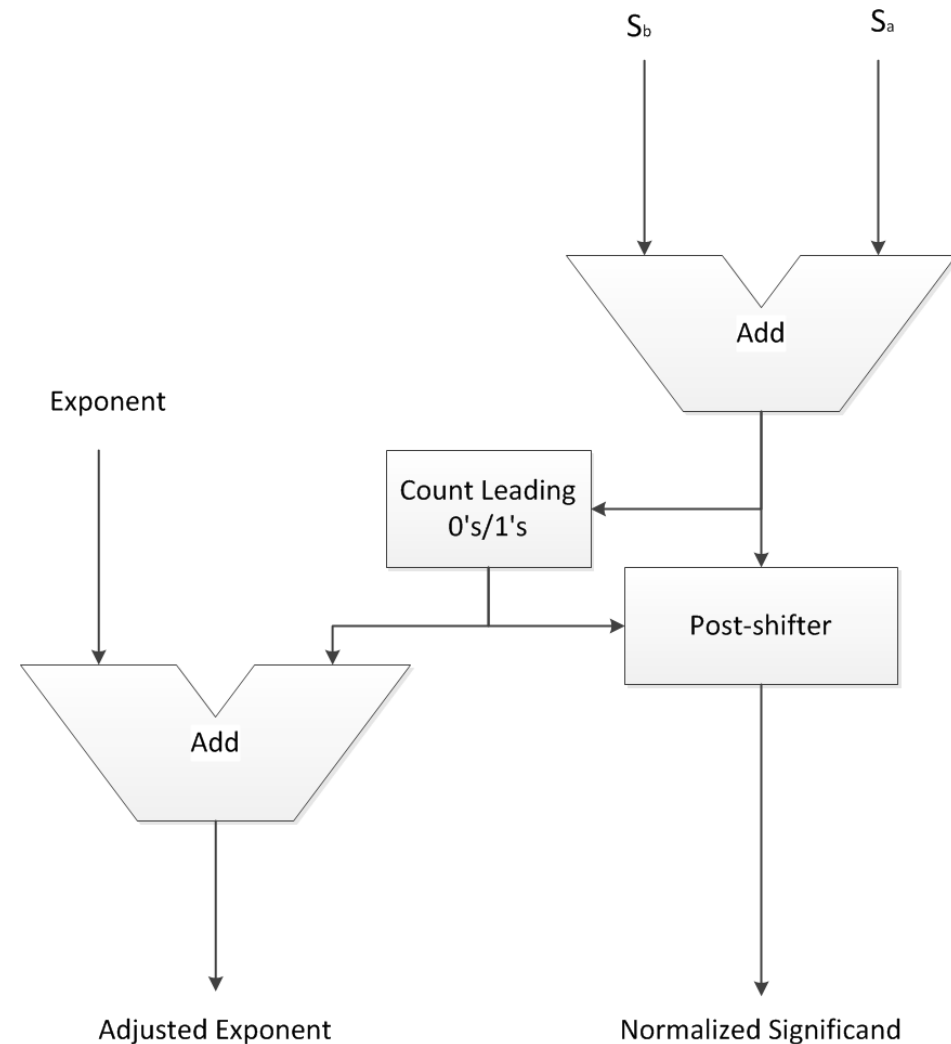


Normalization

- The normalization process is required ensure that significant value is in the range $[1,2)$
- It is done in two stage
 - Calculate the number of leading zeros/ones
 - Shift the significant left/right
- Leading 0's or 1's is calculated using a combinational logic
 - Leading 1's are found in negative two's complement numbers
- To speed it up, 0's or 1's can be predicted before addition so that Post shifting can be done just after addition

Normalization

- Counting leading 0's or 1's in the adder output
- Use this count as input to the barrel shifter
- Post-shifter and Exponent adjustment should wait for the counting process
- Simple implementation using combinational logic approach



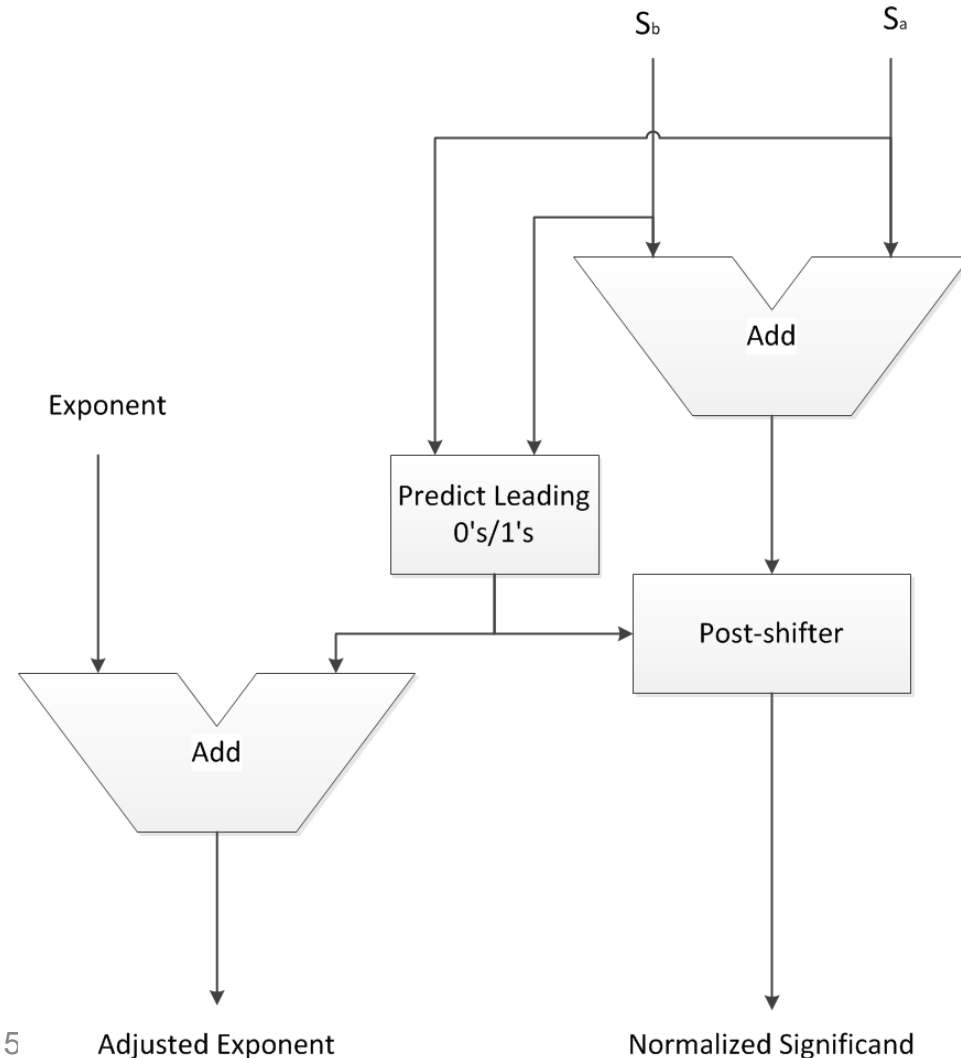
0's or 1's counting

a_0	a_1	a_2	a_3	a_4	a_5	a_6	a_7	s_2	s_1	s_0
1	x	x	x	x	x	x	x	0	0	0
0	1	x	x	x	x	x	x	0	0	1
0	0	1	x	x	x	x	x	0	1	0
0	0	0	1	x	x	x	x	0	1	1
0	0	0	0	1	x	x	x	1	0	0
0	0	0	0	0	1	x	x	1	0	1
0	0	0	0	0	0	1	x	1	1	0
0	0	0	0	0	0	0	1	1	1	1

- 8 x 3 priority encoder for 8 bits string
- For 24 bits (significant length in single word type), we need 24x5 encoder
- Can be implemented using a Look-up table

Normalization

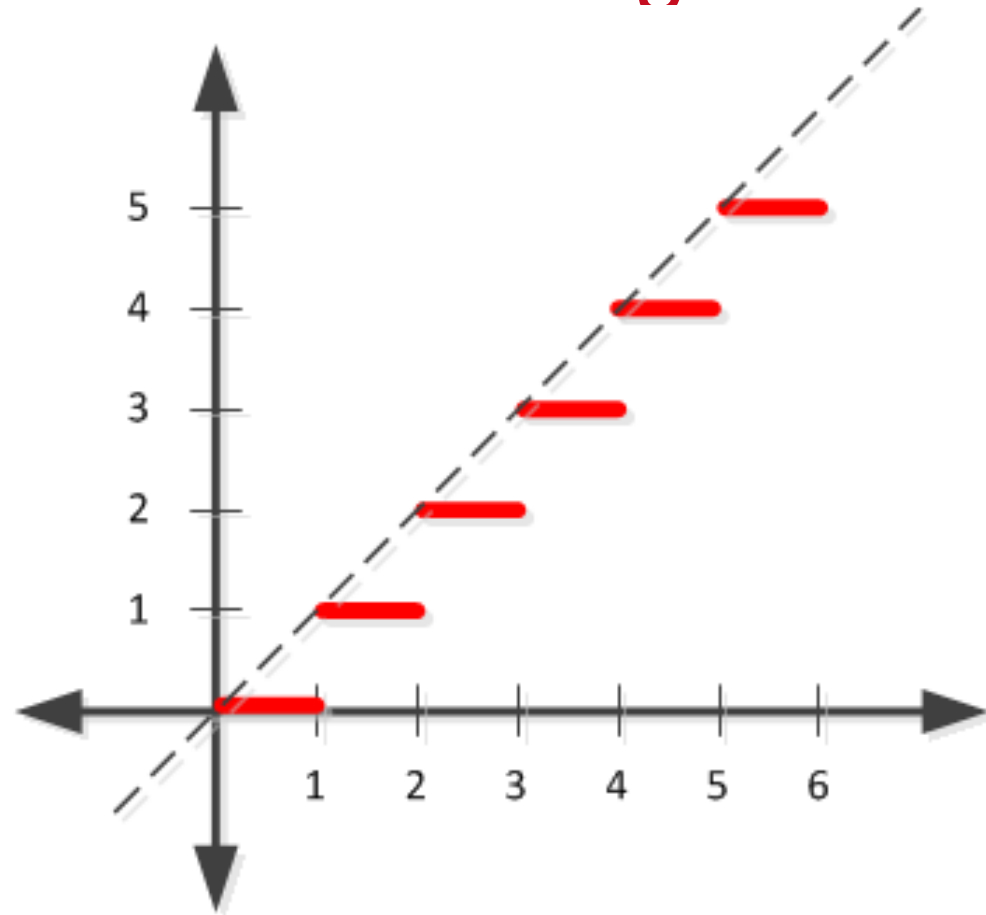
- Predict 0's or 1's before the addition operation
- Post-shifter can start directly after addition





Rounding: Truncation or Rounding

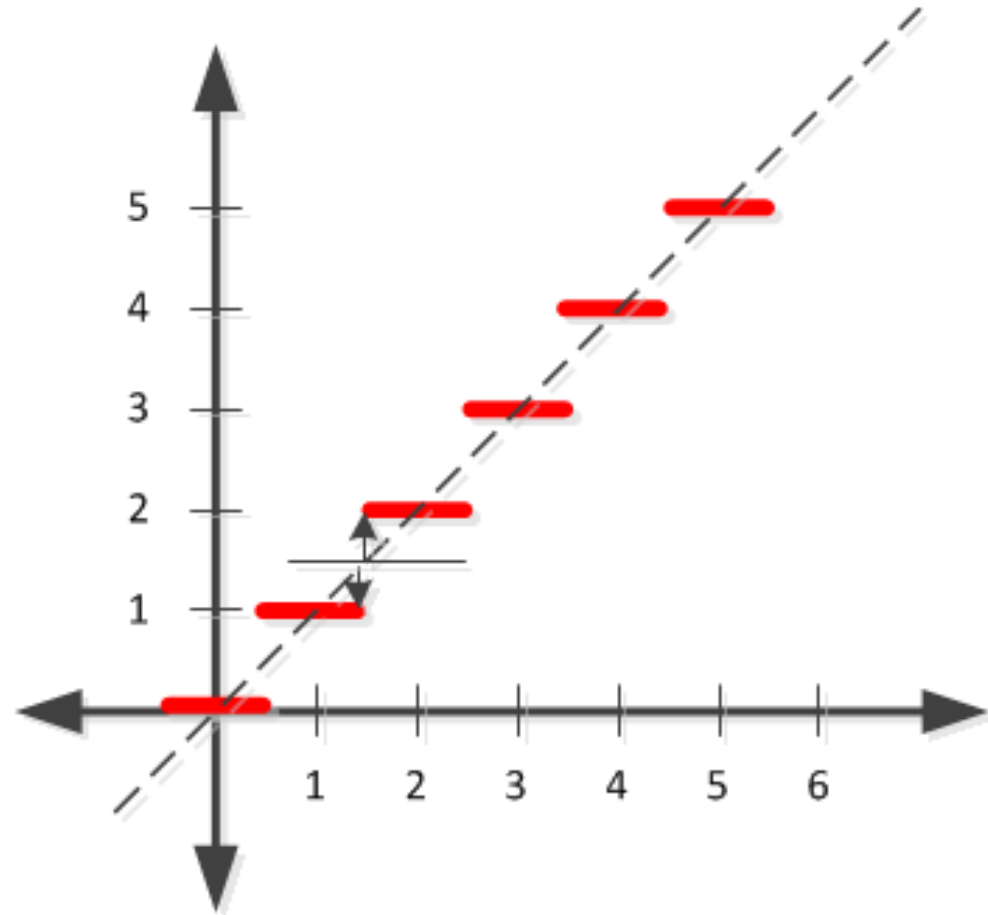
- Rounding or Truncation
- Truncation is simple,
 - remove extra bits.





Rounding

- Rounding can be up, or down to the nearest number



Rounding

- $X = x_2x_1x_0x_{-1}x_{-2}$
- $Y = y_2y_1y_0$
- If two bits to be rounded out:

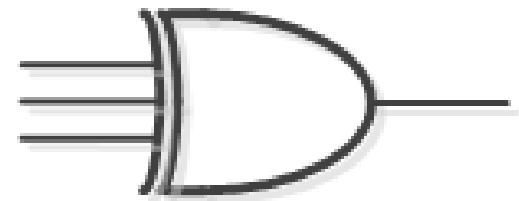
x_{-2}	x_{-1}	Y	Direction	Error
0	0	$Y = \text{int}(x)$	Down	0
0	1	$Y = \text{int}(x)$	Down	-0.25
1	0	$Y = \text{int}(x) + 1$	Up	0.5
1	1	$Y = \text{int}(x) + 1$	Up	0.25

- Combinational logic to encode the extra bits (lookup table can be used)
- Adder circuit is required to adjust the unrounded bits
- If this actual rounding is performed, another stage of normalization may be required

Sign Control

- Required to build adder/subtractor circuit
- Compare the sign of the operands to decide the function of the significant adder taking into account the required operation
- From the table it is an XOR function
- Pre-complement and Post complement are controlled by this signal

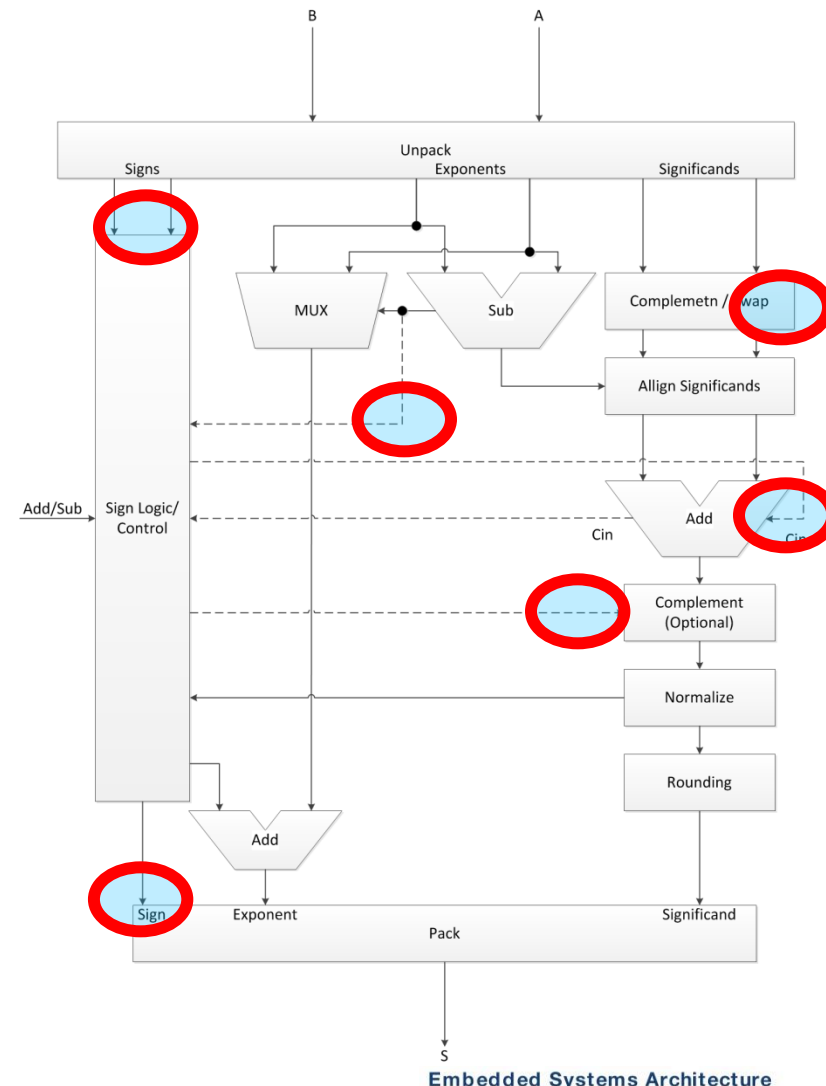
0 = add 1 = sub	sn_b	sn_a	adder operation
0	0	0	0
0	0	1	1
0	1	0	1
0	1	1	0
1	0	0	1
1	1	0	0
1	0	1	0
1	1	1	1

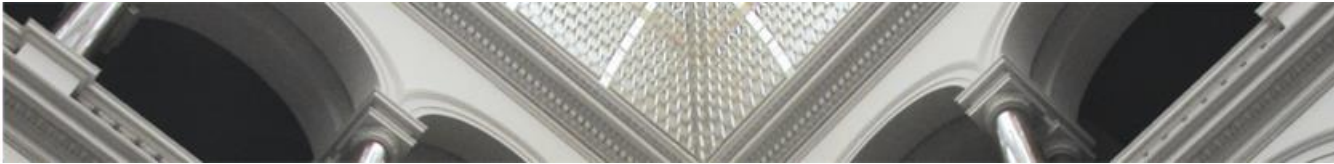


Sign Control

- The sign of the exponent difference may be used to swap the input
- The output sign depends on the operation being performed, and the sign of the adder result

Operation	Output Sign
Addition of two positive number	Positive
Addition of two negative numbers	Negative
Addition of different sign numbers	The sign of the adder output
Subtraction of two numbers with the same sign	The sign of the adder output
Subtraction of different sign numbers	The sign of the first operand





Exceptions in Adder Circuit

- Overflow/Underflow can be detected during exponent adjustment
 - Check the value of the exponent after normalization adjustment (the adder circuit)
- NaN values can be detected in the unpacking and packing stages
 - Check the values of the exponent and significant at the input and output
 - Same circuit as the zero detection circuit we used before
 - Same circuit can also be used to detect the Zero flag
- The Sign flag is the sign of the output



Final Notes on the Adder Circuit

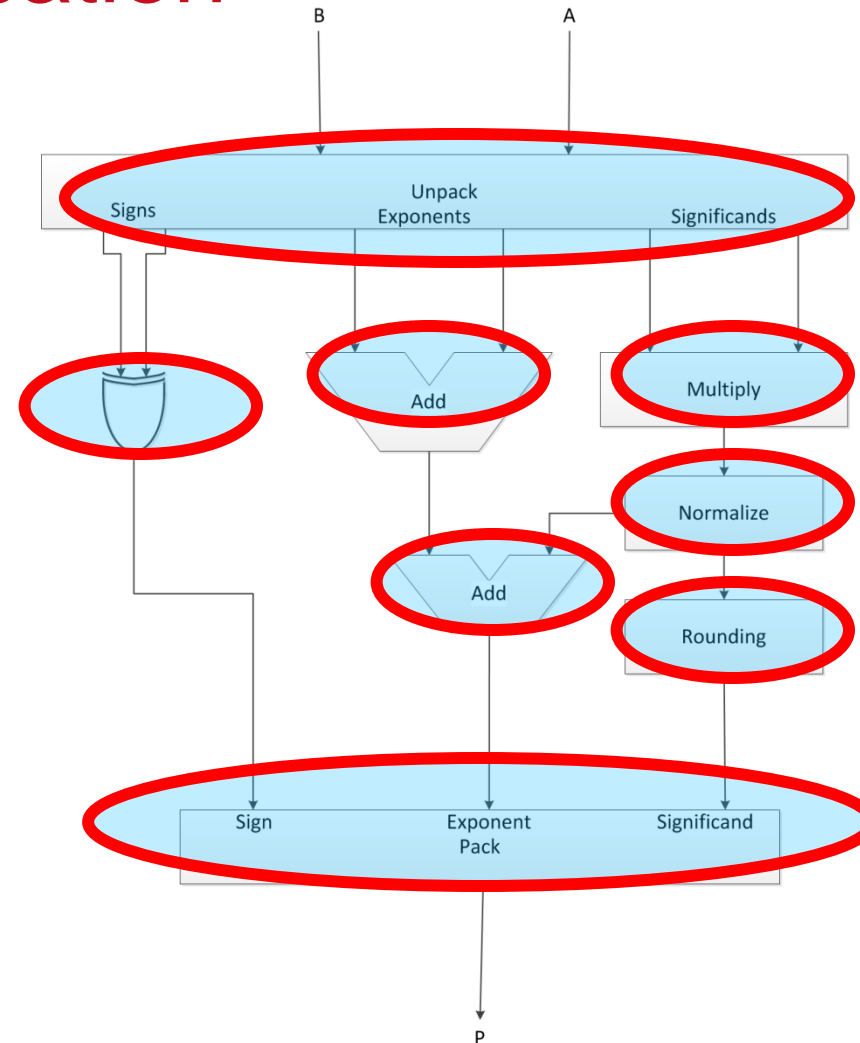
- Complex operation
- Several computational stages are required
- Pipelining is important here – 10 stages
 - Several Pipeline stages are required to minimize the critical path delay and increase the clock frequency
 - This is the reason that modern micro-processors are featured with deep pipelining
- Modern micro-processors have more than data path
 - A data path for integer operations which is faster and commonly used
 - A data path for floating point operation which is slower (in number of clock cycles) but less common to use
- Barrel shifting is a common instruction in many micro-processor
 - The existence of such instruction reduce the software size when shifting is required, which is common in many application
 - Resource sharing can be used to reduce the total processor area

Floating Point Multiplication

- Simple circuit compared to the adder
- Significant component are multiplied directly
- Exponents are added
- Normalization and Rounding are performed to generate the output
 - Barrel shifter is used as shown in the adder circuit
- Exceptions are checked the same way as discussed in the adder circuit

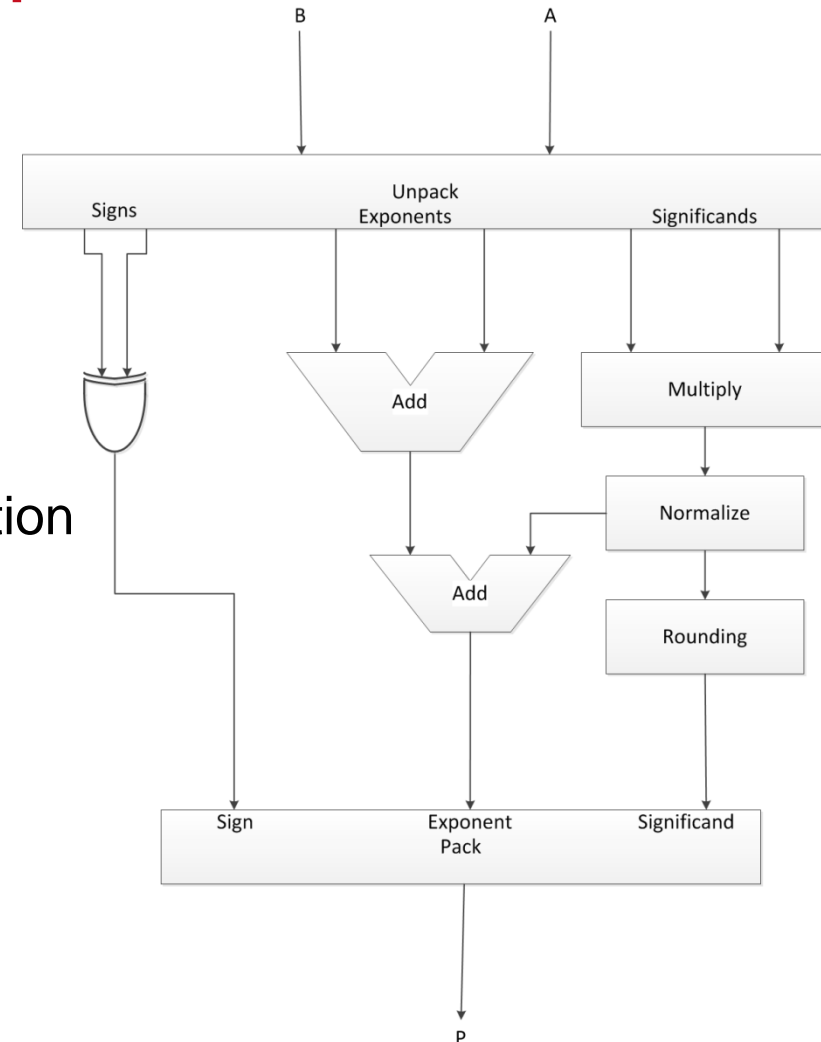
Floating Point Multiplication

- Unpack operand (check exceptions)
- Multiply Significant(s)
- Add the exponents
- Normalize Significant(s)
- Adjust Exponent
- Round the significant
 - Extra normalization may be needed
- Generate the sign
- Pack the output (check exceptions)



Floating Point Multiplier

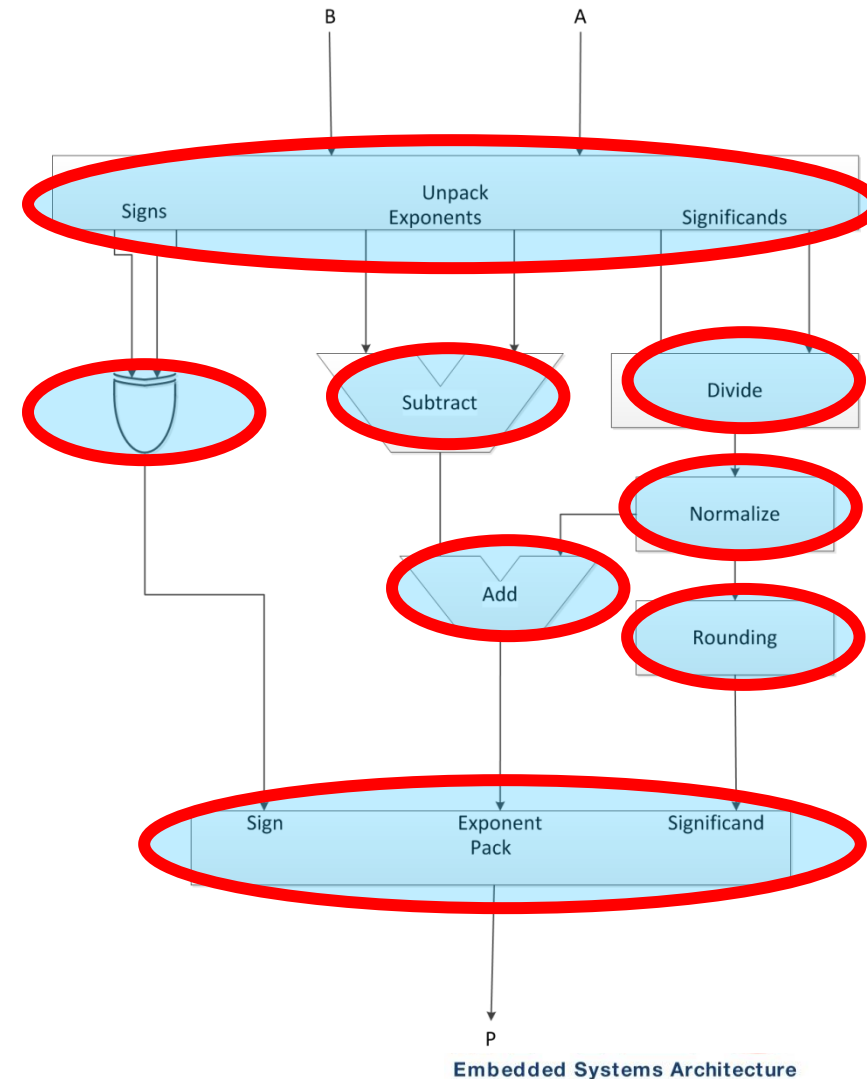
- Pipelining is also important
- 6 stages is required here
- Unsigned multiplier is required
- Barrel shifter is used in the normalization stage (refer to slides 12 – 15)





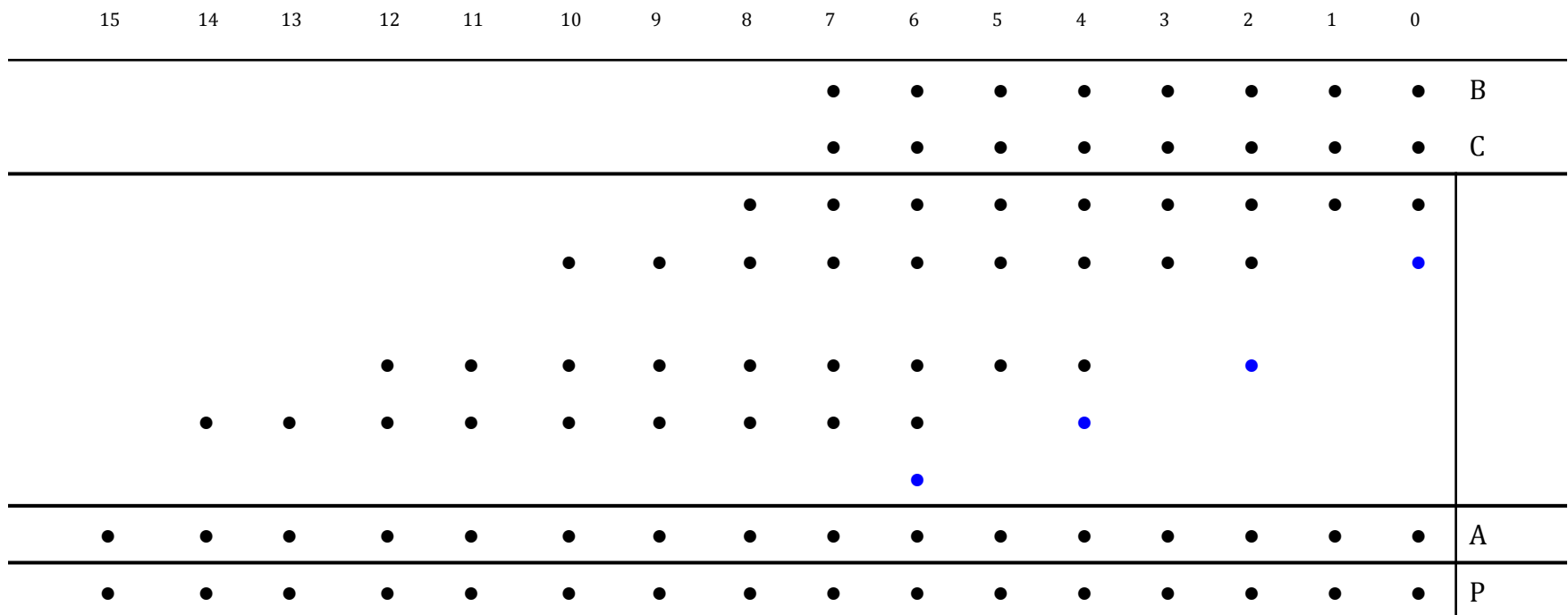
Floating Point Divider

- Multiplication is replaced with division
- Exponents are subtracted not added
- All other components remain the same



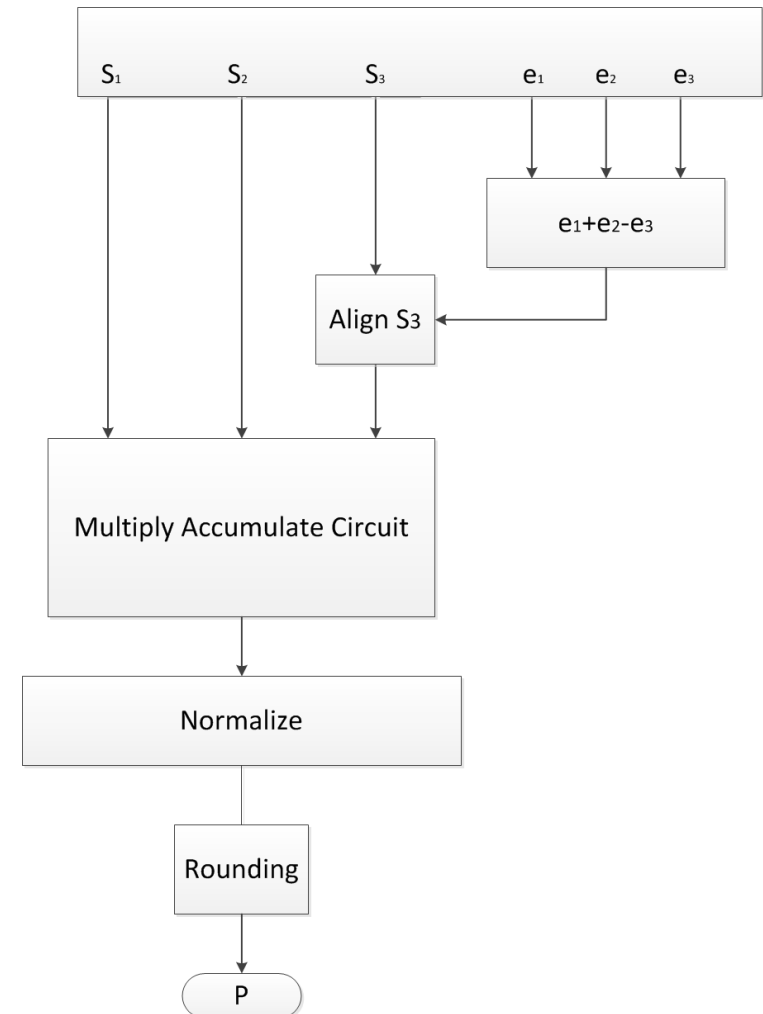
Multiply Accumulate

- Multiply accumulate operation is a basic operation in several applications
- Multiplication can be fused with addition in the same circuit



Multiply Accumulate Circuit

- Modified version of the adder circuit
- The adder is replaced by the fused MAC circuit
- Exponents are calculated from the maximum value of e_1+e_2 and e_3





Error Sources in Floating Point Arithmetic Circuits

- There are several sources in floating point arithmetic circuits
 - Representation Errors: Error due to the limited number of bits in the significant
 - Computation Errors: Error due to the rounding stage in all arithmetic circuit
- This error accumulate after several operations
- Increasing the significant part reduce representation error
- Intermediate representation during calculation can reduce final error
- The accumulated error violates the laws of algebra
 - $A + (B + C) \neq (A + B) + C$

Verify ?

Summary

- Adder Circuit
 - Controlled Shifter – Significant alignment
 - Normalization
 - Rounding
 - Sign Control and Exception Generation
- Multiplier Circuit
- Divider Circuit
- Multiply Accumulate Circuit
- Error Analysis



Questions ?

Questions ?