

# AI ASIC: Design and Practice (ADaP) Fall 2023

**Digital Arithmetic & Circuits** 

燕博南





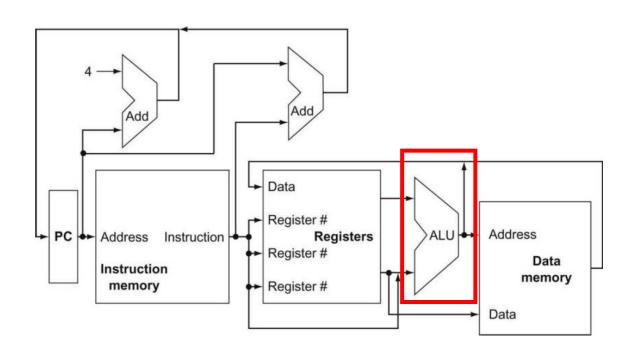
- Number Systems
  - Integer
  - Fixed-Point
  - Floating-Point
- Arithmetic
- Circuits & Implementation

Numbers Arithmetic Circuits



#### Why We Need to Introduce Arithmetic





- Arithmetic Logic Unit (ALU): heart of von Neumann architecture
- Deal with various precisions: decimals, fractions, integers, ...



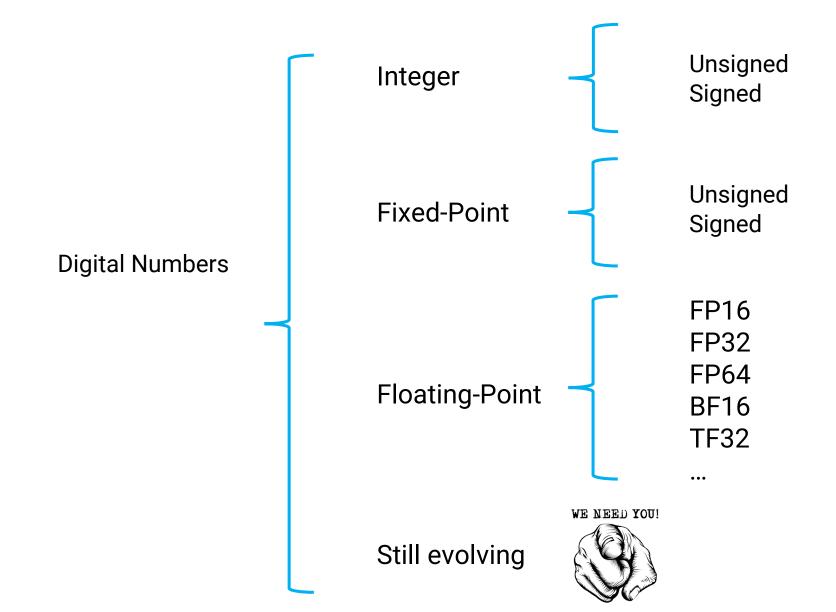
# Part 1

**Unsigned Integers** 



## **Number System of Digital Computers**





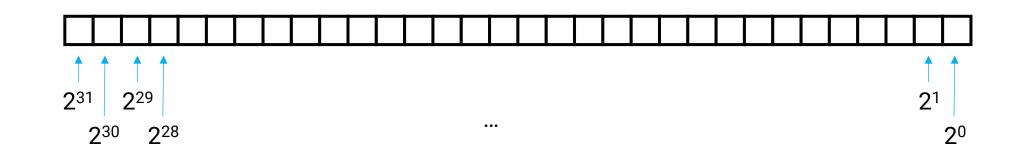


#### **Unsigned Integer**



Numbers Arithmetic Circuits

Unsigned INT32



Significance:

- INT16, INT8, ...
- Example:

32'd7 (=32'h0000\_0007) 8'b1100\_1101 (=8'hCD) Recommend tool:

programmer's calculator

### **Unsigned Integer Arithmetic - Add & Subtract**

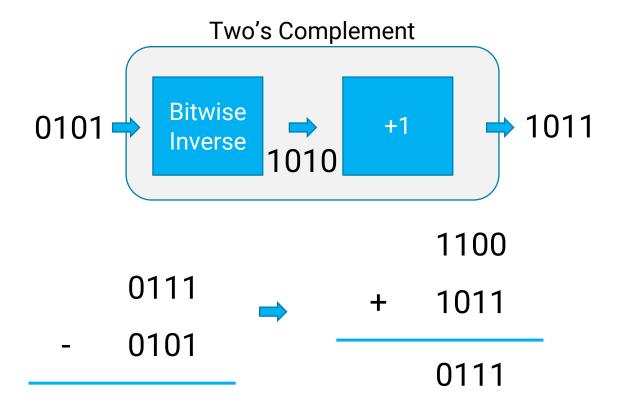


INT4 as example:

Add: 
$$4'd7 + 4'd5 = 4'd12$$



Subtract: 4'd12 - 4'd5 = 4'd7



### Unsigned Integer Arithmetic - Multiply & Divide



Numbers Arithmetic Circuits

• INT4 as example:

Multiply: 
$$4'd3 * 4'd5 = 4'd15$$



#### **Unsigned Integer Circuits - Adder**



**Numbers** 

> Arithmetic

Circuits

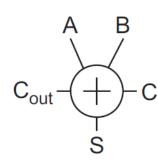
Half-Adder

$$C_{out}$$
 $A$ 
 $B$ 
 $C_{out}$ 

| Α | В | <b>C</b> out | S |
|---|---|--------------|---|
| 0 | 0 | 0            | 0 |
| 0 | 1 | 0            | 1 |
| 1 | 0 | 0            | 1 |
| 1 | 1 | 1            | 0 |

$$S = A \oplus B$$
$$C_{\text{out}} = A \cdot B$$

Full-Adder



| Α | В | С | G | P | K | <b>C</b> out | S |
|---|---|---|---|---|---|--------------|---|
| 0 | 0 | 0 | 0 | 0 | 1 | 0            | 0 |
| Ü |   | 1 | Ü |   |   | 0            | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 | 0            | 1 |
| Ü | 1 | 1 | O | 1 | O | 1            | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 | 0            | 1 |
| 1 |   | 1 |   |   |   | 1            | 0 |
| 1 | 1 | 0 | 1 | 0 | 0 | 1            | 0 |
| 1 | 1 | 1 | 1 | J | J | 1            | 1 |

$$S = A\overline{B}\overline{C} + \overline{A}B\overline{C} + \overline{A}\overline{B}C + ABC$$

$$= (A \oplus B) \oplus C = P \oplus C$$

$$C_{\text{out}} = AB + AC + BC$$

$$= AB + C(A + B)$$

$$= \overline{A}\overline{B} + \overline{C}(\overline{A} + \overline{B})$$

$$= MAJ(A, B, C)$$

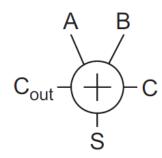
### **Unsigned Integer Circuits - Adder (Cont.)**



**Numbers** 

> Arithmetic

Circuits



$$S = A\overline{B}\overline{C} + \overline{A}B\overline{C} + \overline{A}\overline{B}C + ABC$$

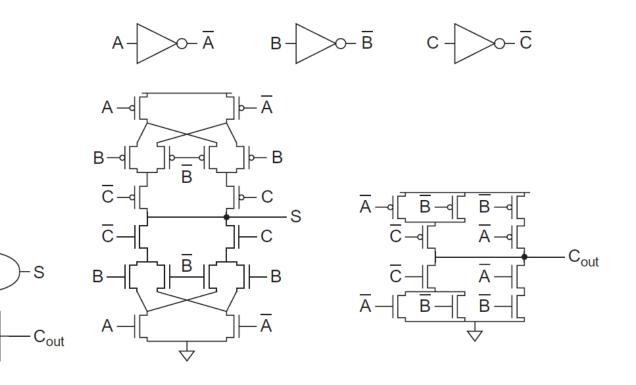
$$= (A \oplus B) \oplus C = P \oplus C$$

$$C_{\text{out}} = AB + AC + BC$$

$$= AB + C(A + B)$$

$$= \overline{AB} + \overline{C}(\overline{A} + \overline{B})$$

$$= MAJ(A, B, C)$$



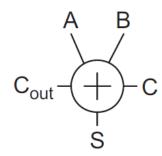
#### **Unsigned Integer Circuits - Adder (Cont.)**



Numbers

> Arithmetic

Circuits



$$S = A\overline{B}\overline{C} + \overline{A}B\overline{C} + \overline{A}\overline{B}C + ABC$$

$$= (A \oplus B) \oplus C = P \oplus C$$

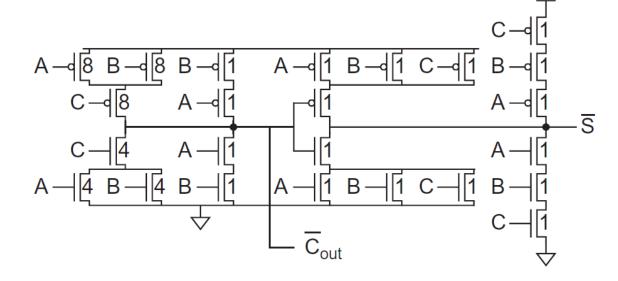
$$C_{\text{out}} = AB + AC + BC$$

$$= AB + C(A + B)$$

$$= \overline{AB} + \overline{C}(\overline{A} + \overline{B})$$

$$= MAJ(A, B, C)$$

#### Improved:



$$S = ABC + (A + B + C)\overline{C}_{out}$$

Idea behind:

Reuse Cout compute circuits to obtain both S



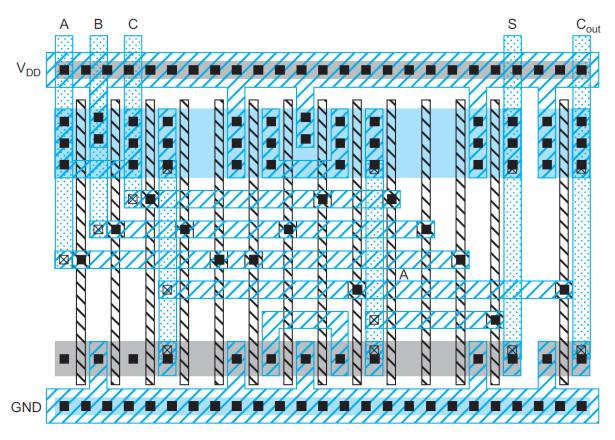
#### **Unsigned Integer Circuits - Adder (Cont.)**



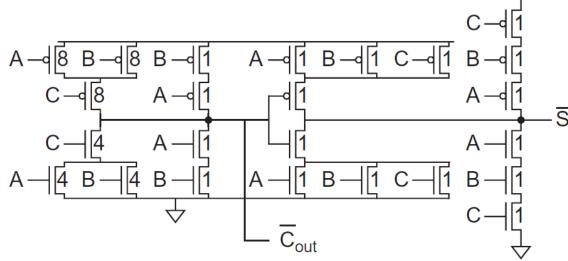
Numbers A

Arithmetic

Circuits



#### Improved:



$$S = ABC + (A + B + C)\overline{C}_{out}$$

Idea behind:

Reuse Cout compute circuits to obtain both S



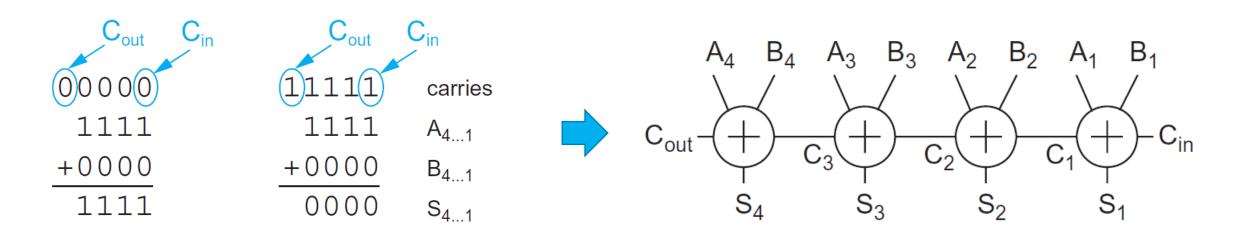
#### Adder Family - Carry-Ripple Adder



Numbers

> Arithmetic

Circuits

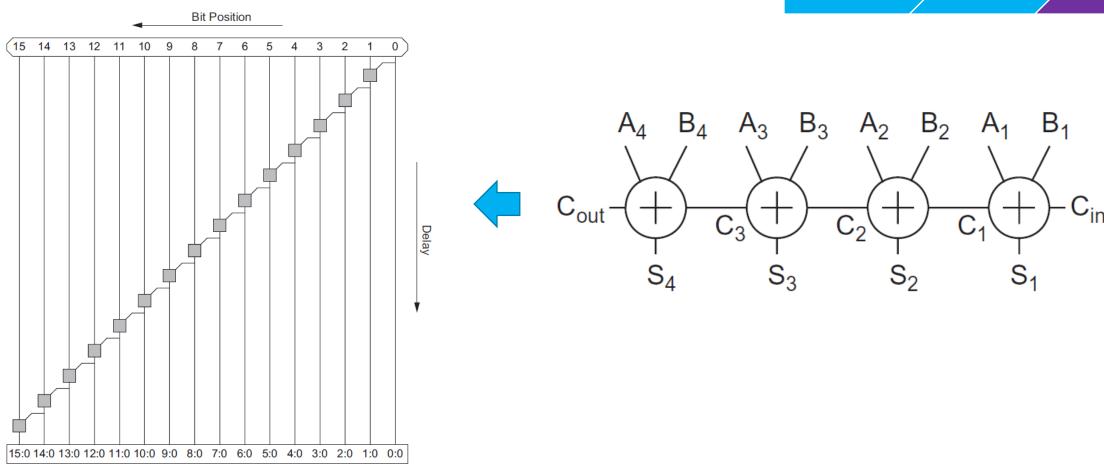


Natural & Intuitive However, carry propagation path too long

## Adder Family - Carry-Ripple Adder



Numbers Arithmetic Circuits



Weste, Neil HE, and David Harris. CMOS VLSI design: a circuits and systems perspective. Pearson Education India, 2015.



## Adder Family - Carry-Skip Adder

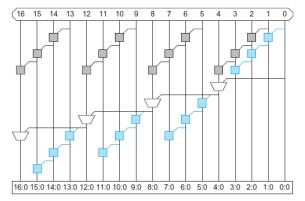


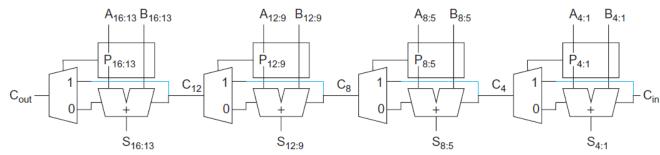
Numbers

> Arithmetic

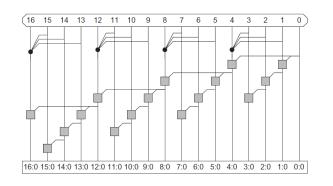
Circuits

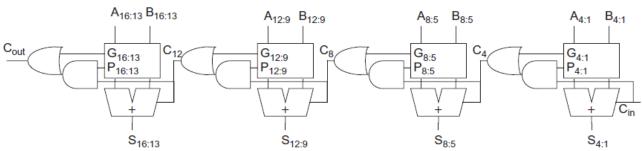
Carry-Skip Adder





Carry-Lookahead Adder





Idea Behind: Group and Divide!



### Adder Family - Big Family!

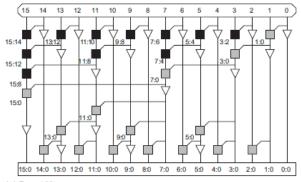
#### 和桌头等 PEKING UNIVERSITY

**Numbers** 

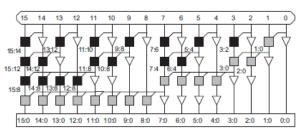
Arithmetic

Circuits

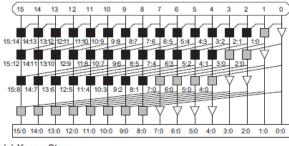
#### Tree Adder Family



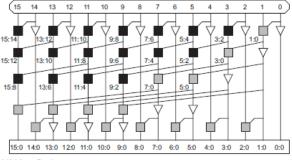
(a) Brent-Kung



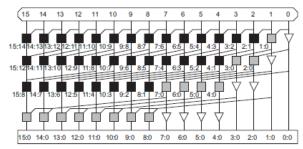
(b) Sklansky



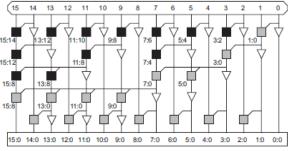
(c) Kogge-Stone



(d) Han-Carlson

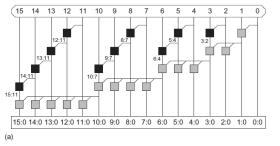


(e) Knowles [2,1,1,1]

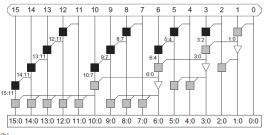


(f) Ladner-Fischer

15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0 13:12 14:12 15:0 14:0 13:0 12:0 11:0 10:0 9:0 8:0 7:0 6:0 5:0 4:0 3:0 2:0 1:0 0:0



**Sarry-Incremental Adder** 

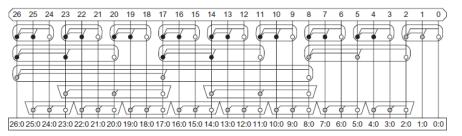


### Adder Family - Big Family! (Cont.)

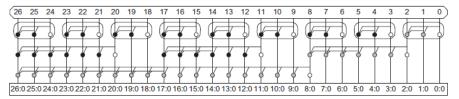


#### Sparse Tree Adders

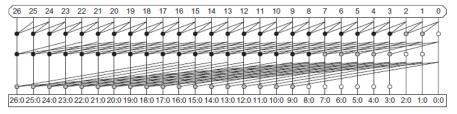
Numbers Arithmetic Circuits



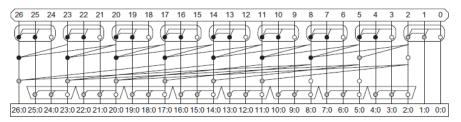
#### (a) Brent-Kung

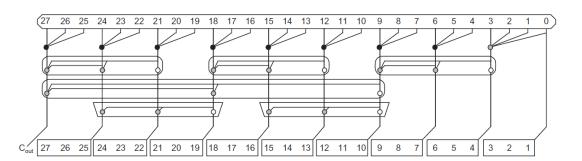


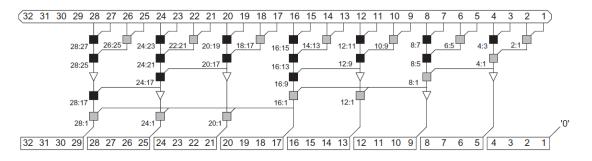
#### (b) Sklansky

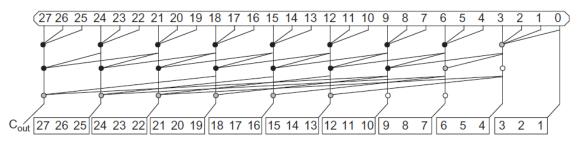


#### (c) Kogge-Stone





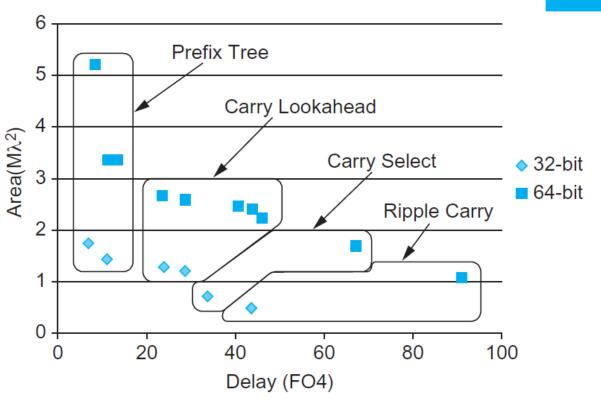




# Adder Family - Choose Wisely







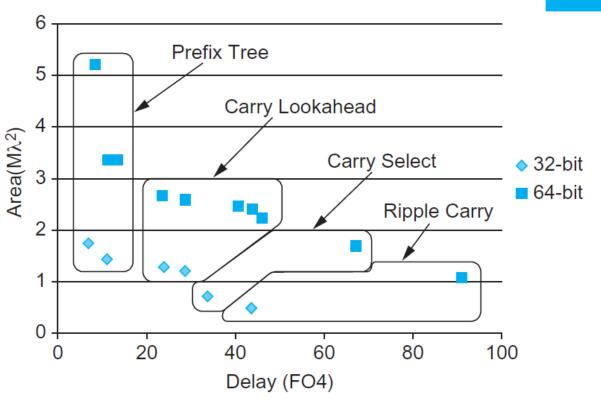
FO: Fan-Out

Everything has trade-off!

# Adder Family - Choose Wisely







FO: Fan-Out

Everything has trade-off!

### **Unsigned Integer Circuits - Subtractor**

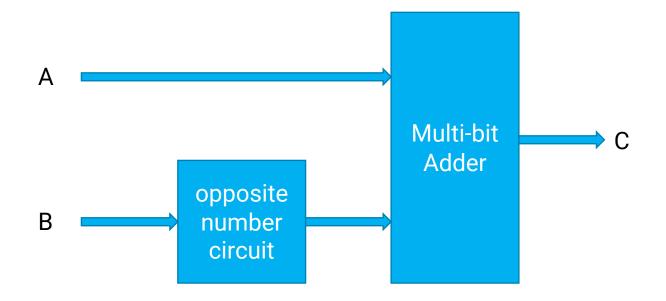


Numbers

Arithmetic

Circuits

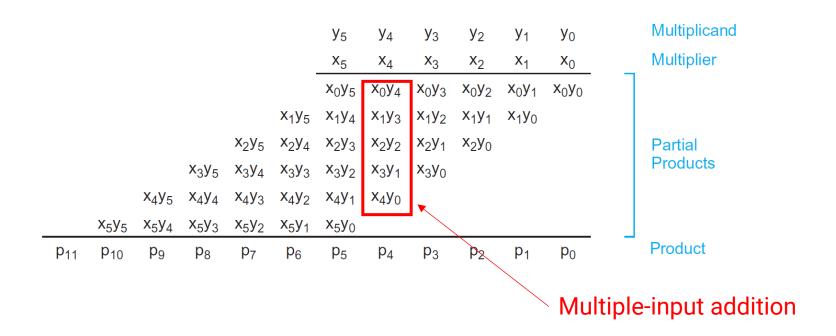
$$C = A-B = A+ (-B)$$



What is "opposite number circuit" though? [Save for a moment later]



Numbers Arithmetic Circuits



How do you implement it?

Yes! Adders!



Numbers

Arithmetic

 $C_3$   $S_3$ 

 $C_2 S_2$ 

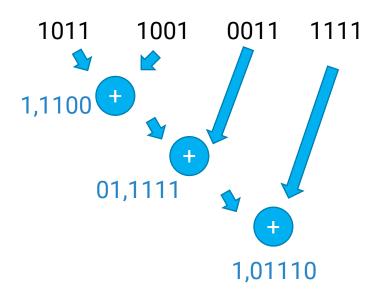
 $X_{N...1}$   $Y_{N...1}Z_{N...1}$ 

n-bit CSA

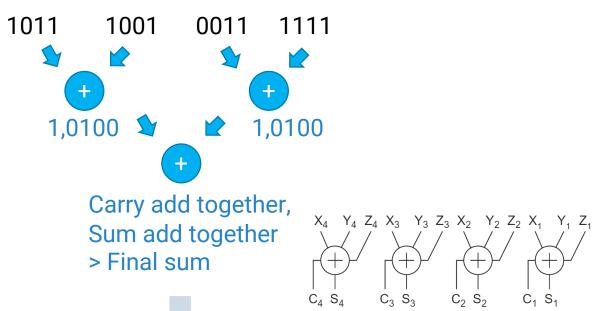
Circuits

Carry-Save Adder (CSA) & "carry-save redundant format"

Example: Sum of 1011 1001 0011 1111



(Carry, Sum)



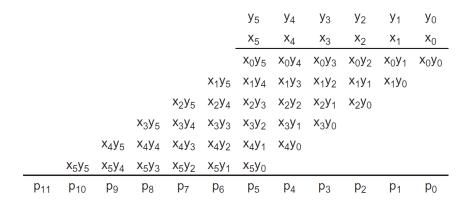


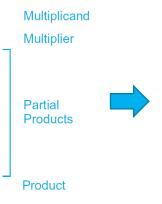


**Numbers** 

> Arithmetic

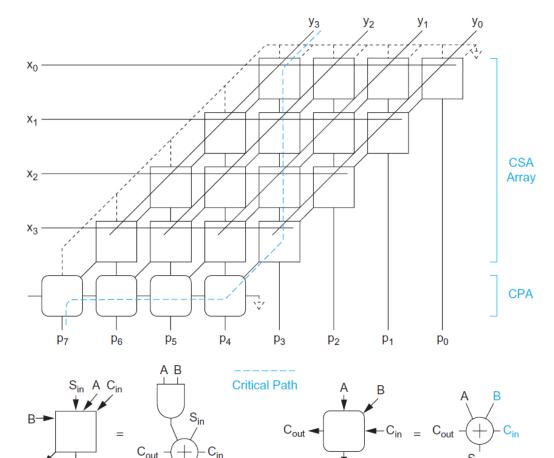
Circuits





#### Partial Product: Single-bit multiplication is equivalent to "AND"

| X | у | Partial product |
|---|---|-----------------|
| 0 | 0 | 0               |
| 0 | 0 | 0               |
| 1 | 0 | 0               |
| 1 | 1 | 1               |

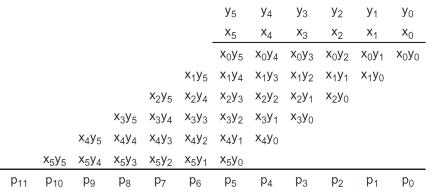


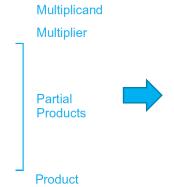




Numbers Arithmetic C

Circuits

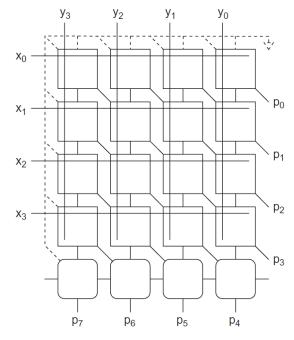




#### Partial Product: Single-bit multiplication is equivalent to "AND"

| X | у | Partial product |
|---|---|-----------------|
| 0 | 0 | 0               |
| 0 | 0 | 0               |
| 1 | 0 | 0               |
| 1 | 1 | 1               |





$$S_{\text{in}} \land C_{\text{in}}$$

$$S_{\text{in}} \land C_{\text{in}}$$

$$C_{\text{out}} \land C_{\text{in}}$$

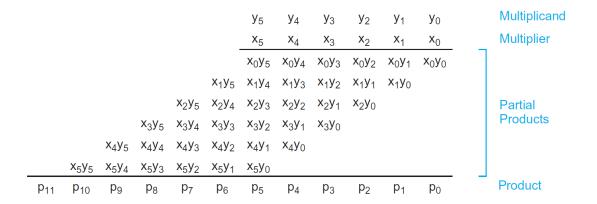




Numbers

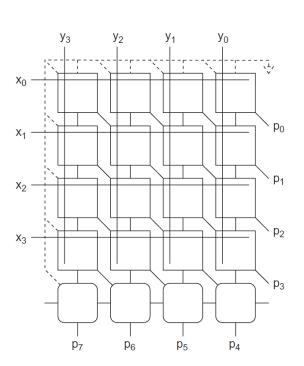
Arithmetic

Circuits



### **Anyway to Optimize This?**

**Booth Encoding** 





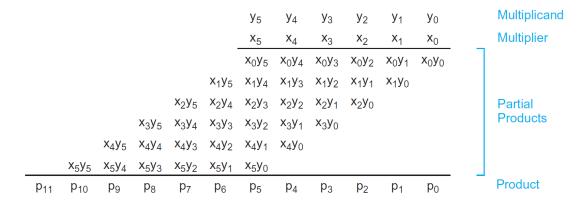
#### **Unsigned Integer Circuits - Multiply-Add**



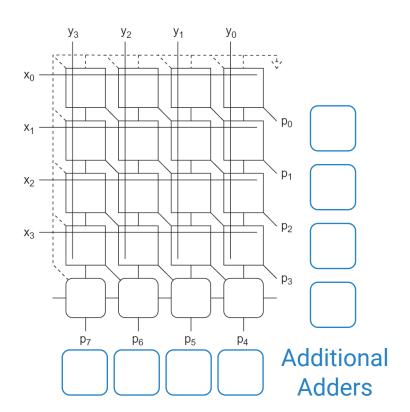
Numbers

Arithmetic

Circuits



# **Application-Specific Optimization**Fused-Multiply-Add (FMA)

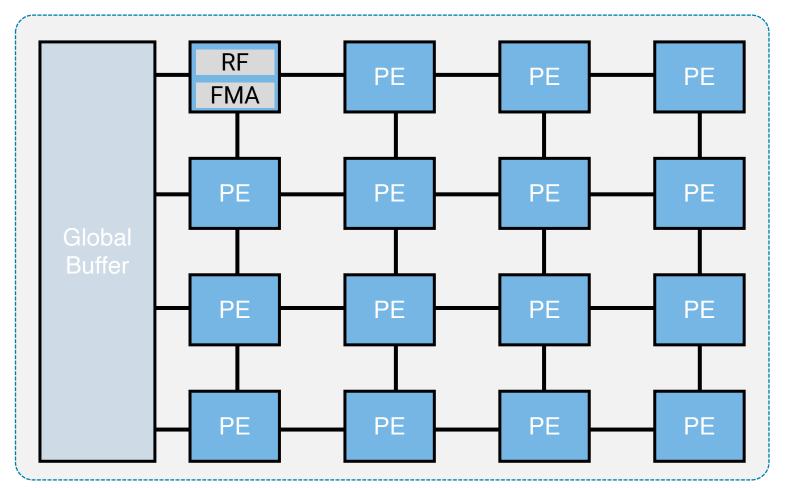




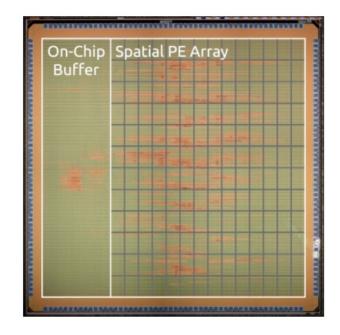
#### Multiply-Add Application Example



Eyeriss: CNN Accelerator



- Numbers Arithmetic Circuits
- FMA as processing elements
- local register file (RF)



Chen, Yu-Hsin, et al. "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks." *IEEE journal of solid-state circuits* 52.1 (2016): 127-138.

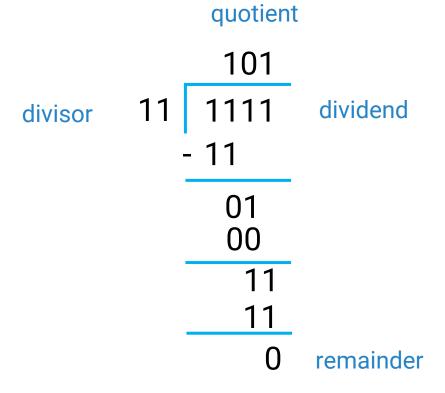
#### **Unsigned Integer Circuits - Divider**

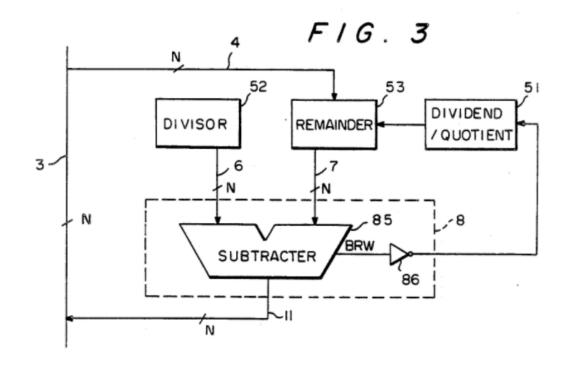


Numbers

> Arithmetic

Circuits





Yamahata, Hitoshi. "Integer division circuit provided with a overflow detector circuit." U.S. Patent No. 4,992,969. 12 Feb. 1991.



# Part 2

Signed Integers

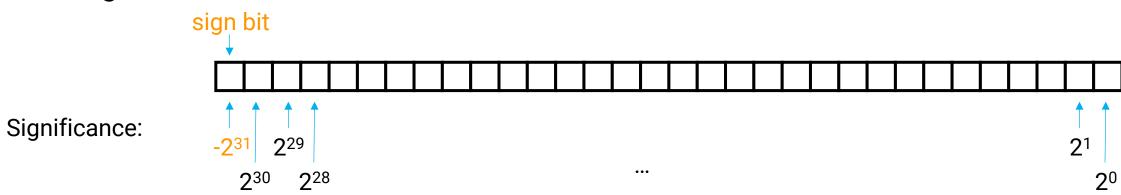


### Signed Integer



Numbers Arithmetic Circuits

Signed INT32



- INT16, INT8, ...
- Example:

-32'd7 (=32'hff\_ff\_ff\_f9) 8'b0100\_1101 (=8'hCD)

| Signed Bit | Meaning  |
|------------|----------|
| 0          | Positive |
| 1          | Negative |

A Useful Tool: Cryptii



#### **Signed Integer**



Problem: What is the range of signed vs. unsigned integers?

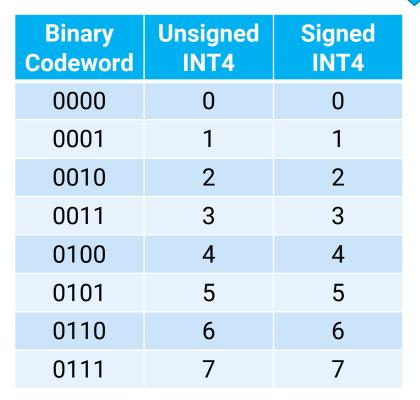
Numbers Arithmetic Circuits

For UINTn: 0~2<sup>n</sup>-1

For INTn:  $-2^{n-1} \sim 2^{n-1}-1$ 

# 

#### Here is the answer of "opposite number circuits"!



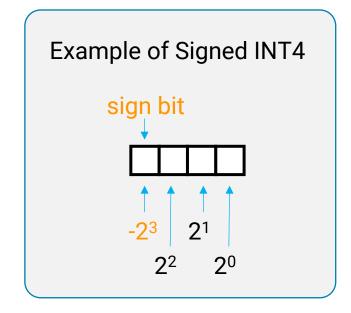
| Binary<br>Codeword | Unsigned<br>INT4 | Signed<br>INT4 |
|--------------------|------------------|----------------|
| 1000               | 8                | -8             |
| 1001               | 9                | -7             |
| 1010               | 10               | -6             |
| 1011               | 11               | -5             |
| 1100               | 12               | -4             |
| 1101               | 13               | -3             |
| 1110               | 14               | -2             |
| 1111               | 15               | -1             |
|                    |                  |                |



#### Signed Integer - Two Types of Shift



- Verilog HDL supports 2 types of shift:
- Logic Shift Operators (<<, >>)
- Arithmetic Shift Operators (<<<, >>>)



- Numbers Arithmetic Circuits
- Logic shift >> <<: filling with zeros</li>
- 3'b100 >> 1'd1 gives 3'b010
- 3'b101 >> 1'd1 gives 3'b010
- 3'b101 << 1'd2 gives 3'b100

Not really stable rule: <<1: multiply 2; >>2: divided by 2

- Arithmetic shift:
  - <<: Shift left specified number of bits, filling with zero.
  - >>>: Shift right specified number of bits, fill with value of sign bit if expression is signed, othewise fill with zero.



#### Signed Integer - Two Types of Shift



Verilog HDL supports 2 types of shift:

| • | Logic Shift Operators | (<<, >>) | ) |
|---|-----------------------|----------|---|
|---|-----------------------|----------|---|

Arithmetic Shift Operators (<<<, >>>)

| Binary Codewor | d Unsigned INT4 | Signed INT4 |
|----------------|-----------------|-------------|
| 1000           | 8               | -8          |
| 1001           | 9               | -7          |
| 1010           | 10              | -6          |
| 1011           | 11              | -5          |
| 1100           | 12              | -4          |
| 1101           | 13              | -3          |
| 1110           | 14              | -2          |
| 1111           | 15              | -1          |

4'b1110 >>> 1'd1 gives 4'b1111 4'b0110 >>> 1'd1 gives 4'b0011

- Numbers Arithmetic Circuits
- Logic shift >> <<: filling with zeros</li>
- 3'b100 >> 1'd1 gives 3'b010
- 3'b101 >> 1'd1 gives 3'b010
- 3'b101 << 1'd2 gives 3'b100

Not really stable rule: <<1: multiply 2; >>1: divided by 2

- Arithmetic shift:
  - <<: Shift left specified number of bits, filling with zero.
  - >>>: Shift right specified number of bits, fill with value of sign bit if expression is signed, othewise fill with zero.



# Part 3 About Fraction

Fixed-Point

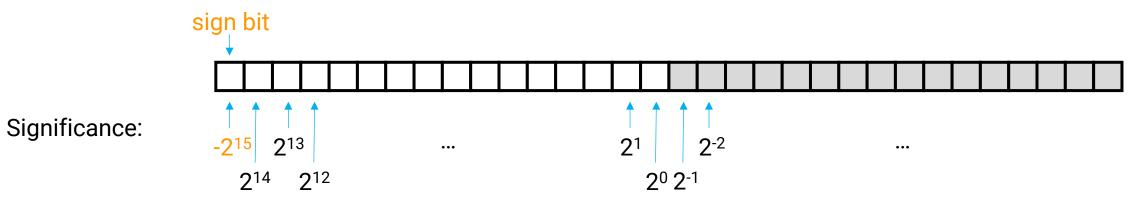


#### **Fixed-Point Number**



Numbers Arithmetic Circuits

• Fixed32



• Example:

$$(1.10)_2 = (1*2^0+1*2^{-1}+0*2^{-2})_{10} = (1.50)_{10}$$

| Signed Bit | Meaning  |
|------------|----------|
| 0          | Positive |
| 1          | Negative |

#### **Fixed-Point Number**



Arithmetic Just Works the Same Way!

**Numbers** Arithmetic Circuits

Verilog HDL does not support fixed-point natively



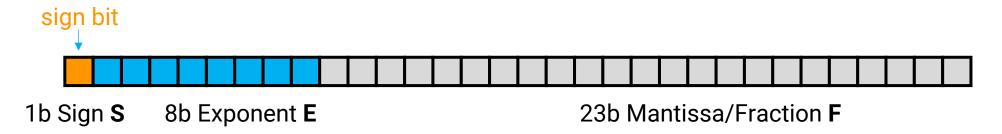
# Part 4

Floating-Point



Numbers Arithmetic Circuits

• FP32

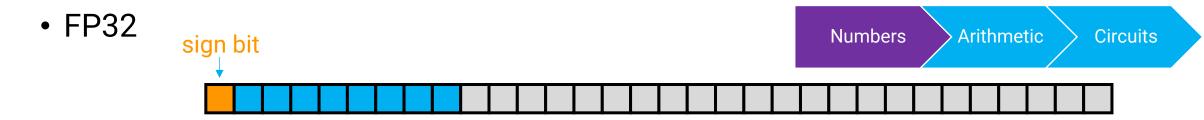


## Meaning:

| Exponent | Fraction | Object                | Value  |
|----------|----------|-----------------------|--|
| 0        | 0        | 0                     |  |
| 0        | Nonzero  | Denormalized number   | $(-1)^{\mathbf{S}} \times (0.\mathbf{F}) \times 2^{\mathbf{E}-\mathbf{B}}$ |
| Nonzero  | Anything | Floating-point number | $(-1)^{\mathbf{S}} \times (1.\mathbf{F}) \times 2^{\mathbf{E}-\mathbf{B}}$ |
| All "1"  | 0        | infinity              |  |
| All "1"  | Nonzero  | NaN (not a number)    |  |

1b Sign **S** 





23b Mantissa/Fraction F

| Exponent | Fraction | Object                | Value                              |
|----------|----------|-----------------------|------------------------------------|
| 0        | 0        | 0                     |                                    |
| 0        | Nonzero  | Denormalized number   | $(-1)^{S}\times(0.F)\times2^{E-B}$ |
| Nonzero  | Anything | Floating-point number | $(-1)^{S}\times(1.F)\times2^{E-B}$ |
| All "1"  | 0        | infinity              |                                    |
| All "1"  | Nonzero  | NaN (not a number)    |                                    |

$$S=1$$
,  $E=0$ ,  $F=0$ , what is the value?  $-0=0$ 

8b Exponent E





• FP32

Numbers Arithmetic Circuits

| Exponent | Fraction | Object                | Value                                  |
|----------|----------|-----------------------|--|
| 0        | 0        | 0                     |  |
| 0        | Nonzero  | Denormalized number   | $(-1)^{S} \times (0.F) \times 2^{1-B}$ |
| Nonzero  | Anything | Floating-point number | $(-1)^{S}\times(1.F)\times2^{E-B}$     |
| All "1"  | 0        | infinity              |  |
| All "1"  | Nonzero  | NaN (not a number)    |  |

| Туре                 | Sign | Exponent | Exponent bias | significand | total |
|----------------------|------|----------|---------------|-------------|-------|
| Half (IEEE 754-2008) | 1    | 5        | 15            | 10          | 16    |
| Single               | 1    | 8        | <b>B</b> 127  | 23          | 32    |
| Double               | 1    | 11       | 1023          | 52          | 64    |
| Quad                 | 1    | 15       | 16383         | 112         | 128   |

S=0, E=0, F=23'b10000\_00000\_00000\_0000 what is its decimal value?

$$(0.1)_2 \times 2^{-127} = 0.5 \times 2^{-127}$$





• FP32

| Numbers | Arithmetic | Circuits |  |
|---------|------------|----------|--|
|         |            |          |  |

| Exponent | Fraction | Object                | Value  |
|----------|----------|-----------------------|--|
| 0        | 0        | 0                     |  |
| 0        | Nonzero  | Denormalized number   | $(-1)^{S} \times (0.F) \times 2^{1-B}$                                     |
| Nonzero  | Anything | Floating-point number | $(-1)^{\mathbf{S}} \times (1.\mathbf{F}) \times 2^{\mathbf{E}-\mathbf{B}}$ |
| All "1"  | 0        | infinity              |  |
| All "1"  | Nonzero  | NaN (not a number)    |  |

| Туре                 | Sign | Exponent | Exponent bias | significand | total |
|----------------------|------|----------|---------------|-------------|-------|
| Half (IEEE 754-2008) | 1    | 5        | 15            | 10          | 16    |
| Single               | 1    | 8        | <b>B</b> 127  | 23          | 32    |
| Double               | 1    | 11       | 1023          | 52          | 64    |
| Quad                 | 1    | 15       | 16383         | 112         | 128   |

S=0, E=8'b127, F=23'b10000\_00000\_00000\_0000 what is its decimal value?

$$(1.1)_2 \times 2^{127-127} = 1.5$$





• Range

| Numbers | Arithmetic | Circuits |  |
|---------|------------|----------|--|
|         |            |          |  |

| Exponent | Fraction | Object                | Value                                  |
|----------|----------|-----------------------|--|
| 0        | 0        | 0                     |  |
| 0        | Nonzero  | Denormalized number   | $(-1)^{S} \times (0.F) \times 2^{1-B}$ |
| Nonzero  | Anything | Floating-point number | $(-1)^{S}\times(1.F)\times2^{E-B}$     |
| All "1"  | 0        | infinity              |  |
| All "1"  | Nonzero  | NaN (not a number)    |  |





Range

| Numbers | Arith | nmetic | Circuits |  |
|---------|-------|--------|----------|--|
|         |       |        |          |  |

| Exponent | Fraction | Object                | Value                              |
|----------|----------|-----------------------|------------------------------------|
| 0        | 0        | 0                     |                                    |
| 0        | Nonzero  | Denormalized number   | $(-1)^{S}\times(0.F)\times2^{1-B}$ |
| Nonzero  | Anything | Floating-point number | $(-1)^{S}\times(1.F)\times2^{E-B}$ |
| All "1"  | 0        | infinity              |                                    |
| All "1"  | Nonzero  | NaN (not a number)    |                                    |

#### Denormalized:





Numbers Arithmetic Circuits

### Add

```
123456.7 = 1.234567 * 10^5

101.7654 = 1.017654 * 10^2 = 0.001017654 * 10^5

Hence:

123456.7 + 101.7654 = (1.234567 * 10^5) + (1.017654 * 10^2)

= (1.234567 * 10^5) + (0.001017654 * 10^5)

= (1.234567 + 0.001017654) * 10^5

= 1.235584654 * 10^5
```

```
E=5; F=1.234567 (123456.7)

+ E=2; F=1.017654 (101.7654)

E=5; F=1.234567

+ E=5; F=0.001017654 (after shifting) Round-off error

E=5; F=1.235584654 (true sum: 123558.4654)
```

Actually, result is: e=5; s=1.235585 (final sum: 123558.5)

Try by yourself: (E=5, F=1.234567) + (E=-3, F=9.876543) = ??





Numbers Arithmetic Circuits

### Subtract

```
Try by yourself: (E=5, F=1.234571) - (E=5, 1.234567) = ??
```

```
E=5; F=1.234571
- E=5; F=1.234567
------
E=5; F=0.000004
E=-1; F=4.000000 (after rounding/normalization)
```

Change to normalized form of FP numbers





Numbers Arithmetic Circuits

## Multiply:

```
E=3; F=4.734612

× E=5; F=5.417242

E=8; F=25.648538980104 (true product)

E=8; F=25.64854 (after rounding)

E=9; F=2.564854 (after normalization)
```

**Exponent: Sum Operation** 

Mantissa: Multiply Operation

Don't forget normalization

• Divide:

**Exponent: Subtract Operation** 

Mantissa: Divide Operation

Don't forget normalization

Q: What if normalized number multiplies denormalized number?





Numbers Arithmetic Circuits

# **Incompleteness of Floating-Point Arithmetic**

May not associative:

May not distributive:

$$1234.567 + 45.67844 = 1280.245$$

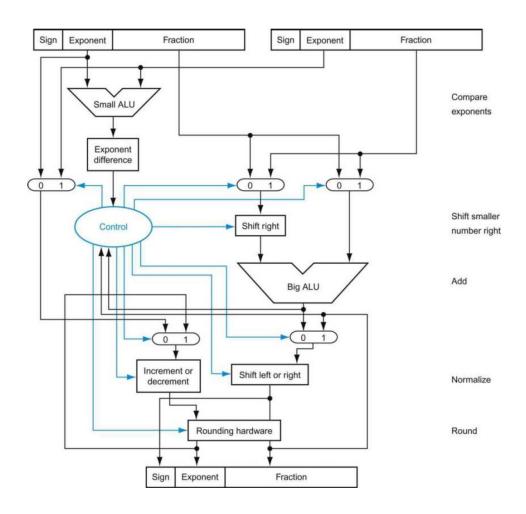
$$1280.245 + 0.0004 = 1280.245$$
but
$$45.67840 + 0.00004 = 45.67844$$

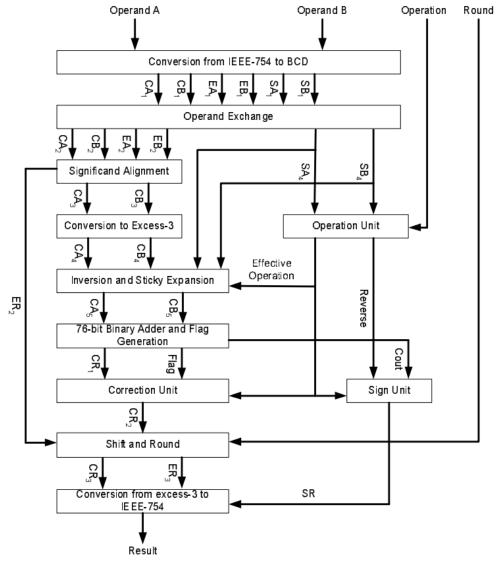
$$45.67844 + 1234.567 = 1280.246$$

Arithmetic



#### Adder:

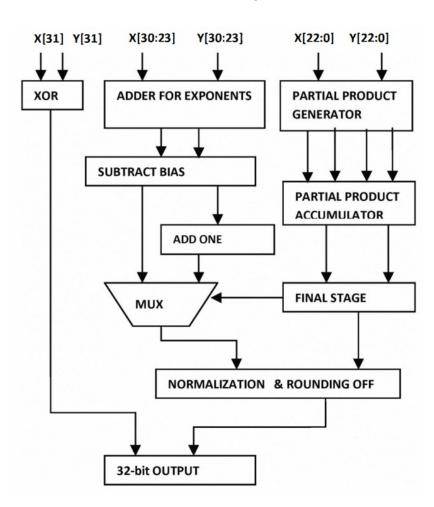






Multiply:





Sign: XOR

Exponent: Add

Mantissa: Multiply

Jain, Anna, et al. "FPGA design of a fast 32-bit floating point multiplier unit." 2012 International Conference on Devices, Circuits and Systems (ICDCS). IEEE, 2012.





Numbers Arithmetic Circuits

## Floating Point ALU supports +-\*/

#### **RISC-V floating-point assembly language**

| Category      | Instruction                   | Example           | Meaning                    | Comments                           |
|---------------|-------------------------------|-------------------|----------------------------|------------------------------------|
|               | FP add single                 | fadd.s f0, f1, f2 | f0 = f1 + f2               | FP add (single precision)          |
|               | FP subtract single            | fsub.s f0, f1, f2 | f0 = f1 - f2               | FP subtract (single precision)     |
|               | FP multiply single            | fmul.s f0, f1, f2 | f0 = f1 * f2               | FP multiply (single precision)     |
|               | FP divide single              | fdiv.s f0, f1, f2 | f0 = f1 / f2               | FP divide (single precision)       |
| Arithmetic    | FP square root single         | fsqrt.s f0, f1    | f0 = √f1                   | FP square root (single precision)  |
|               | FP add double                 | fadd.d f0, f1, f2 | f0 = f1 + f2               | FP add (double precision)          |
|               | FP subtract double            | fsub.d f0, f1, f2 | f0 = f1 - f2               | FP subtract (double precision)     |
|               | FP multiply double            | fmul.d f0, f1, f2 | f0 = f1 * f2               | FP multiply (double precision)     |
|               | FP divide double              | fdiv.d f0, f1, f2 | f0 = f1 / f2               | FP divide (double precision)       |
|               | FP square root double         | fsqrt.d f0, f1    | f0 = √f1                   | FP square root (double precision)  |
|               | FP equality single            | feq.s x5, f0, f1  | x5 = 1 if f0 == f1. else 0 | FP comparison (single precision)   |
|               | FP less than single           | flt.s x5, f0, f1  | x5 = 1 if f0 < f1, else 0  | FP comparison (single precision)   |
| Comparison    | FP less than or equals single | fle.s x5, f0, f1  | x5 = 1 if f0 <= f1, else 0 | FP comparison (single precision)   |
| Comparison    | FP equality double            | feq.d x5, f0, f1  | x5 = 1 if f0 == f1, else 0 | FP comparison (double precision)   |
|               | FP less than double           | flt.d x5, f0, f1  | x5 = 1 if f0 < f1, else 0  | FP comparison (double precision)   |
|               | FP less than or equals double | fle.d x5, f0, f1  | x5 = 1 if f0 <= f1, else 0 | FP comparison (double precision)   |
|               | FP load word                  | flw f0, 4(x5)     | f0 = Memory[x5 + 4]        | Load single-precision from memory  |
| Data transfer | FP load doubleword            | fld f0, 8(x5)     | f0 = Memory[x5 + 8]        | Load double-precision from memory  |
|               | FP store word                 | fsw f0, 4(x5)     | Memory[x5 + 4] = f0        | Store single-precision from memory |
|               | FP store doubleword           | fsd f0, 8(x5)     | Memory[x5 + 8] = f0        | Store double-precision from memory |