

Fixed-Point Number (Addition/Multiplication)

Integers

□ Generally we use 8-bits, 16-bits, 32-bits or 64-bits to store integers.

□ $20 = 16 + 4 = (0001\ 0100)_2$ 8-bits

□ $20 = 16 + 4 = (0000\ 0000\ 0001\ 0100)_2$ 16-bits
←pad zeros→

□ We use 2's complement format for the notation of negative signed numbers:

$20 = (0\dots01\ 0100)_2$

$-20 = (1110\ 1100)_2$ 8-bits

$-20 = (1111\ 1111\ 1110\ 1100)_2$ 16-bits

←pad ones→
Sign bit

Integers

- ❑ How to store integers in registers?
- ❑ Consider that we have 8-bit registers.
- ❑ $20 = (10100)_2$
- ❑ As 8-bit integer: (r1)
 - $r1 = 20 = (0001\ 0100)_2$
- ❑ As 16-bit integer: (r1 r2)
 - $r1 = 0 = (0000\ 0000)_2$
 - $r2 = 20 = (0001\ 0100)_2$
 - $(r1\ r2) = 20 = (\underbrace{0000\ 0000}_{r_1}\ \underbrace{0001\ 0100}_{r_2})_2$
- ❑ As 32-bit integer: (r1 r2 r3 r4)
 - $r1 = r2 = r3 = 0 = (0000\ 0000)_2$
 - $r4 = 20 = (0001\ 0100)_2$
 - $(r1\ r2\ r3\ r4) = 20 = (\underbrace{0000\ 0000}_{r_1}\ \underbrace{0000\ 0000}_{r_2}\ \underbrace{0000\ 0000}_{r_3}\ \underbrace{0001\ 0100}_{r_4})_2$

Integers

□ Represent 123456789 in 32-bit integer:

- $123456789 = (111\ 0101\ 1011\ 1100\ 1101\ 0001\ 0101)_2$
- Convert to 32-bits: $\approx 7\text{-bits}$
- $0000\ 0111\ 0101\ 1011\ 1100\ 1101\ 0001\ 0101$ $r_1\ r_2\ r_3\ r_4$
- $r1 = (0000\ 0111)_2 = 0x07 = 7$
- $r2 = (0101\ 1011)_2 = 0x5b = 91$
- $r3 = (1100\ 1101)_2 = 0xcd = 205$
- $r4 = (0001\ 0101)_2 = 0x15 = 21$

- $(r1\ r2\ r3\ r4) = 0x075bcd15$
 $= (0000\ 0111\ 0101\ 1011\ 1100\ 1101\ 0001\ 0101)_2$
 $= 123456789$

Integers (2's complement)

- Given following values of registers, find the value of (r4 r3 r2 r1)?
- r1 = 72, r2 = 100, r3 = 250, r4 = 255

$$r1 = 72 = (0100\ 1000)_2$$

$$r2 = 100 = (0110\ 0100)_2$$

$$r3 = 250 = (1111\ 1010)_2$$

$$r4 = 255 = (1111\ 1111)_2$$

$$(r4\ r3\ r2\ r1) = (1111\ 1111\ 1111\ 1010\ 0110\ 0100\ 0100\ 1000)_2$$

The number is negative! Take 2's complement:

$$(0000\ 0000\ 0000\ 0101\ 1001\ 1011\ 1011\ 1000)_2 = 367\ 544$$

$$= -367544$$

8-bits
↓
↘
↘
↘

negative \Rightarrow what's the negative number?

Integer Additon

- Assume that you have an operator that adds only two digits:

$$\begin{array}{r} \\ + \\ \hline \end{array}$$

Sum is a digit.
But carry is just a bit, can either be zero or one.

carry → C S ← *sum*

Each digit is a number in a base b .

$b=10$ \Rightarrow numbers: 0-9 ←

$b=2$ \Rightarrow numbers: 0-1 ←

$b=2^8=256$ \Rightarrow numbers: 0-255 ←

Note that the sum of two single-digit yields one digit and extra one bit at most!

Integer Additon

- Assume that we have an operator that adds only two digit. How can we add two numbers with multiple digits?

$$\begin{array}{r} 5639 \\ + 1427 \\ \hline \end{array}$$

?

Solution: Add digits individually
 Also add carry!

Integer Additon

$$\begin{array}{r} 5639 \\ + 1427 \\ \hline 6 \text{ (Carry=1)} \end{array}$$

Integer Additon

$$\begin{array}{r} 56\mathbf{3}9 \\ + 14\mathbf{2}7 \\ \hline \mathbf{6}6 \text{ (Carry=0)} \end{array}$$

Integer Additon

$$\begin{array}{r} 0 \\ 5639 \\ + 1427 \\ \hline 066 \text{ (Carry=1)} \end{array}$$

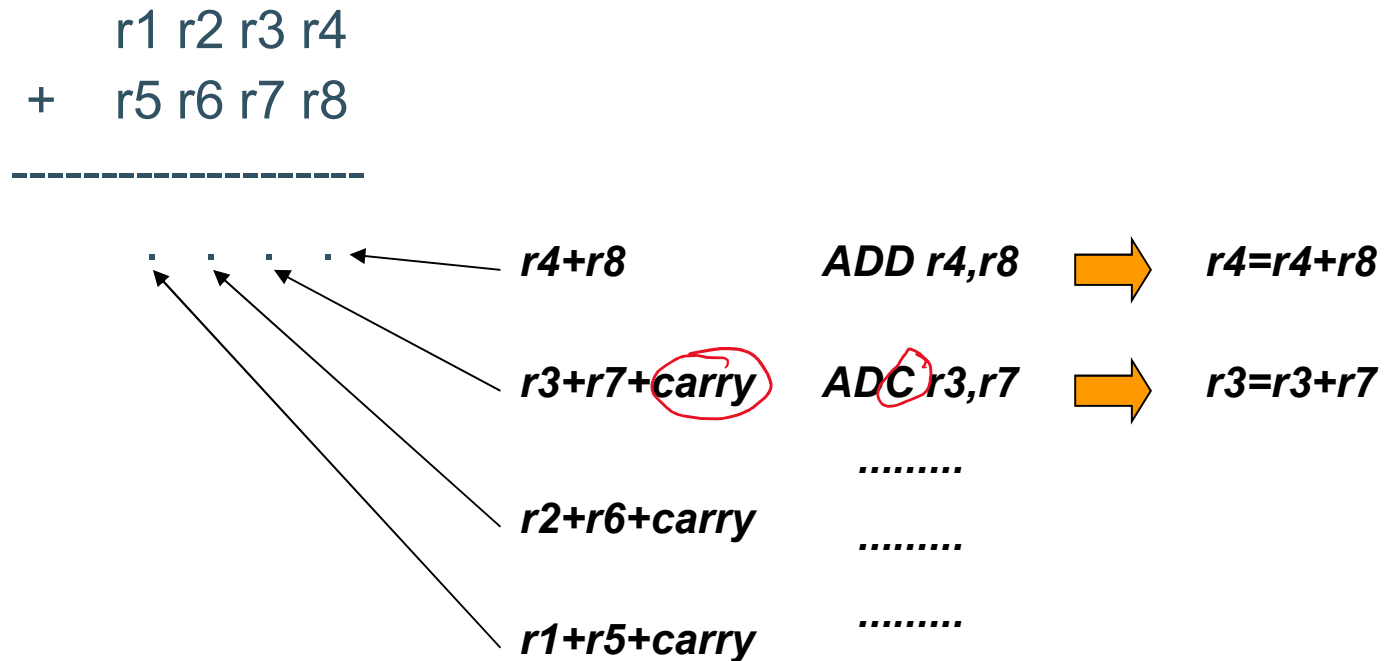
Integer Additon

$$\begin{array}{r} 1 \\ 5639 \\ + 1427 \\ \hline 7066 \text{ (Carry=0)} \end{array}$$

Integer Additon

Now consider that we are working in base 256.

Put each digit in a register so that we'll have 4 register for each 32-bit number.



Integer Additon

- ❑ What about signed numbers?
- ❑ Use 2's complement for negative numbers and just add! **Ignore the last produced carry.**
- ❑ How does it work? Explained later...
- ❑ What about subtraction?
- ❑ Subtraction can easily be implemented by taking 2's complement of the second operand first and then applying addition:
 - $A - B = A + (-B)$

Integer Multiplication

- Assume that you have an operator that multiplies only two digits:

$$\begin{array}{r} \\ x \\ \hline C \end{array}$$

Each digit is a number in a base b .

$b=10 \quad \Rightarrow \text{ numbers: } 0-9$

$b=2 \quad \Rightarrow \text{ numbers: } 0-1$

$b=2^8=256 \quad \Rightarrow \text{ numbers: } 0-255$

Note that the product of two single-digit yields two digits at most!

Integer Multiplication

- You have more than one digit to multiply:

$$\begin{array}{r} 58 \\ \times 37 \\ \hline \end{array}$$

?

By using the operator, we can calculate:

$$7 \times 8 = 56$$

$$7 \times 5 = 35$$

$$3 \times 8 = 24$$

$$3 \times 5 = 15$$

Integer Multiplication

- How can we use these values to calculate the result?

$$\begin{array}{r} 58 \\ \times 37 \\ \hline 56 \\ 35 \\ 24 \\ + 15 \\ \hline \end{array}$$

Integer Multiplication

- We can use integer addition to find the result.

$$\begin{array}{r}
 58 \\
 \times 37 \\
 \hline
 56 \\
 35 \\
 24 \\
 + 15 \\
 \hline
 ?
 \end{array}$$

$$\begin{array}{r}
 + \quad \boxed{056} \\
 \hline
 406 \leftarrow \text{Sum 1} \\
 + \quad \boxed{24} \\
 \hline
 0646 \leftarrow \text{Sum 2} \\
 + \quad \boxed{15} \\
 \hline
 2146 \leftarrow \text{Sum 3}
 \end{array}$$

- *This operation is equivalent to 16-bit multiplication using 8-bit multiplication and 8-bit addition.*
- *Note that the number of digits in the result is equal to the sum of the number of input digits.*

Integer Multiplication

- Now assume that the digits are in base $2^8 = 256$.
(8-bit are necessary for each digit)
- Then, **16-bit multiplication is done by using 8-bit multiplication and 8-bit addition**. Each 8-bit register can hold only one digit!

$$\begin{array}{r} r1 \ r2 \\ x r3 \ r4 \\ \hline r5 \ r6 \\ r7 \ r8 \leftarrow \\ r9 \ r10 \leftarrow \\ + r11 \ r12 \leftarrow \\ \hline r13 \ r14 \ r15 \ r16 \end{array}$$

In fact, we do not need 16 registers to accomplish 16-bit multiplication.

If we compute the partial sums, we can re-use the registers which hold the values that are unnecessary.

Fixed-Point Numbers

- ❑ Fixed-point numbers are generally stored in “**In.Qm**” format (sometimes referred as Qn.m format)
- ❑ **n** = number of bits in integer part.
- ❑ **m** = number of bits in fractional part.
- ❑ Example: **I8.Q16**

2 ⁷	2 ⁶	2 ⁵	2 ⁴	2 ³	2 ²	2 ¹	2 ⁰		2 ⁻¹	2 ⁻²	2 ⁻³	2 ⁻⁴	2 ⁻⁵	2 ⁻⁶	2 ⁻⁷	2 ⁻⁸	2 ⁻⁹	2 ⁻¹⁰	2 ⁻¹¹	2 ⁻¹²	2 ⁻¹³	2 ⁻¹⁴	2 ⁻¹⁵	2 ⁻¹⁶
0	0	1	0	1	1	1	0	.	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0

$$= 32 + 8 + 4 + 2 + 1/2 + 1/4 + 1/16$$

$$= 46.8125$$

Signed Fixed-Point Numbers

- Positive fixed-point numbers are the same as unsigned fixed-point numbers.
- Negative fixed-point numbers are obtained by simply calculating 2's complement as they are integers.

unsigned 18.Q8: $01000110.11000000 = 70.75$
signed 2's comp. $10111001.01000000 = -70.75$

negative \downarrow 2's

$$\begin{array}{r}
 01000110.11000000 \\
 + \\
 01000110.11000000 \\
 \hline
 10111001.01000000 = 70.75
 \end{array}$$

Fixed-Point Addition

- ❑ Fixed-point addition is the same as integer addition!
- ❑ Align two fixed point number and apply integer addition:

$$\begin{array}{r} r1\ r2\ .\ r3\ r4 \\ +\ r5\ r6\ .\ r7\ r8 \\ \hline \end{array}$$

_____ . _____

Unsigned Fixed-Point Multiplication

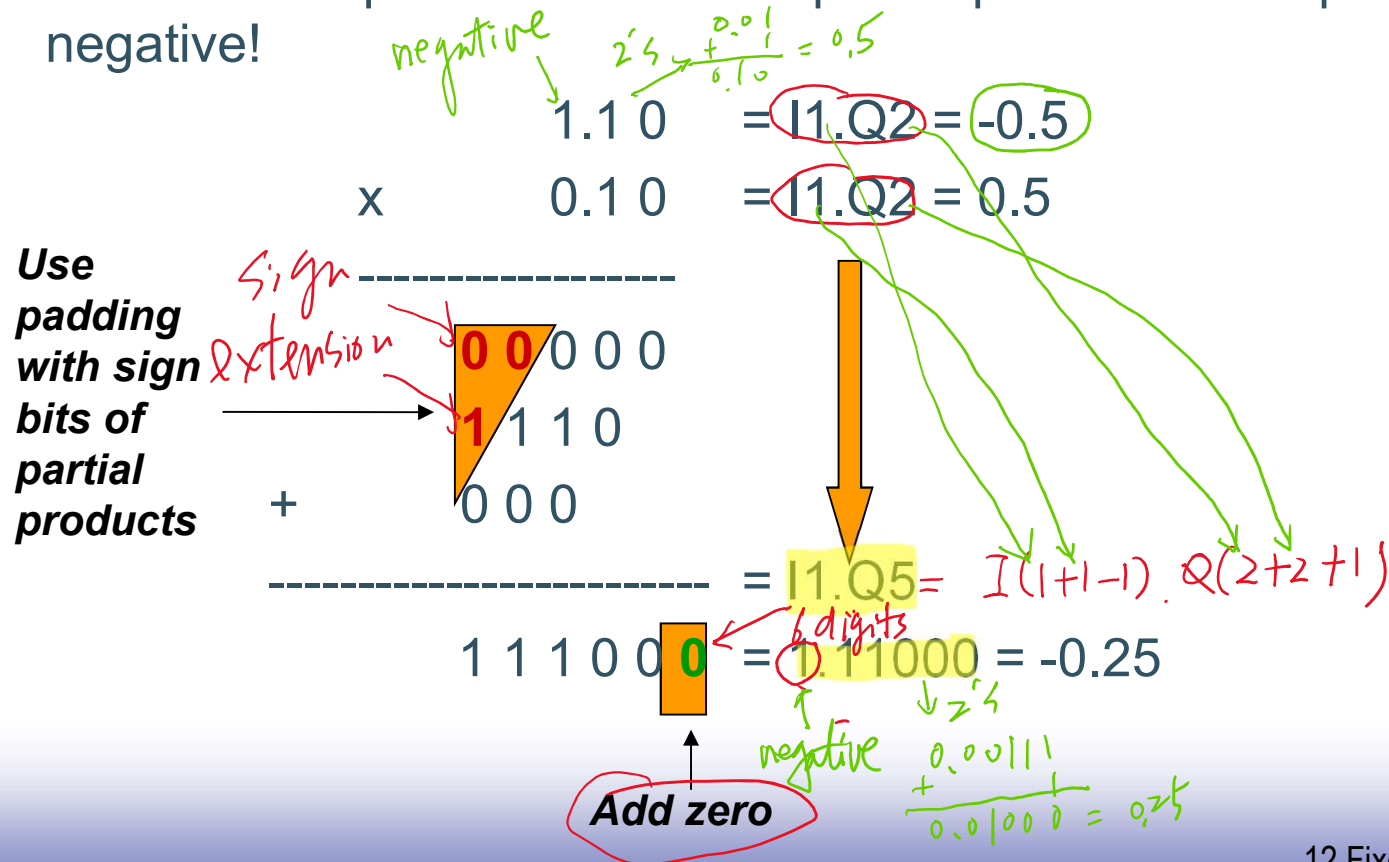
- Unsigned fixed-point multiplication is similar to integer multiplication.
- Consider the following multiplications:

58	5.8	I1.Q1	Ia.Qb
x 37	x 3.7	I1.Q1	Ic.Qd
-----	-----		
2146	21.46	I2.Q2	
$I(a+c).Q(b+d)$			

Just multiply like integer multiplication. Align the numbers according to the point (.)

Signed Fixed-Point Multiplication

- Use 2's complement format for fixed-point numbers.
- ~~□~~ $(Ia.Qb) * (Ic.Qd) = I(a+c-1).Q(b+d+1)$
- Take 2's complement of the last partial product if multiplier is negative!



Signed Fixed-Point Multiplication

□ Example:

negative

1 1.0 1 = 12.Q2 = -0.75

x 1.1 0 1 = 11.Q3 = -0.375

1 1 1 1 1 0 1

0 0 0 0 0 0 0

1 1 1 0 1

+ 0 0 1 1

0 0 0 1 0 0 1 0 = 12.Q6 = 00.010010 = 0.28125

discarded carry-out

8-digits

Add zero

2's complement of the last partial product

$I(2+1-1).Q(2+3+1) = I(2).Q(6)$
8-digits

$2's(1101) = \begin{array}{r} 0010 \\ + 1101 \\ \hline 0011 \end{array}$

The diagram illustrates the multiplication of two 4-bit signed fixed-point numbers: 11.01 (12.Q2) and 1.101 (11.Q3). The first number is negative, indicated by a green arrow and the word 'negative'. The second number is also negative. The partial products are shown as a 3x7 grid of bits. The first partial product is 1111101, the second is 0000000, and the third is 11101. These are added to get a sum of 0001001. A red arrow points to the first partial product, and a green arrow points to the third partial product. A green circle highlights the last bit of the third partial product (1), and a green arrow points to the 2's complement of this bit (11), which is added to the sum. The final result is 00010010, which is 12.Q6. The first two bits (00) are labeled 'discarded carry-out'. The last six bits (010010) are labeled '8-digits'. A green arrow points to the last bit of the result (0), which is labeled 'Add zero'.