INTEGER ARITHMETIC CAS CS 210 Computer System	
CSAPP 2.3	





	×	у	x+y	result	3bit
	0	0	000	0	000=0
w=3 $2^3=8$	1	1	1 001 001 010	2	010=2
min=000 max= 2^3-1	6	1	110 001 111	7	111=7
=1000-1 =111	7	1	111 111 001 1000	8	000=0
	6	4	1 110 100 1010	10	010=2
	7	7	111 111 111 111 1110	14	110=6



Unsigned A	ddition
Operands: <i>w</i> bits True Sum: <i>w</i> +1 bits Discard Carry: <i>w</i> bits	$\begin{array}{c c} u & & & \\ \hline \\ + v & & \\ \hline \\ u+v & \\ \hline \\ u+v & \\ \hline \\ UAdd_{u}(u,v) & \\ \hline \\ \end{array}$
Standard Additio	on Function
implements Mod	$\begin{aligned} \mathbf{u}_{\mathbf{u}} = \mathbf{u}_{\mathbf{v}} \\ \mathbf{u}_{\mathbf{u}} = \mathbf{u}_{\mathbf{v}} \\ \mathbf{v}_{\mathbf{v}} = \mathbf{u}_{\mathbf{v}} \\ \mathbf{v}_{\mathbf{v}} $

OVERFLOW

- OVERFLOW: when the full integer result cannot fit within the word size limits of the data type
- occurs when $x+y \ge 2^w$
- C does **not** signal overflow errors they occur silently. Unsigned overflow can be detected if (x+y) < x
- We will see that most processors do indicate if an overflow occurs
- Both addition and multiplication can suffer from it







					3bit
	0	0	000	0	000=0
w=3	2	1	010 001 011	3	011=3
2 ³ =8 Signed: min=100=-4	3	1	11 011 001 	4	100=-4
max=011=3	-3	-1	111 101 111 1100	-4	100=-4
	-4	-1	1 100 111 1011	-5	011=3
	-4	-4	1 100 100 1000	-8	000=0

Two's Complement Addition
Operands: <i>w</i> bits <i>u · · · · · · · · · · · · · · · · · · ·</i>
True Sum: w+1 bits $u + v$ Discard Carry: w bits TAdd _w (u, v)
TAdd and UAdd have Identical Bit-Level Behavior Signed vs. unsigned addition in C:
<pre>int s, t, u, v; s = (int) ((unsigned) u + (unsigned) v); t = u + v</pre>
• Will give s == t

OVERFLOW & "UNDER-FLOW"

- OVERFLOW: when the result exceeds 2's Complement Max
- UNDERFLOW: when result exceeds 2's Complement Min
- C does **<u>not</u>** signal either errors they occur silently.
- Again most processors do indicate if either occurs
- Both addition and multiplication can suffer from it





MULTIPLYING BY A CONSTANT (2.3.6)

- + HW support for multiplication but slow compared to $\sim, \&, |, +, >>, <<$
- So when it can the compiler tries to replace multiplication with equivalent computations using the more primitive ops
- key rule used is that both for unsigned and signed
- $x * 2^{k} = x \le k$ for $0 \le k \le w$ (see 2.3.6 for proof):
- eg. 3 * 2 = 3 << 1 and 3 * 4 = (3 << 1) << 1 = 3 << 2

 $x*14 \xrightarrow{\text{compiler } x} (2^3 + 2^2 + 2^1) = (x << 3) + (x << 2) + (x << 1)$ or even better $x * (2^4 - 2^1) = (x << 4) - (x << 1)$

DIVIDING BY PWRS OF 2 (2.3.7) Integer division defined to always round towards zero

The generation defined to always round towards zero 7/2 = 3 and -7/2 = -3

- like multiplication we would like to exploit right shifts to replace divisions by powers of two: x / 2^k → x >>k :0>k<w: eg.
 7/2 →7 >> 1 →[0111] >> 1 →[0011]=3
- this works for positive x as truncation will round towards 0 but it does not work for negative x that requires rounding eg.:

-7/2 = [1001]/2 supposed be -3 but

$$|00| >> | = [|00| = -4$$

• truncation operates like floor but for negative division we need ceiling.





REAL NUMBERS : FLOATING POINT

IEEE Standard 754

- Established in 1985 as uniform standard for floating point arithmetic
- Before that, many idiosyncratic formats
- Supported by all major CPUs
- Driven by Numerical Concerns
 - Nice standards for rounding, overflow, underflow
 - Hard to make go fast
 - Numerical analysts predominated over hardware types in defining standard









values	ne	ale		uie	Exponent
	Exp	exp	Е	2 ⁸	
	0	0000	-6	1/64	(denorms)
	1	0001	-6	1/64	
	2	0010	-5	1/32	
	3	0011	-4	1/16	
	4	0100	-3	1/8	
	5	0101	-2	1/4	
	6	0110	-1	1/2	
	7	0111	0	1	
	8	1000	+1	2	
	9	1001	+2	4	
	10	1010	+3	8	
	11	1011	+4	16	
	12	1100	+5	32	
	13	1101	+6	64	
	14	1110	+7	128	
	15	1111	n/a		(inf. NaN)

Dyna	imic	Ra	nge	
	s exp	frac	E	Value
	0 0000	000	-6	0
	0 0000	001	-6	1/8*1/64 = 1/512 Closest to zero
Denormalized numbers	0 0000	010	-6	2/8*1/64 = 2/512
	0 0000	110	-6	$6/8 \times 1/64 = 6/512$
	0 0000	111	-6	7/8*1/64 = 7/512 ← largest denorm
	0 0001	000	-6	8/8*1/64 = 8/512 ← smallest norm
	0 0001	001	-6	9/8*1/64 = 9/512
		110		14/9+1/2 = 14/16
	0 0110	110	11	$14/6^{-1/2} = 14/16$
Normalized	0 0110	111	-1	15/8*1/2 = 15/16
numbers	0 0111	000	0	8/8*1 = 1
	0 0111	001	0	
	0 0111	010	0	$10/8 \times 1 = 10/8$
		110	-	14/0+100 - 004
	0 1110	110	4	14/8*128 = 224 ← largest norm
	0 1110	111		15/8*128 = 240
	0 1111	000	n/a	INT

←closest to zero			
← largest denorm			
← smallest norm			
← closest to 1 below			
← closest to 1 above			
← largest norm			



tribution of Values ose-up view)	
EEE-like format e = 3 exponent bits	
' = 2 fraction bits ∃ias is 3	
-0.5 0 0.5 Denormalized A Normalized II Infinity	1

EXTRA	

WHY DOES 2'S COMPLEMENT	
WORK	

















|--|



TWO'S-COMPLEMENT MULTIPLICATION (2.3.5)

bit-level representation of unsigned and 2's complement multiplication the same (see 2.3.5 for simple proof):

$x\ast^{\mathbf{t}}_{w}y=U2T_{w}((x\cdot y)\mathrm{mod}2^{w})$

some simple 3 bit examples

mode	×		у		х • у		trunc x • y	
unsigned	5	[101]	3	[011]	15	[001111]	7	[111]
2's comp.	-3	[101]	3	[011]	-9	[110111]	-	[111]
unsigned	4	[100]	7	[111]	28	[011100]	4	[100]
2's comp.	-4	[100]	-	[111]	4	[000100]	-4	[100]
unsigned	3	[011]	3	[011]	9	[001001]	I	[001]
2's comp.	3	[011]	3	[011]	9	[001001]	I	[001]

FLOATING POINT	













Special Properties of Encoding

FP Zero Same as Integer Zero

All bits = 0

- Can (Almost) Use Unsigned Integer Comparison
 - Must first compare sign bits
 - Must consider -0 = 0
 - NaNs problematic
 - Will be greater than any other valuesWhat should comparison yield?
 - Otherwise OK
 - Denorm vs. normalized
 - Normalized vs. infinity