

$50000 * 50000 \rightarrow ???$ negative

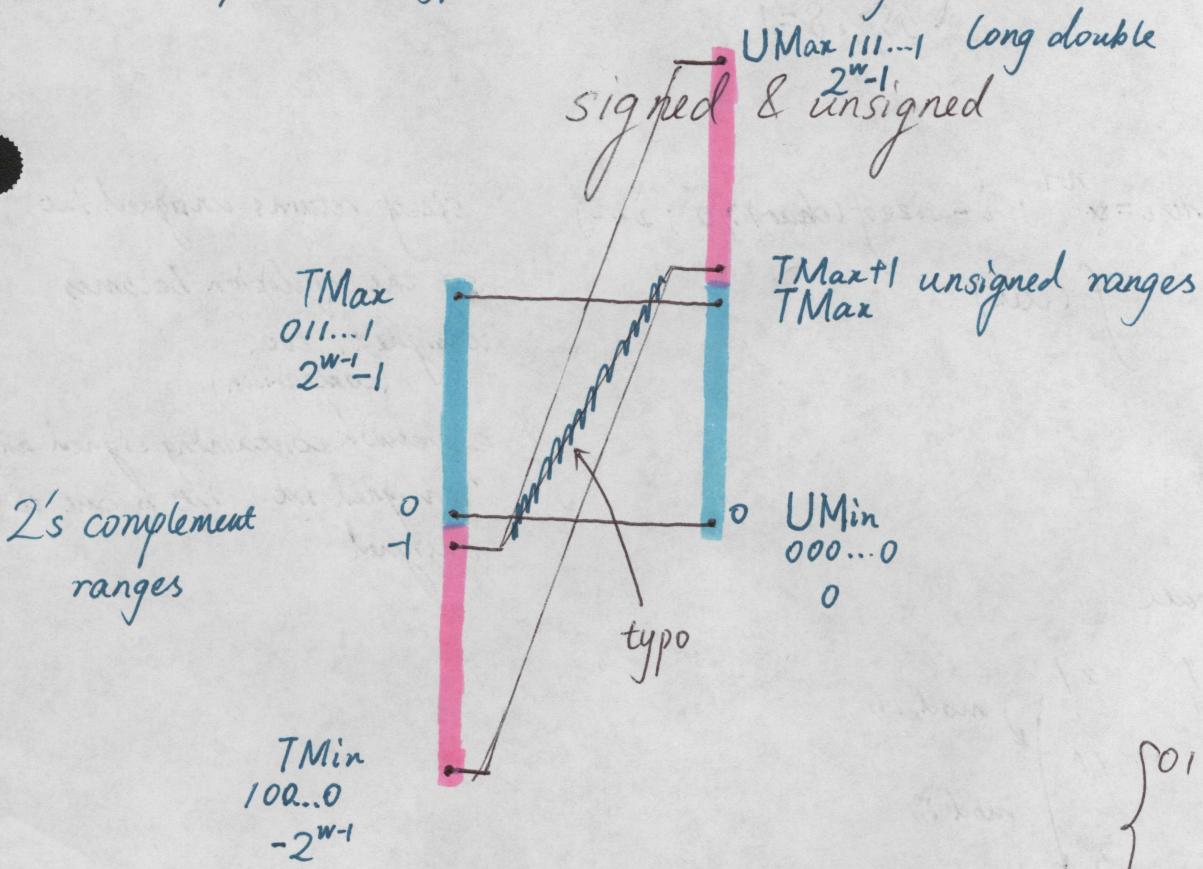
$300 * 400 * 500 * 600 \rightarrow ???$ 1640261632 overflow

fact: we finite representation to represent things that are potentially infinite

so there are compromises: overflow, round-off problems

Java primitive types: byte short int long float double boolean char

C primitive types: ^{data} short int long float double char



$$|T_{Max}| = |T_{Min}| - 1$$

$$UMax = 2 * T_{Max} + 1$$

$$\left\{ \begin{array}{l} 011\dots 1 \\ \dots \\ 000\dots 0 \\ \hline 111\dots 1 \\ \dots \\ 100\dots 0 \end{array} \right.$$

2^{w-1}

2^{w-1}

reason of asymmetry

if ($x < 0$) give TMin as input, then $\overset{\text{not}(x)}{x} \rightarrow \overset{+1}{\sqrt{+1}}$ "complement?
 return $-x$; which is -2147483648
 else bit pattern is maintained,
 return x ; but reinterpreted

~~for (unsigned int i = $\overset{n-1}{\alpha}$; i > 0; i--) f(a[i])~~

Loop will go forever, will not
 stop as you want
 an unsigned int is never negative

int $x = 4$
 printf("%d", $x \gg 65$)
 $\uparrow 65 \% 8 = 1$

the compiler warns you
 but will print 2

~~for (int i = $\overset{n-1}{\alpha}$; i - sizeof(char) > 0; i--) f(a[i])~~

sizeof returns unsigned int,
 so the condition becomes
 unsigned, too.
 (comparison)

~~f expression containing signed and
 unsigned int, int is cast to
 unsigned~~

unsigned truncate

$\begin{array}{r} 1 \\ 16 \\ 8 \\ 4 \\ 2 \\ 1 \end{array}$	$\begin{array}{r} 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1 \end{array}$	$\begin{array}{r} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array}$	$\begin{array}{r} 27 \\ 11 \\ 3 \\ 2 \\ 1 \end{array}$	$\left. \begin{array}{l} \mod 16 \\ \mod 8 \end{array} \right\}$
--	---	---	--	--

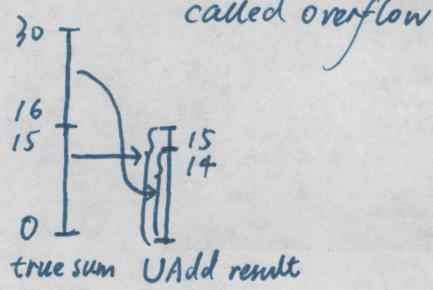
unsigned addition $W=4$ (to get the intuition) 0...15

$$\begin{array}{r}
 1101 \\
 0101 \\
 \hline
 \times 0010 \quad 18 \rightarrow 2
 \end{array}$$

just discard
without warnings
messages, errors

overflow, differs by 16
4 bits can't represent 18

W bits can represent up to $2^W - 1$ for unsigned int, so the sum of two ints would be up to $(2^W - 1) * 2 = 2^{W+1} - 2$, for sum larger than or equal to 2^W , we subtract it by 2^W . e.g. if $W=4$



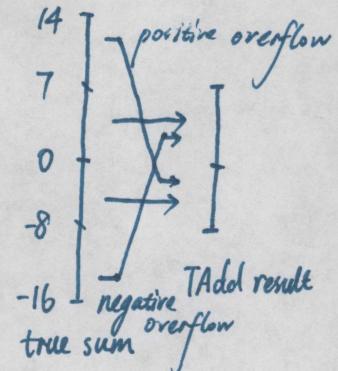
signed addition $W=4$ -8...7

$$\begin{array}{r}
 1101 -3 \\
 0101 5 \\
 \hline
 \times 0010 \quad 2
 \end{array}$$

discard carry

$$\begin{array}{r}
 1011 -5 \\
 0011 3 \\
 \hline
 1110 -2
 \end{array}$$

OK



$$\begin{array}{r}
 1101 -3 \\
 1010 -6 \\
 \hline
 \times 0111 +7 \\
 -9
 \end{array}$$

negative overflow
differs by 16
4 bits can't represent
-9

$$\begin{array}{r}
 0111 7 \\
 0101 5 \\
 \hline
 1100 -4
 \end{array}$$

positive overflow
differs by 16
signed bit represents -8,
not +8

unsigned multiplication

$$\begin{array}{r}
 0011 \quad 3 \\
 \times 0101 \quad 5 \\
 \hline
 111101 \quad 15 \text{ OK}
 \end{array}$$

signed multiplication

$$\begin{array}{r}
 0101 \quad 5 \\
 \times 0101 \quad 5 \\
 \hline
 \cancel{0} \cancel{0} 011001 \quad 25 \rightarrow 9 = 25 \% 16
 \end{array}$$

5*4

$$\begin{array}{r}
 \cancel{0} \cancel{0} 110100 \\
 \times 4 \\
 \hline
 20 \rightarrow
 \end{array}$$

5*5

$$\begin{array}{r}
 \cancel{0} \cancel{0} 011001 \\
 \times -7 \\
 \hline
 25 \rightarrow
 \end{array}$$

$$1101 \quad -3$$

$$\begin{array}{r}
 \times 1110 \quad -2 \\
 \hline
 \cancel{0} \cancel{1} 110 \quad 6 \text{ OK}
 \end{array}$$

trick: $-3 \rightarrow 13u$
 $-2 \rightarrow 14u$

$$13 \times 14 = 182.66$$

Power-of-2 Multiply with Shift for both signed & unsigned

$$x = \sum_{i=0}^{w-1} 2^i$$

$$x \cdot 2^k = 2^k \cdot \sum_{i=0}^{w-1} 2^i = \sum_{i=0}^{w-1} 2^{i+k}$$

$w \rightarrow w+k \rightarrow$ discard k bits

unsigned Power-of-2 Divide with Shift

$$\begin{array}{r}
 0110 \quad 6 \\
 \gg 1 \text{ logical shift} \\
 0011 \quad 3 \\
 \gg 1 \quad 1.5 \rightarrow 1 \\
 0001 \quad 1
 \end{array}$$

signed Power-of-2 Divide

$$1010 \quad -6$$

$\gg 1$ arithmetic shift

$$1101 \quad -3$$

$$\begin{array}{r}
 \gg 1 \\
 1110 \quad -2 \\
 \hline
 1110 \\
 \gg 1 \\
 1111 \quad -1
 \end{array}$$

so we should add
a bias ↑
computer doing this

Why should we use unsigned? Java banished unsigned numbers
and only use 2's complement to represent numbers

Should do explicit cast ^{as much}
^{as possible}

Rule of Thumb

$$2^{10} = 1024 \sim 10^3$$

$$2^{20} \sim 10^6$$

$$2^{30} \sim 10^9$$

$$2^{40} \sim 10^{12}$$

$$2^{47} \sim 128\text{TB} \quad (\text{TB} = 10^{12})$$

approximation

from programmer's perspective, memory is just a big array of bytes \rightarrow virtual memory for 64 bit machine, an address is represented in 64 bits. But the maximum address you're allowed to use is 47 bits. When you run a program the operating system only allows certain region of memory to be referenced. Access other regions will create segmentation fault error. Idea of big array of bytes is taken care by the operating system and computer hardware, invisible to programmers, even machine-level.

gcc {
-m32 32 bit code
-m64 64 bit code

the hardware itself doesn't necessarily define what the word size is. it's the combination of hardware and compiler to determine what the word size to use in the program.

most significant ← least significant
两个数字: 0x12345678, 0x11223344

0x78 | 0x56 | 0x34 | 0x12 | 0x44 | 0x33 | 0x22 | 0x11
0x12 | 0x34 | 0x56 | 0x78 | 0x11 | 0x22 | 0x33 | 0x44

little endian x86, ARM processors
least significant byte has lowest address
内存地址增长方向 变量指针转换时地址保持不变
big endian Sun, PPC Mac, Internet (TCP/IP)
least significant byte has highest address

内存顺序和书写顺序一致
for negative numbers so sign bit first in big endian machine

puzzles

if $x < 0$, then $x * 2$ may not be negative
example: $2 * T_{\text{Min}} = 0$

casting it to unsigned $ux > 0$

if $x \& 7 = 7$ ((are 3 bits are all 1s))

then $x \ll 30$ 00...0111

11...000 \rightarrow negative

if $x > y$ $-x < -y$: false

negative T_{Min} is still T_{Min}

$$\begin{array}{r} 100\dots0 \\ \sim 011\dots1 \\ +1100\dots0 \end{array}$$

$x * x \geq 0$: false

if $x > 0 \& y > 0$ $x + y > 0$: false

if $x > 0$ $-x \leq 0$: true

if $x \leq 0$ $-x > 0$: false T_{Min} . Any positive number can be represented as negative, but T_{Min} as a negative number cannot be represented as positive

$ux \gg 3 == ux / 8$: true

$x \gg 3 == x / 8$: false -1 right shift 3 bits is still -1

$x \& (x-1) != 0$: false $T_{\text{Min}} - 1 = T_{\text{Max}}$ so $T_{\text{Min}} \& (T_{\text{Min}} - 1) == 0$
get rid of the rightmost 1 in x

$ux > -1$: false signed integer compared with unsigned integer,
promoted to unsigned
(unsigned int) - 1 = UMax