# CMPT 295

Unit - Data Representation

Lecture 4 – Representing integral numbers in memory – Arithmetic operations

# Warm up question

➡ What is the value of …

   ➡ `TMin` (in hex) for `signed char` in C: _____

   ➡ `TMax` (in hex) for `signed int` in C: _____

   ➡ `TMin` (in hex) for `signed short` in C: _____

# Last Lecture

- Interpretation of bit pattern B into either unsigned value U or signed value T
  - *B2U(X)* and *U2B(X)* encoding schemes (conversion)
  - *B2T(X)* and *T2B(X)* encoding schemes (conversion)
    - Signed value expressed as two's complement => *T*
- Conversions from unsigned <-> signed values
  - *U2T(X)* and *T2U(X)* => adding or subtracting $2^w$
- Implication in C: when converting (implicitly via promotion and explicitly via casting):
  - Sign:
    - Unsigned <-> signed (of same size) -> Both have same bit pattern, however, this bit pattern may be interpreted differently
      - Can have unexpected effects -> producing a different value
  - Size:
    - Small -> large (for signed, e.g., `short` to `int` and for unsigned, e.g., `unsigned short` to `unsigned int`)
      - sign extension: For unsigned -> zeros extension, for signed -> sign bit extension
      - Both yield expected result –> resulting value unchanged
    - Large -> small (for signed, e.g., `int` to `short` and for unsigned, e.g., `unsigned int` to `unsigned short`)
      - truncation: Unsigned/signed -> most significant bits are truncated (discarded)
      - May not yield expected results -> original value may be altered
  - Both (sign and size): 1) **size** conversion is first done then 2) **sign** conversion is done

# Today's Menu

- Representing data in memory – Most of this is review
  - "Under the Hood" - Von Neumann architecture
  - Bits and bytes in memory
    - How to diagram memory -> Used in this course and other references
    - How to represent series of bits -> In binary, in hexadecimal (conversion)
    - What kind of information (data) do series of bits represent -> Encoding scheme
    - Order of bytes in memory -> Endian
  - Bit manipulation – bitwise operations
    - Boolean algebra + Shifting
- Representing integral numbers in memory
  - Unsigned and signed
  - Converting, expanding and truncating
  - Arithmetic operations
- Representing real numbers in memory
  - IEEE floating point representation
  - Floating point in C – casting, rounding, addition, …

Let's first illustrate what we covered last lecture with a demo!

# Demo – Looking at *size* and *sign* conversions in C

- What does the demo illustrate?
  - Size conversion:
    - Converting to a larger (wider) data type -> Converting `short` to `int`
    - Converting to a smaller (narrower) data type -> Converting `short` to `char`
  - Sign conversion:
    - Converting from signed to unsigned -> Converting `short` to `unsigned short`
    - Converting from unsigned to signed -> Converting `unsigned short` to `short`
  - Size and Sign conversion:
    - Converting from signed to unsigned larger (wider) data type -> Converting `short` to `unsigned int`
    - Converting from signed to unsigned smaller (narrower) data type -> Converting `short` to `unsigned char`
- This demo (code and results) posted on our course web site

# Integer addition (unlimited space)

- What happens when we add two decimal numbers?

$$\begin{array}{r} 1 \quad \text{<- carry in} \\ 107_{10} \\ +\ 938_{10} \\ \hline \text{carry out ->}\ 1045_{10} \end{array}$$

- Same thing happens when we add two binary numbers:

$$\begin{array}{r} 1\ 1 \quad \text{<- carry in} \\ 101100_2 \\ +\ \ 101110_2 \\ \hline \text{carry out ->}\ 1011010_2 \end{array}$$

# Unsigned addition (limited space, i.e., fixed size in memory)

➧ What happens when we add two unsigned values:

$w = 8$
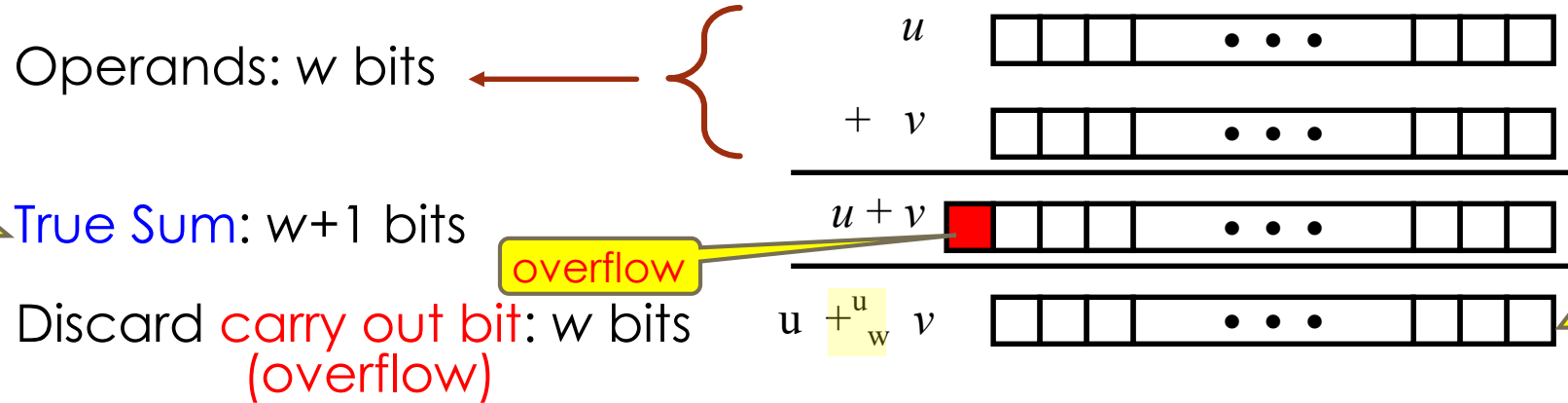
a)    $00111011_2$      $59_{10}$
       $+ \, 01011010_{\,2}$    $+ \, 90_{10}$
                        $149_{10}$

b)    $10101110_2$      $174_{10}$
       $+ \, 11001011_2$    $+ \, 203_{10}$
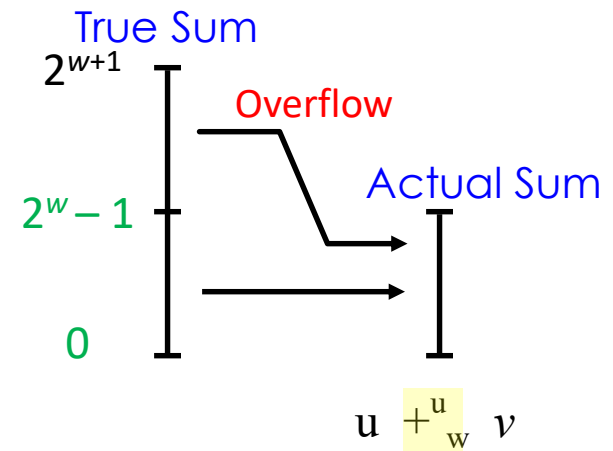                        $377_{10}$

# Unsigned addition ($+^u_w$) and overflow

Would be the result of **integer addition** with unlimited space: expected sum

Result of **unsigned addition** with limited space: actual sum

Operands: $w$ bits

$u$

$+\ v$

True Sum: $w+1$ bits

$u + v$

overflow

Discard carry out bit: $w$ bits (overflow)

$u +^u_w v$

- Discarding carry out bit has same effect as applying modular arithmetic

$$s\ =\ u +^u_w v\ =\ (u + v)\ \text{mod}\ 2^w$$

True Sum

$2^{w+1}$

Overflow

Actual Sum

$2^w - 1$

0

$u +^u_w v$

# Closer look at unsigned addition overflow

w = 8 -> [0..255]

$255_{10} = 11111111_2$
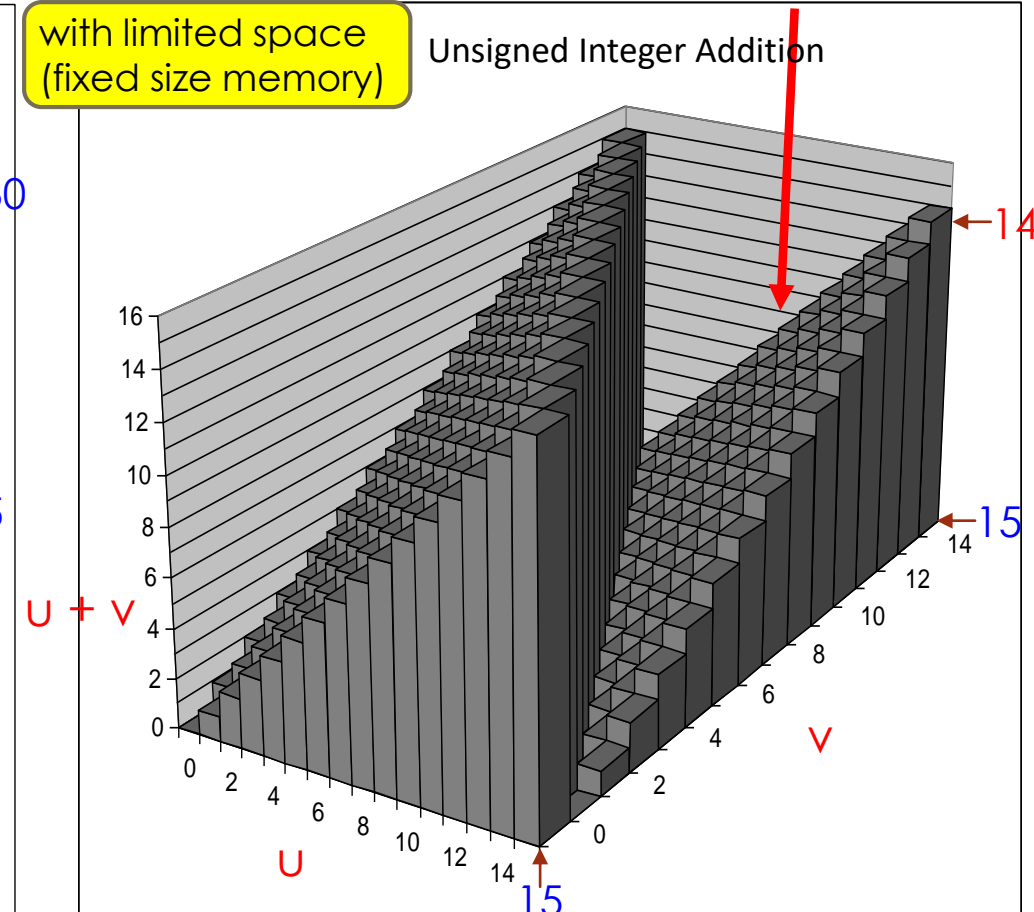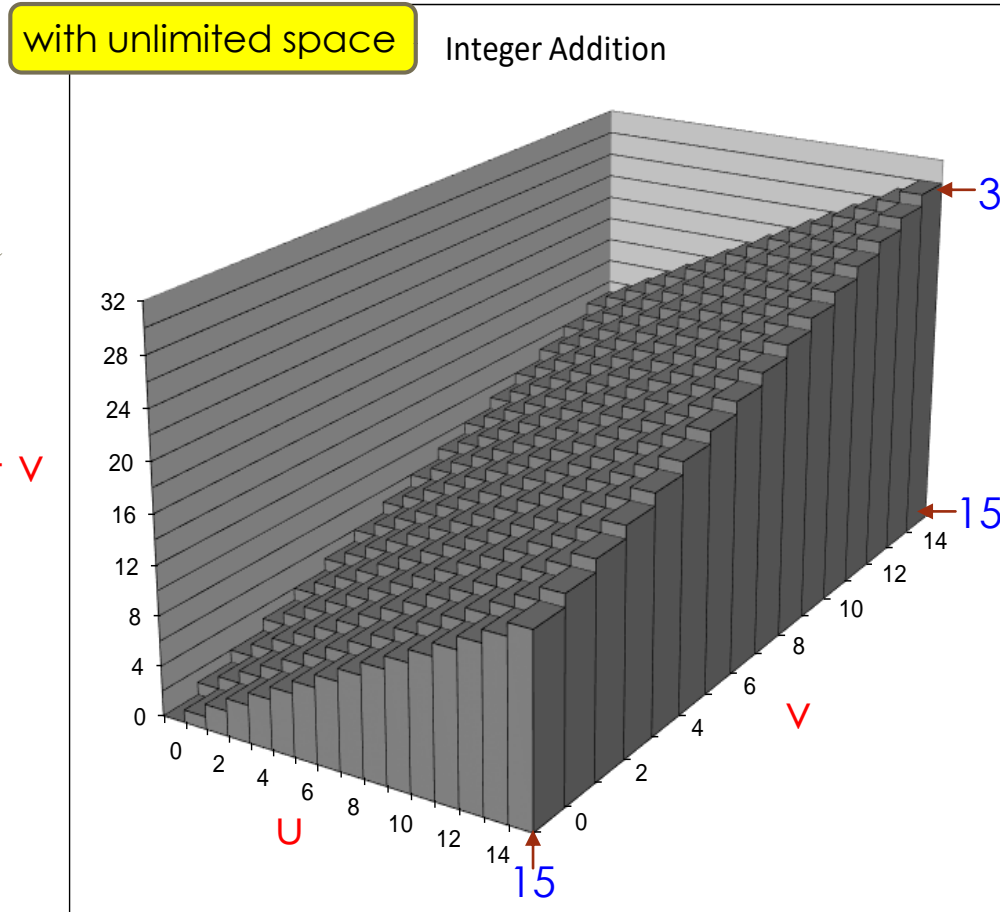
$90_{10} = 01011010_2$

$45_{10} = 00101101_2$

carry in

$90_{10}$      $01011010_2$

$+ \; 45_{10}$      $+ \; 00101101_2$

$135_{10}$      $10000111_2$

carry out

carry in

$255_{10}$      $11111111_2$

$+ \; 45_{10}$      $+ \; 00101101_2$

$300_{10}$      $100101100_2$

carry out

True Sum

511

Overflow

w = 9

300

Actual Sum

255

135

range where w = 8

0

44

9

# Comparing integer addition with unsigned addition (w = 4)

**Overflow**: Effect of fixed size memory

with unlimited space — Integer Addition

with limited space (fixed size memory) — Unsigned Integer Addition



An overflow occurs when there is a <u>carry out</u>

For example: 15 (1111₂) + 15 (1111₂) = 30 (11110₂ <- true sum) and = 14 (11110₂ <- actual sum)

# Signed addition (limited space, i.e., fixed size in memory)

➡ What happens when we add two signed values:

$w = 8$

a) $\quad 00111011_2 \qquad 59_{10}$
$\quad\quad + 01011010_2 \qquad + 90_{10}$
$\quad\quad\overline{\phantom{+ 01011010_2}} \qquad \overline{\phantom{+ 90_{10}}}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad 149_{10}$

b) $\quad 10101110_2 \qquad -82_{10}$
$\quad\quad + 11001011_2 \qquad + \text{-}53_{10}$
$\quad\quad\overline{\phantom{+ 11001011_2}} \qquad \overline{\phantom{+ \text{-}53_{10}}}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad \text{-}135_{10}$

Observation: Unsigned and signed additions have identical behavior @ the bit level, i.e., their sum have the same bit-level representation, but their interpretation differs
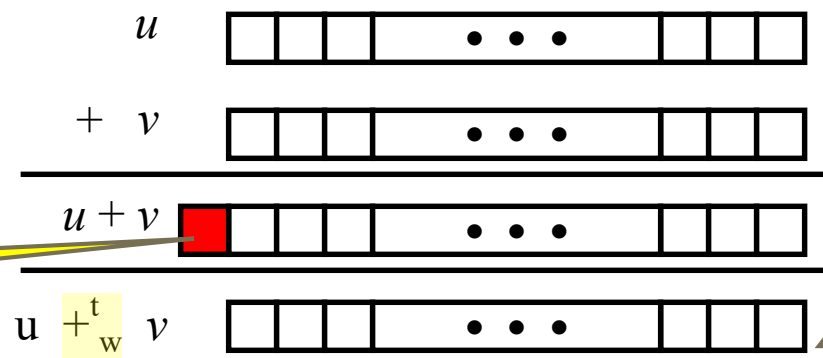
# Signed addition ($+^t_w$) and overflow

Would be the result of **integer addition** with unlimited space: expected sum
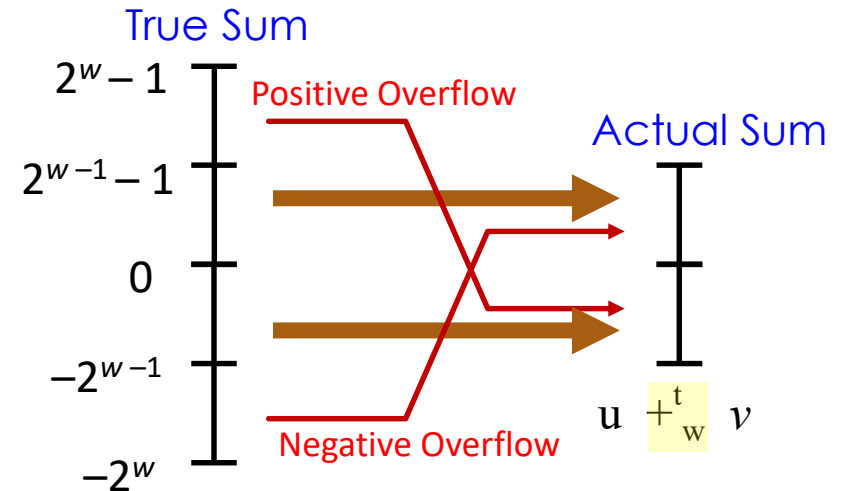
Operands: $w$ bits

True Sum: $w+1$ bits

overflow

Discard carry out bit : $w$ bits (overflow)

Result of **signed addition** with limited space: actual sum

$u$

$+\ v$

$u + v$

$u\ +^t_w\ v$

- Discarding carry out bit has same effect as applying modular arithmetic

$$s\ =\ u\ +^t_w\ v\ =\ U2T_w\,[(u + v)\ \text{mod}\ 2^w]$$

True Sum

$2^w - 1$

Positive Overflow

Actual Sum

$2^{w-1} - 1$

$0$

$-2^{w-1}$

Negative Overflow

$u\ +^t_w\ v$
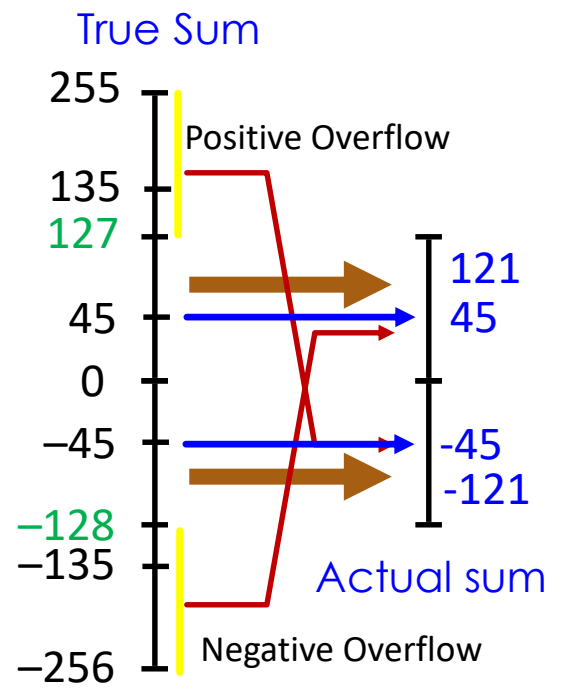
$-2^w$

# Closer look at signed addition overflow

w = 8 -> [-128..127]

$90_{10} = 01011010_2$

$45_{10} = 00101101_2$

$-45_{10} = 11010011_2$

$- 90_{10} = 10100110_2$

carry in

$90_{10}$     $01011010_2$

$+ \ 45_{10}$     $+ \ 00101101_2$

$135_{10}$     $010000111_2$   <= -121

carry out

carry in

$-90_{10}$     $10100110_2$

$+ \ -45_{10}$     $+ \ 11010011_2$

$-135_{10}$     $101111001_2$   <= 121

carry out

carry in

$-90_{10}$     $10100110_2$

$+ \ 45_{10}$     $+ \ 00101101_2$

$-45_{10}$     $011010011_2$

carry out

carry in

$90_{10}$     $01011010_2$

$+ \ -45_{10}$     $+ \ 11010011_2$

$45_{10}$     $100101101_2$

carry out

True Sum

Positive Overflow

Actual sum

Negative Overflow

255

135

127

45

0

−45

−128

−135

−256

121

45

-45

-121

# Visualizing signed addition overflow ($w = 4$)

For example: 7 ($0111_2$) + 1 ($0001_2$) = 8 ($1000_2$ <- true sum) and = -8 ($1000_2$ <- actual sum)

# What about subtraction? -> Addition

$x + (-x) = 0$

- Subtracting a number is equivalent to adding its additive inverse
  - Instead of subtracting a positive number, we could add its negative version:

$$107 \qquad\qquad 107$$
$$\underline{-\ 118} \qquad => \qquad \underline{+\ (-118)}$$
$$-\ 11$$

- Let 's try:    $107_{10}$  ->    $01101011_2$  ->   $01101011_2$

$$\underline{-\ 118_{10}} \ -> \underline{-\ 01110110_2} \ -> \underline{+\ 10001010_2}$$
$$-\ 11 \qquad\qquad\qquad\qquad\qquad 11110101_2 => -11_{10}$$

$T2B(X)$ conversion: **$(\sim(U2B(|X|)))+1$**
**$= (\sim(U2B(|-118|)))+1$**
**$= (\sim(U2B(118)))+1$**
**$= (\sim(01110110_2))+1$**
**$= (10001001_2)+1$**
**$= 10001010_2$**

CHECK: $-128+64+32+16+4+1 = -11_{10}$

15

# Multiplication ($*^u_w$ , $*^t_w$) and overflow



Operands: $w$ bits

True Product: $2w$ bits

Discard: $w$ bits

▶ Discarding high order $w$ bits has same effect as applying modular arithmetic

$$p = u *^u_w v = (u * v) \mod 2^w$$

$$p = u *^t_w v = U2T_w [(u * v) \mod 2^w]$$

▶ Example: $w = 4$

$$
\begin{array}{r}
5_{10} \\
\times\ 5_{10} \\
\hline
25_{10}
\end{array}
$$

$$
\begin{array}{r}
0101_2 \\
\times\ 0101_2 \\
\hline
0101_2 \\
0000_2 \\
0101_2 \\
0000_2 \\
\hline
001\!\!\!\!\!/\,1001_2
\end{array}
$$

16

# Multiplication with power-of-2 versus shifting

- If $x * y$ where $y = 2^k$ then $x << k$
  - For both signed and unsigned

- Example:
  - $x * 8 = x * 2^3 \rightarrow x << 3$
  - $x * 24 = (x * 2^5) - (x * 2^3) = (x * 32) - (x * 8) \rightarrow (x << 5) - (x << 3)$
    (decompose 24 in powers of 2 => 32 − 8)

- Most machines shift and add faster than multiply
  - We'll soon see that compiler generates this code automatically

# Summary

- Demo of size and sign conversion in C: code and results posted!
- Addition:
  - Unsigned/signed:
    - Behave the same way at the bit level
    - Interpretation of resulting bit vector (sum) may differ
  - Unsigned addition -> may overflow, i.e., $(w+1)^{th}$ bit is set
    - If so, then actual sum obtained => $(x + y)$ mod $2^w$
  - Signed addition -> may overflow, i.e., $(w+1)^{th}$ bit is set
    - If so, then true sum may be too +ve -> positive overflow OR too –ve -> negative overflow
    - Then actual sum obtained => $U2T_w [(x + y)$ mod $2^w]$
- Subtraction
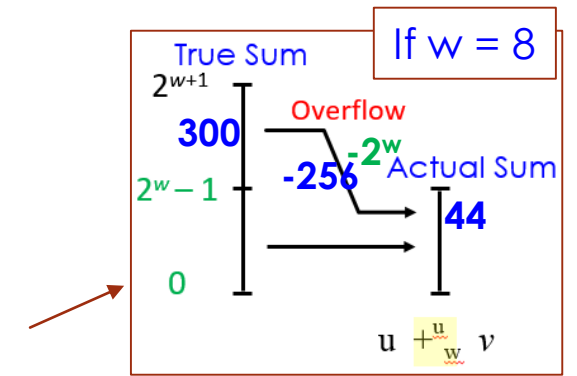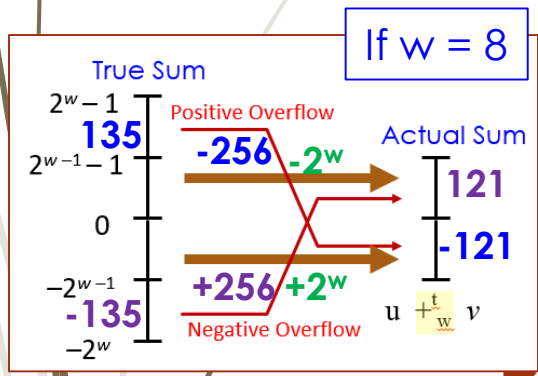  - Becomes an addition where negative operands are transformed into their additive inverse (in two's complement)
- Multiplication:
  - Unsigned: actual product obtained -> $(x * y)$ mod $2^w$
  - Signed: actual product obtained -> $U2T_w [(x * y)$ mod $2^w]$
  - Can be replaced by additions and shifts



If w = 8

True Sum

300

Overflow
$-2^w$
Actual Sum
-256

44

$2^{w+1}$
$2^w - 1$
0

$u \ +_w^u \ v$



If w = 8

True Sum

Positive Overflow
135
-256 $-2^w$
Actual Sum
121
-121
+256 $+2^w$
-135
Negative Overflow

$2^w - 1$
$2^{w-1} - 1$
0
$-2^{w-1}$
$-2^w$

$u \ +_w^t \ v$

18

# Next lecture

- Representing data in memory – Most of this is review
  - "Under the Hood" - Von Neumann architecture
  - Bits and bytes in memory
    - How to diagram memory -> Used in this course and other references
    - How to represent series of bits -> In binary, in hexadecimal (conversion)
    - What kind of information (data) do series of bits represent -> Encoding scheme
    - Order of bytes in memory -> Endian
  - Bit manipulation – bitwise operations
    - Boolean algebra + Shifting
- Representing integral numbers in memory
  - Unsigned and signed
  - Converting, expanding and truncating
  - Arithmetic operations
- Representing real numbers in memory
  - IEEE floating point representation
  - Floating point in C – casting, rounding, addition, …

We'll illustrate what we covered today by having a demo!