# CS-200
# Computer Architecture
# —
## Part Id: Instruction Set Architecture Arithmetic

Paolo Ienne

<paolo.ienne@epfl.ch>

# Notation

- Number (represented on a specific no. of digits/bits)

$$A = A^{(n)} = A^{(m)}$$

- Number (in binary or decimal)

$$A = A_{10} = A_2 = A_{2c}$$

Binary, 2's complement

Binary

- Individual digits (bits)

$$a_{n-1}, \ a_{n-2}, \dots \ a_2, \ a_1, \ a_0$$

Simply 100010
if the digits are known

- Digit string (representation)

$$\langle a_{n-1} a_{n-2} \dots a_2 a_1 a_0 \rangle$$

# Numbers

We usually care for three types of numbers:

- **Integers** (signed and unsigned)

$$0, 1, 2, 3, 4294967295, -2147483648$$

- **Fixed Point**

$$0.\underline{12}, 3.\underline{14}, 1073741823.\underline{75}$$

  – Essentially integers with **implicit $10^k$ or $2^k$ scaling**

  – Extremely important in practice (most signal-processing is fixed point)

- **Floating Point**

$$3.14E3, -2.5E1, 1.0E0, 4.2E-2, -1.5E-3$$

# Unsigned Integers

- Weighted (positional)

- Nonredundant

- Fixed-radix (radix-10 or radix-2)

- Canonical


- Definition:

If R = 2, binary

$$A = \langle a_{n-1} a_{n-2} \ldots a_2 a_1 a_0 \rangle = \sum_{i=0}^{n-1} a_i R^i$$

# Signed Integers

- **Sign-and-Magnitude**

- **2's Complement** (particular choice of True-and-Complement)

- **Biased**
  – Practically used only in Floating Point numbers (mentioned later)

# Sign and Magnitude

- **Human friendly!**
- The first symbol is a sign (+/− for humans, 0/1 for computers)
- The rest is an unsigned number:

$$+100, -2345$$

$$+111_2 = 0111_2^{(4)}$$

$$\boxed{-111_2 = 1111_2^{(4)}}$$

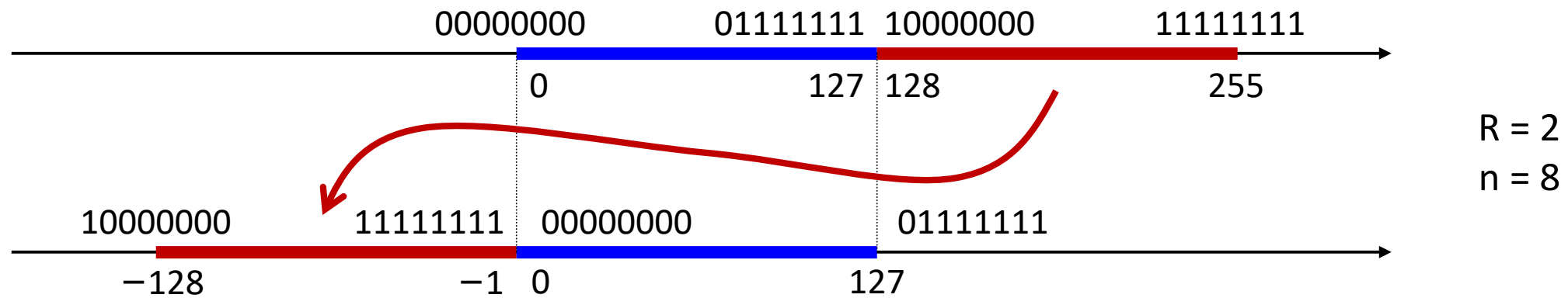If we use 0/1 for the sign, the number of bits matters

- Definition:

If R = 2, binary

$$A = \langle s a_{n-2} \dots a_2 a_1 a_0 \rangle = (-1)^s \cdot \sum_{i=0}^{n-2} a_i R^i$$

0 or 1

# Radix's Complement

- Special form of **True-and-Complement** with C = R<sup>n</sup>



R = 2
n = 8

- Definition:

If R = 2, binary

$$A = \langle a_{n-1} a_{n-2} \ldots a_2 a_1 a_0 \rangle = -a_{n-1} R^{n-1} + \sum_{i=0}^{n-2} a_i R^i$$

# Radix's Complement

- **Not** a **human-friendly** representation
- In **decimal** (10's complement):

$$5,678_{10c}^{(5)} = 05,678_{10c} = +5,678_{10}$$

$$9,999,999_{10c}^{(7)} = -1_{10}$$

$$8,766_{10c}^{(4)} = -1,234_{10}$$

- In **binary** (2's complement):

$$0100,1101,0010_{2c}^{(12)} = 100,1101,0010_2 = +1,234_{10}$$

$$1111,1111_{2c}^{(8)} = -1_{10}$$

$$1011,0000,1110_{2c}^{(12)} = -1,234_{10}$$

# 2's Complement from Subtraction

- Consider a "normal" **paper-and-pencil subtraction**

$$
\begin{array}{cccccccc}
& 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0_2 & \quad 10_{10} \\
- & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1_2 & \quad 17_{10} \\
\hline
\end{array}
$$

# 2's Complement from Subtraction

- Consider a "normal" **paper-and-pencil subtraction**

$$
\begin{array}{ccccccccc}
-1 & -1 & -1 & & & & -1 & & \\
0 & 0 & 0 & 0 & 1 & 0 & 1 & 0_2 & 10_{10} \\
- \; 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1_2 & 17_{10} \\
\hline
\ldots & \ldots & 1 & 1 & 1 & 1 & 0 & 0 & 1_2 \\
\end{array}
$$

$\downarrow$

Stop and "accept" the −1…

$$
\begin{array}{ccccccccc}
-1 & 1 & 1 & 1 & 1 & 0 & 0 & 1_2 \\
-2^7 & +2^6 & +2^5 & +2^4 & +2^3 & & & +2^0 & -7_{10} \\
\end{array}
$$

A sign bit

# Addition Is Unchanged from Unsigned

- Only two instructions (with the immediate version; `subi` is a pseudo)

| Arithmetic | | | | | | |
|---|---|---|---|---|---|---|
| add | rd,rs1,rs2 | $rd \leftarrow rs1 + rs2$ | R | 0x00 | 0x0 | 0x33 |
| addi | rd,rs1,imm | $rd \leftarrow rs1 + \text{sext}(imm)$ | I | | 0x0 | 0x13 |
| sub | rd,rs1,rs2 | $rd \leftarrow rs1 - rs2$ | R | 0x20 | 0x0 | 0x33 |

- Old architectures (MIPS, notably) had distinct `add` and `addu` but it was essentially a misnomer; **ignore** it and do not be confused!

- Instead, addition of Sign-and-Magnitude numbers is a different problem (see later) → this is why **2's complement is the universal representation** of signed integers today
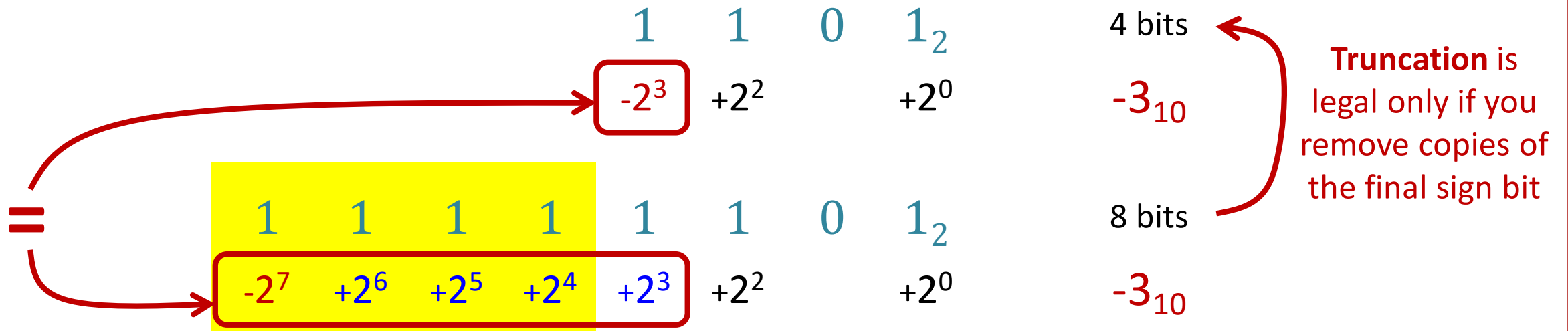
# Sign Extension

- **Unsigned numbers** can be though as having infinite 0s in front

$$-1_{10} = -0001_{10}$$

$$1,0101_2 = 0000,0000,0001,0101_2$$

- Instead, **2's complement numbers** have infinite replicas of the MSB/sign bit in front

# Instructions for Signed Numbers

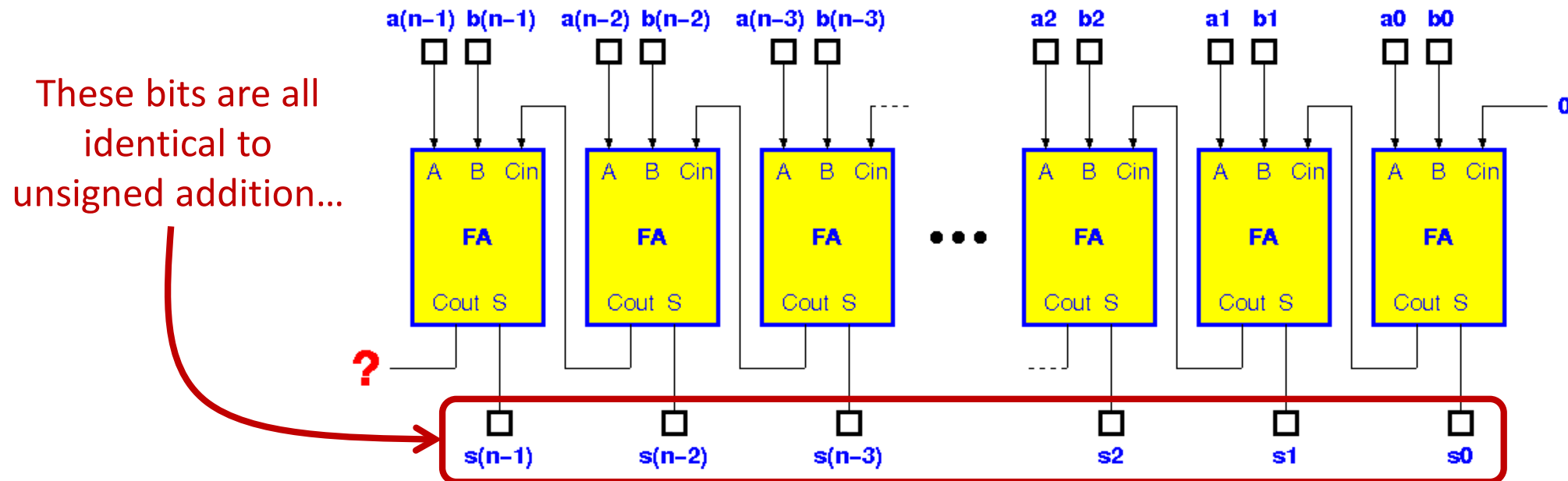Insert zeroes ($1$ = logic → **unsigned**) or sign bits (**a** = arithmetic → **signed**)

**Shift**

| | | | | | | |
|---|---|---|---|---|---|---|
| srl | rd,rs1,rs2 | $\text{rd} \leftarrow \text{rs1} \gg_u \text{rs2}$ | R | 0x00 | 0x5 | 0x33 |
| srli | rd,rs1,imm | $\text{rd} \leftarrow \text{rs1} \gg_u \text{imm}$ | I | 0x00 | 0x5 | 0x13 |
| sra | rd,rs1,rs2 | $\text{rd} \leftarrow \text{rs1} \gg_s \text{rs2}$ | R | 0x20 | 0x5 | 0x33 |
| srai | rd,rs1,imm | $\text{rd} \leftarrow \text{rs1} \gg_s \text{imm}$ | I | 0x20 | 0x5 | 0x13 |

$$1110_2/2 = 0111_2$$
$$\text{but}$$
$$1110_{2c}/2 = 1111_{2c}$$

**Compare**

| | | | | | | |
|---|---|---|---|---|---|---|
| slt | rd,rs1,rs2 | $\text{rd} \leftarrow \text{rs1} <_s \text{rs2}$ | R | 0x00 | 0x2 | 0x33 |
| slti | rd,rs1,imm | $\text{rd} \leftarrow \text{rs1} <_s \text{sext}(\text{imm})$ | I | | 0x2 | 0x13 |
| sltu | rd,rs1,rs2 | $\text{rd} \leftarrow \text{rs1} <_u \text{rs2}$ | R | 0x00 | 0x3 | 0x33 |
| sltiu | rd,rs1,imm | $\text{rd} \leftarrow \text{rs1} <_u \text{sext}(\text{imm})$ | I | | 0x3 | 0x13 |

$$0000_2 < 1111_2$$
$$\text{but}$$
$$0000_{2c} > 1111_{2c}$$

**Branch**

| | | | | | | |
|---|---|---|---|---|---|---|
| blt | rs1,rs2,imm | $\text{pc} \leftarrow \text{pc} + \text{sext}(\text{imm} \ll 1),\ \text{if rs1} <_s \text{rs2}$ | B | | 0x4 | 0x63 |
| bge | rs1,rs2,imm | $\text{pc} \leftarrow \text{pc} + \text{sext}(\text{imm} \ll 1),\ \text{if rs1} \geq_s \text{rs2}$ | B | | 0x5 | 0x63 |
| bltu | rs1,rs2,imm | $\text{pc} \leftarrow \text{pc} + \text{sext}(\text{imm} \ll 1),\ \text{if rs1} <_u \text{rs2}$ | B | | 0x6 | 0x63 |
| bgeu | rs1,rs2,imm | $\text{pc} \leftarrow \text{pc} + \text{sext}(\text{imm} \ll 1),\ \text{if rs1} \geq_u \text{rs2}$ | B | | 0x7 | 0x63 |

**Load**

| | | | | | | |
|---|---|---|---|---|---|---|
| lb | rd,imm(rs1) | $\text{rd} \leftarrow \text{sext}(\text{mem}[\text{rs1} + \text{sext}(\text{imm})][7:0])$ | I | | 0x0 | 0x03 |
| lbu | rd,imm(rs1) | $\text{rd} \leftarrow \text{zext}(\text{mem}[\text{rs1} + \text{sext}(\text{imm})][7:0])$ | I | | 0x4 | 0x03 |
| lh | rd,imm(rs1) | $\text{rd} \leftarrow \text{sext}(\text{mem}[\text{rs1} + \text{sext}(\text{imm})][15:0])$ | I | | 0x1 | 0x03 |
| lhu | rd,imm(rs1) | $\text{rd} \leftarrow \text{zext}(\text{mem}[\text{rs1} + \text{sext}(\text{imm})][15:0])$ | I | | 0x5 | 0x03 |

# Overflows in 2's Complement Addition

- The **sum** is the same as with **unsigned numbers**:

These bits are all identical to unsigned addition…



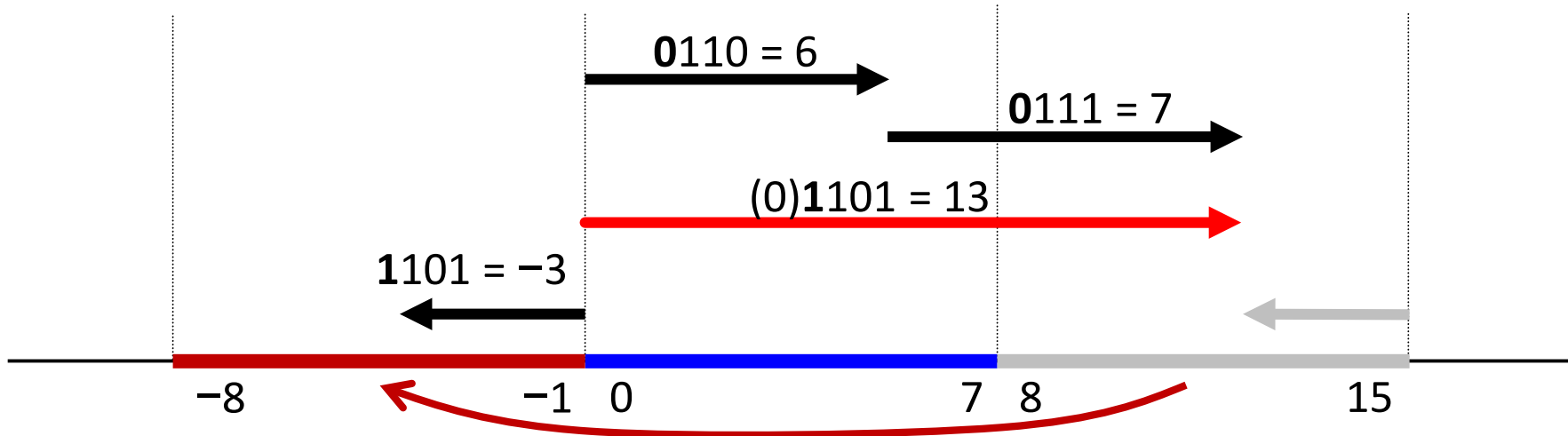…but how to assess **overflows**?

# Overflows in Hardware

- In hardware, **carry out** is the only missing bit from the **complete** result
- We can think of overflows as a **truncation** problem:

For unsigned numbers, the carry bit must be **zero**

| | | |
|---|---|---|
| **No** | 1 | ... |
| **Ok** | 1 | 1 |
| **No** | 1 | 0 |

In 2's complement, the carry bit must be **equal to the next bit**

# Overflow in Software

- Some architectures (e.g., **x86**) give us the **carry bit** in a special "register" (a **flag**)

  → overflow detection is the same as in hardware

- Other (modern) architectures give us **only the result** of the addition (e.g., **RISC-V**)

- Detection usually based on the following observations:

  - If addition of **opposite sign numbers**, magnitude can only reduce → **no overflow possible**
  - If addition of **same sign numbers**, overflow possible but the sign of the result will appear wrong

**0**110 = 6

**0**111 = 7

(0)**1**101 = 13

**1**101 = −3

−8        −1  0                7  8                15

# Detect Addition Overflow in Software

- Add two 32-bit signed integers **and detect overflow**
  - At call time, a0 and a1 contain the two integers
  - On return, a0 contains the result and a1 must be nonzero in case of overflow

# A + Ā = –1

- A "strange" but **very useful property**

$$\boxed{A + \bar{A} = -1} \quad \text{or} \quad -A = \bar{A} + 1$$
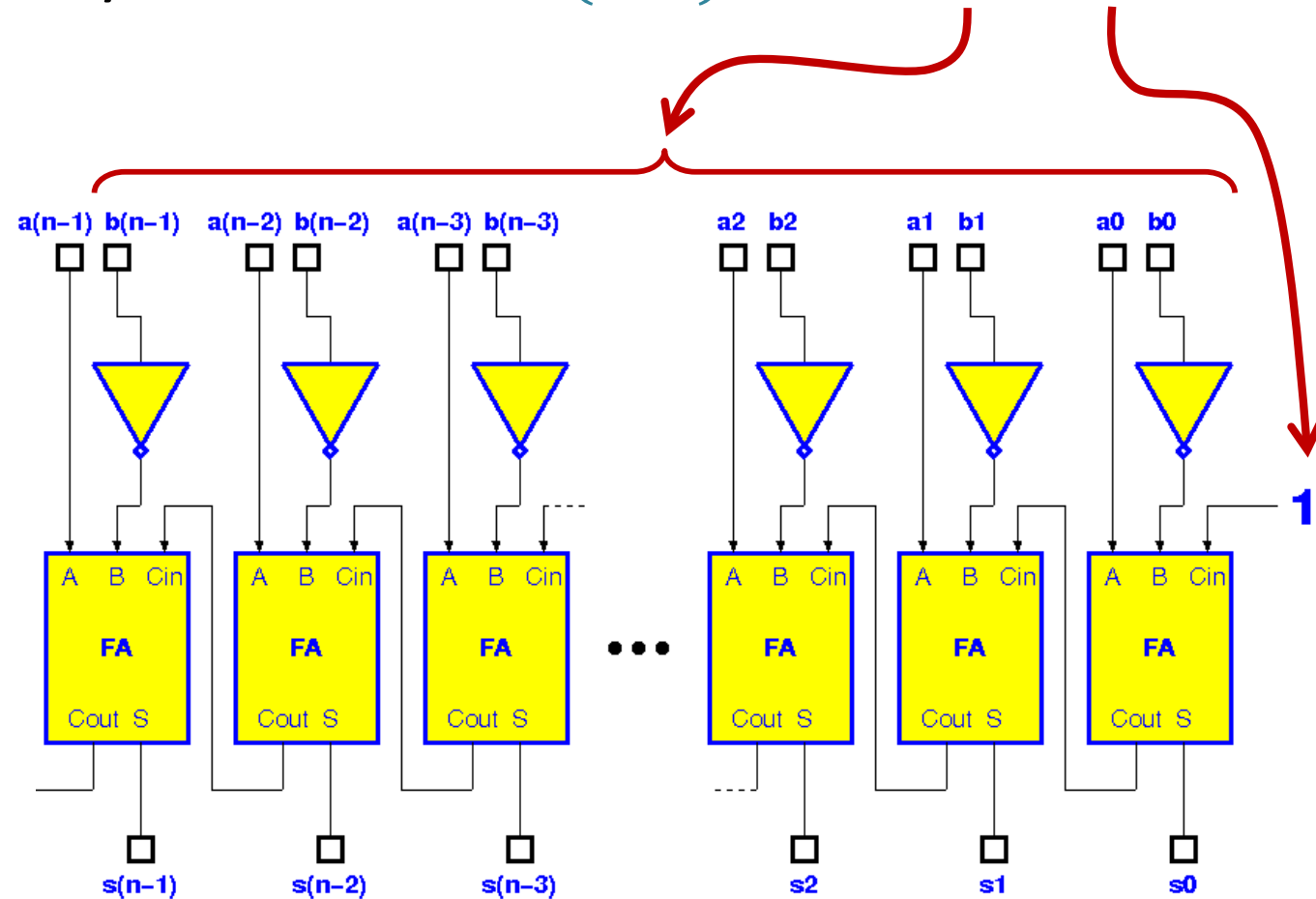
- Not too hard to prove

$$\left(-a_{n-1}2^{n-1} + \sum_{i=0}^{n-2} a_i 2^i\right) + \left(-\overline{a_{n-1}}2^{n-1} + \sum_{i=0}^{n-2} \overline{a_i} 2^i\right) =$$

$$= -(a_{n-1} + \overline{a_{n-1}}) \cdot 2^{n-1} + \sum_{i=0}^{n-2} (a_i + \overline{a_i}) \cdot 2^i = -2^{n-1} + \sum_{i=0}^{n-2} 2^i = -1$$

- Also somehow intuitive

$A$ → 0 1 0 0 1 1 0 0 +

$\bar{A}$ → 1 0 1 1 0 0 1 1 =

1 1 1 1 1 1 1 1

$-1$

# Two's Complement Subtractor

- Using this property, $A - B = A + (-B) = A + \overline{B} + 1$

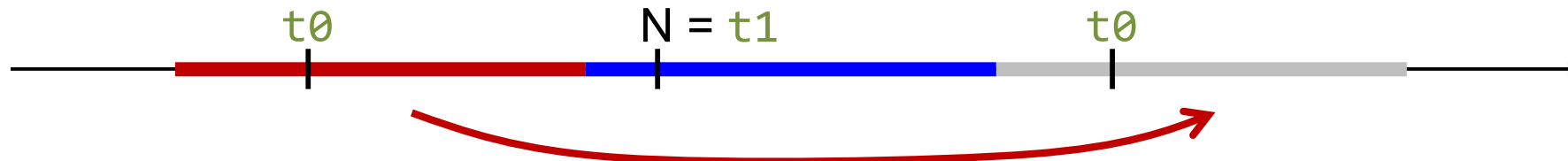# Two's Complement Add/Subtract Units

# Fun Stuff: Bounds Check

- Check for a signed number t0 (e.g., an array **index**) to be **within the bounds 0..N−1** where N is t1

$$\texttt{bgeu t0, t1, out\_of\_bound}$$

Unsigned!

- **Two checks with a single branch!**
  - If t0 ≥ 0, bgeu is like bge and the right behaviour is evident
  - If t0 < 0, as an unsigned t0 looks like larger than any signed positive

# Floating Point

- Corresponds to our **everyday habits**

| | | | | |
|---|---|---|---|---|
| .18 μm | → | $.18 \cdot 10^{-6}$ m | → | $1.8 \cdot 10^{-7}$ m |
| 75 km | → | $75 \cdot 10^{3}$ m | → | $7.5 \cdot 10^{4}$ m |
| 35 mm | → | $35 \cdot 10^{-3}$ m | → | $3.5 \cdot 10^{-2}$ m |
| 2.5 m | → | $2.5 \cdot 10^{0}$ m | → | $2.5 \cdot 10^{0}$ m |

- A significand (or **mantissa**) and an **exponent** of the base, for instance

2's complement exponent

$$X = \langle s a_{n-1} \ldots a_2 a_1 a_0 e_{m-1} \ldots e_1 e_0 \rangle = (-1)^s \cdot \sum_{i=0}^{n-1} a_i 2^i \cdot 2^{-e_{m-1} 2^{m-1} + \Sigma_{j=0}^{m-2} e_j 2^j}$$

Sign-and-Magnitude significand

# Floating Point

- Large **dynamic range** but **variable accuracy**

- Redundant unless **normalized**

- Not real numbers: **not associative!**

- Often exponent in **biased** signed representation
  - Zero can be represented by 0000…0000
  - Easier for comparisons and hardware implementations

- Often **normalized mantissa** $1 \leq m < 2$ with **hidden bit** (<u>1.</u>xxxxx)

- Today the **IEEE 754 standard** is almost universally adopted

- **x86/x64** supports FP through SSE/AVX extensions (since 1999)

- **RISC-V** supports FP through ISA extensions (not used in CS-200)

# Example
# Sign-and-Magnitude Addition

- Write a function in RISC-V assembler to sum two **32-bit signed numbers** represented in **sign-and-magnitude (S&M) format** and produce the result also in sign-and-magnitude format

- The two operands are in registers a0 and a1 on entry and the result should be placed in register a0

- Ignore overflows

...or think about them as an additional exercise

# References

- Patterson & Hennessy, COD – RISC-V Edition
  - **Chapter 2** and, in particular, **Section 2.4**
  - **Chapter 3** and, in particular, **Section 3.2**