

Title: Scientific Computation (PHYS 608) - Lecture 4

Date: Oct 29, 2009 10:30 AM

URL: <http://pirsa.org/09100180>

Abstract:

Binary Representations of Integers

Binary Representations of Integers

Decimal 6

1 byte = 8 bits

Binary Representation of Integers

Decimal 6

1 byte = 8 bits

0000 0 1 1 0

2^2 2^1 2^0

$$4 + 2 = 6$$

Signed Magnitude

Signed Magnitude

$$\textcircled{1}000 \quad 0110 = -6$$

sign bit

$$1000 \quad 0000 = -0$$

$$0000 \quad 0000 = +0$$

Signed Magnitude

$$\textcircled{1}000 \quad 0110 = -6$$

sign bit

$$1000 \quad 0000 = -0$$
$$0000 \quad 0000 = +0$$

One's Complement

$$0000 \quad 0110$$

Signed Magnitude

$$\textcircled{1}000 \quad 0110 = -6$$

sign bit

$$1000 \quad 0000 = -0$$

$$0000 \quad 0000 = +0$$

One's Complement

$$0000 \quad 0110$$

$$1111 \quad 1001$$

Signed Magnitude

$$\textcircled{1}0000 \quad 0110 = -6$$

sign bit

$$\begin{array}{l} 1000 \quad 0000 = -0 \\ 0000 \quad 0000 = +0 \end{array} \quad \left. \vphantom{\begin{array}{l} 1000 \\ 0000 \end{array}} \right\}$$

One's Complement

$$0000 \quad 0110 = 6$$

$$1111 \quad 1001 = -6$$

$$1111 \quad 1111$$

Signed Magnitude

$$\textcircled{1}0000 \quad 0110 = -6$$

sign bit

$$1000 \quad 0000 = -0$$
$$0000 \quad 0000 = +0$$

One's Complement

$$0000 \quad 0110 = 6$$

$$1111 \quad 1001 = -6$$

$$1111 \quad 1111$$

$$0000 \quad 0000$$
$$-1111 \quad 1111$$

Signed Magnitude

① 0000 0110 = -6
↓
sign bit

1000 0000 = -0
0000 0000 = +0

One's Complement

0000 0110 = 6
1111 1001 = -6

0000 0000
-1111 1111

Binary Representations of Integers

two's Complement

One's complement

0000 0110

add one

Binary Representations of Integers

two's Complement

One's complement

$$11111001 = -6$$

add one

$$00000001$$

$$11111010$$

Binary Representations of Integers

two's Complement

One's complement

$$11111001 = -6$$

add one

$$00000001$$

$$11111010 = -6 \text{ in two's complement}$$

0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1
 1

0 0 0 0 0 0 0 0

0
-0 in One's complement

→ 0 0 0 0 0 0 0 0

1 1 1 1 1 1 1 1

1

→ 0 0 0 0 0 0 0 0

0

-0 in one's complement

two's complement of 0

→ 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1

0
- 0 in one's complement

→ 0 0 0 0 0 0 0 0

two's complement of 0

0 1 1 1 1 1 1 1
1 0 0 0 0 0 1 0

127
- 126 in two's complement

→ 0 0 0 0 0 0 0 0
 1 1 1 1 1 1 1 1

0
 - 0 in one's complement

→ 0 0 0 0 0 0 0 0

two's complement of 0

0 1 1 1 1 1 1 1
 1 0 0 0 0 0 1 0

127
 - 126 in two's complement

→ 0 0 0 0 0 0 0 0
 1 1 1 1 1 1 1 1
 1

0
 - 0 in one's complement

→ 0 0 0 0 0 0 0 0

two's complement of 0

{ 0 1 1 1 1 1 1 1
 1 0 0 0 0 0 1 0

127
 - 126 in two's complement

0 0 0 0 0 0 0 1

Binary binary arithmetic

8-Bit Pattern

0000 0000

8-Bit Pattern

0000 0000

Two's

0

8-Bit Pattern

0000 0000

0000 0001

⋮

0111 1111

Two's

0

1

127

Unsigned

0

1

127

8-Bit Pattern

0000 0000

0000 0001

⋮

0111 1111

1000 0000

Two's

0

1

127

-128

Unsigned

0

1

127

128

8-Bit Pattern

0000 0000
0000 0001

Two's

0
1

Unsigned

0
1

1 1 1 1
0 0 0 0
0 0 0 1

127
-128
-127

127
128
129

8-Bit Pattern

Two's

Unsigned

0000 0000

0

0

0000 0001

1

1

⋮

0111 1111

127

127

1000 0000

-128

128

1000 0001

-127

129

1111 1111

-1

8-Bit Pattern

Two's

Unsigned

0000 0000

0

0

0000 0001

1

1

⋮

0111 1111

127

127

1000 0000

-128

128

1000 0001

-127

129

⋮

1111 1111

-1

255

Excess - N

$$N = 127$$

Excess - N

$$N = 127$$

Binary pattern for number + N

Excess-N

$N = 127$

Binary pattern for number + N
unsigned

0 0 0 0 0 0 0 0

0



Excess - N

$N = 127$

0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 1
⋮

Binary pattern for number + 127
Excess Unsigned

-127 0
-126 1

| | | | | | | |

Excess - N

$N = 127$

0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 1
⋮

Binary pattern for number + 127
Excess Unsigned

-127 0
-126 1

1 1 1 1 1 1 1 1

128

255

Floating point number

$$x = (0.b_1 b_2 b_3 \dots)_2$$

Floating point number

$$x = (0.b_1 b_2 b_3 \dots)_2$$

$$= b_1 \times 2^{-1} + b_2 \times 2^{-2} + \dots$$

$$= (4b_1 + 2b_2 + b_3) \times 8^{-1} + (4b_4 + 2b_5 + b_6) \times$$

↑ Octal form

Floating point number

$$= (0.b_1 b_2 b_3 \dots)_2$$

$$= b_1 \times 2^{-1} + b_2 \times 2^{-2} \dots$$

$$= (b_1 + 2b_2 + b_3) \times 8^{-1} + (4b_4 + 2b_5 + b_6) \times 8^{-2}$$

↑ Octal form

Floating point number

$$= (0.b_1 b_2 b_3 \dots)_2$$

$$= b_1 \times 2^{-1} + b_2 \times 2^{-2} \dots$$

$$(4b_1 + 2b_2 + b_3) \times 8^{-1} + (4b_4 + 2b_5 + b_6) \times 8^{-2}$$

↑ Octal form

Floating point number

$$= (0.b_1 b_2 b_3 \dots)_2$$

$$= b_1 \times 2^{-1} + b_2 \times 2^{-2} \dots$$

$$= (4b_1 + 2b_2 + b_3) \times 8^{-1} + (4b_4 + 2b_5 + b_6) \times 8^{-2}$$

↑ Octal form

$$(101 \ 011 \ 110 \ 011)_2$$

Normalized floating points

2

Normalized floating points

37.218

→

3.7218×10^1

0.37218×10^2

Normalized floating point

37.218

→

3.7218×10^1

0.37218×10^2

Binary normalized floating point

2
x

Normalized floating points

37.218

→

3.7218×10^1
 0.37218×10^2

Binary normalized floating point

$$x = \pm q \times 2^m$$

$$q = (0.b_1 b_2 \dots)_2$$

Normalized floating point

37.218

→

3.7218×10^1
 0.37218×10^2

Binary normalized floating point

$$x = \pm q \times 2^m$$

$$q = (0.b_1 b_2 \dots)_2 \quad b_1 \neq 0$$

Normalized floating point

37.218

→

3.7218×10^1
 0.37218×10^2

Binary normalized floating point

$$x = \pm q \times 2^m$$

$$q = (0.b_1 b_2 \dots)_2 \quad b_1 \neq 0 \quad \frac{1}{2} \leq q < 1$$

or

Normalized floating point

37.218

→

3.7218×10^1
 0.37218×10^2

Binary normalized floating point

$$x = \pm q \times 2^m$$

$$q = (0.b_1 b_2 \dots)_2 \quad b_1 \neq 0 \quad \frac{1}{2} \leq q < 1$$

$$\text{or } q = (1.b_1 b_2 \dots)_2$$

Normalized floating points

37.218

→

3.7218×10^1

0.37218×10^2

Binary normalized floating point

$$x = \pm q \times 2^m$$

$$q = (0.b_1 b_2 \dots)_2$$

$$b_1 \neq 0 \quad \frac{1}{2} \leq q < 1$$

or

$$q = (1.b_1 b_2 \dots)_2$$

$$1 \leq q < 2$$

$$\frac{1}{10} = (0,1)_{10} = \overline{(\quad)}$$

$$\frac{1}{10} = (0.1)_{10} = (0.0631463146314\overline{6314})_{16}$$

$$\frac{1}{10} = (0,1)_{10} = (0,0631463146314\dots)_{10}$$

$$= 0.$$

-2

$$\frac{1}{10} = (0.1)_{10} = (0.0631463146314\dots)_8$$
$$= (0.000110011001100)_2$$

Truncated

$$\frac{1}{10} = (0.1)_{10} = (0.0631463146314\dots)_8$$
$$= (0.000110011001100)_2$$

Truncated

IEEE Standard Hoisting Part



IEEE Standard Floating Point

$$\pm q \times 2^m$$

sign $q \rightarrow$ 1 bit

$|m| \rightarrow$ 8 bits

IEEE Standard Floating Point

$$\pm q \times 2^m$$

sign $q \rightarrow$ 1 bit

$|m| \rightarrow$ 8 bits

number $q \rightarrow$ 23 bits

IEEE Standard Floating Point

$$\pm q \times 2^m$$

sign $q \rightarrow$ 1 bit

$|m| \rightarrow$ 8 bits

number $q \rightarrow$ 23 bits



IEEE Standard Floating Point

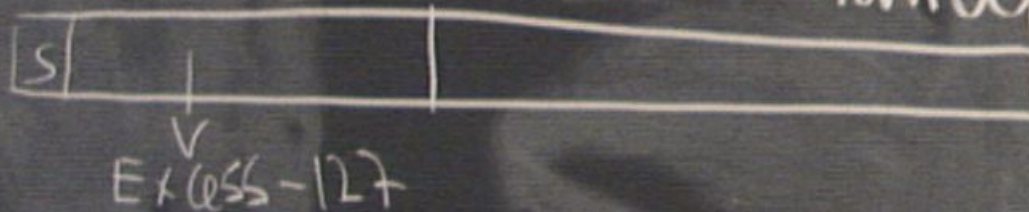
$$\pm q \times 2^m$$

sign $q \rightarrow$ 1 bit

$|m| \rightarrow$ 8 bits

number $q \rightarrow$ 23 bits

\leftarrow 32 bits



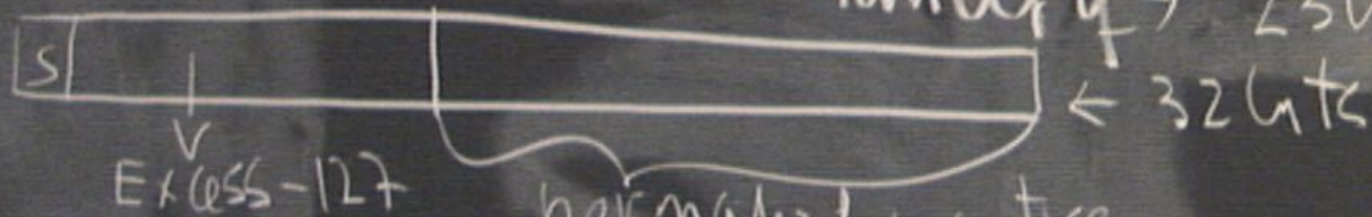
IEEE Standard Floating Point

$$\pm q \times 2^m$$

sign $q \rightarrow$ 1 bit

$|m| \rightarrow$ 8 bits

number $q \rightarrow$ 23 bits



normalized mantissa

$1.b_1 \dots b_{23}$

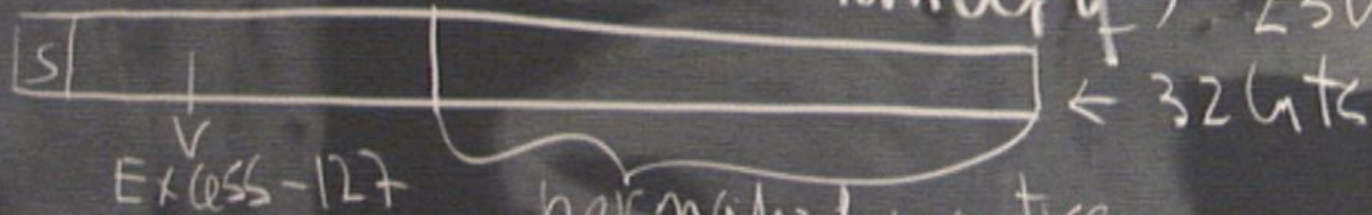
IEEE Standard Floating Point

$$\pm q \times 2^m$$

sign $q \rightarrow$ 1 bit

$|m| \rightarrow$ 8 bits

number $q \rightarrow$ 23 bits



① $b_1 \dots b_{23}$
 \rightarrow Not written

Exponent

C-127

Exponent

$$-126 \leq C - 127 \leq 127$$

Largest number

$$(1.1111 \dots 1)_2 \cdot 2^{127}$$

Exponent

$$-126 \leq C - 127 \leq 127$$

Largest number

$$(1.1111 \dots 1)_2 \cdot 2^{127} = (2 - 2^{-23}) \cdot 2^{127} \approx 2^{128}$$

Exponent

$$-126 \leq C-127 \leq 127$$

Largest number

$$(1.1111 \dots 1)_2$$

$$2^{127} = (2 - 2^{-23}) \times 2^{127}$$
$$\approx 2^{128} \approx 3.4 \times 10^{38}$$

Exponent

$$-127 \leq C - 127 \leq 127$$

Largest number

$$(1.1111 \dots)_2 \cdot 2^{127} = (2 - 2^{-23}) \cdot 2^{127} \\ \approx 2^{128} \approx 3.4 \times 10^{38}$$

Smallest number

$$2^{-127} \approx 10^{-38}$$

10⁻³⁸

Smallest number

$$2^{-127} \approx 10^{-38}$$

Machine Epsilon

$$1.b_1 \dots b_{23}$$

then

$$1.0 + 2^{-23} \neq 1.0$$

$$1.0 + 2^{-24}$$

10³⁸

Smallest number

$$2^{-127} \approx 10^{-38}$$

Machine Epsilon

$$1.b_1 \dots b_{23}$$

then

$$\left. \begin{aligned} 1.0 + 2^{-23} &\neq 1.0 \\ 1.0 + 2^{-24} &= 1.0 \end{aligned} \right\} \text{Epsilon} = 2$$

10⁻³⁸

Smallest number

$$2^{-127} \approx 10^{-38}$$

Machine Epsilon

$$1.b_1 \dots b_{23}$$

then

$$\left. \begin{aligned} 1.0 + 2^{-23} &\neq 1.0 \\ 1.0 + 2^{-24} &= 1.0 \end{aligned} \right\}$$

$$\text{epsilon} = 2^{-23}$$

$$\underline{\underline{1.2 \times 10^{-7}}}$$

10³⁸

Smallest number

$$2^{-127} \approx 10^{-38}$$

Machine Epsilon

$$1.b_1 \dots b_{23}$$

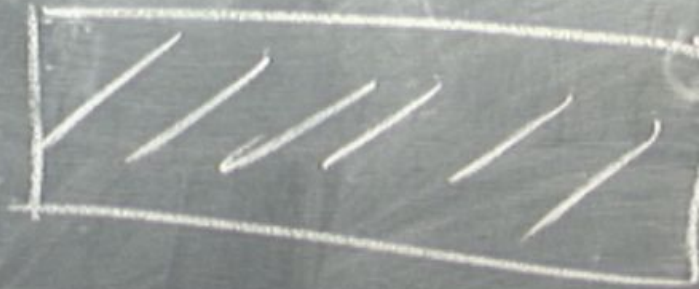
then

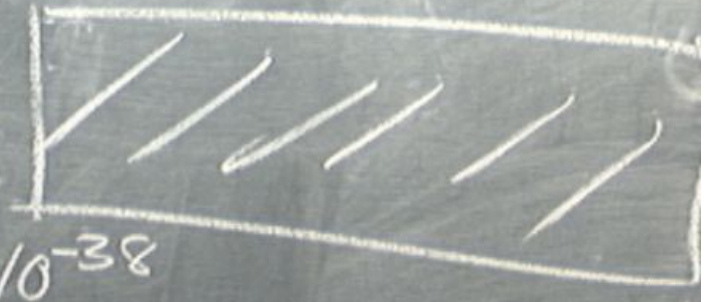
$$\left. \begin{aligned} 1.0 + 2^{-23} &\neq 1.0 \\ 1.0 + 2^{-24} &= 1.0 \end{aligned} \right\}$$

6-7 digits
of precision

$$\text{Epsilon} = 2^{-23} = 1.2 \times 10^{-7}$$

10³⁸







-10^{38}

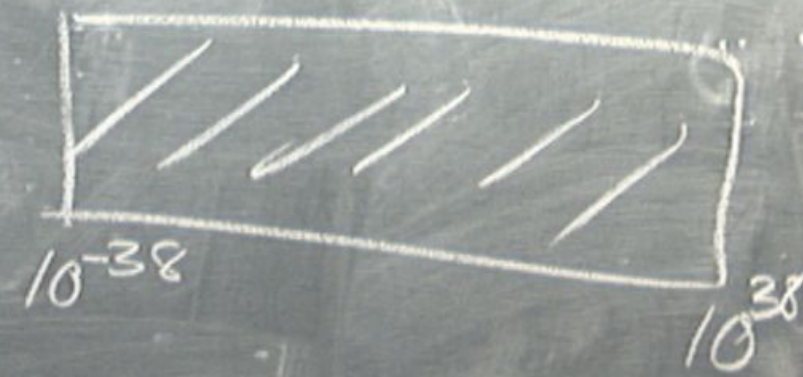
-10^{-38}



10^{-38}

10^{38}

new



↓
Ove

new



↓
Over r

Under floor
Set to 0.0

new



negative overflow
negative underflow



Underflow
Set to 0.0

↓
Overflow

nan



-10^{38}

-10^{38}



10^{-38}

10^{38}

negative overflow
negative underflow

Underflow
Set to 0.0

↓
Overflow

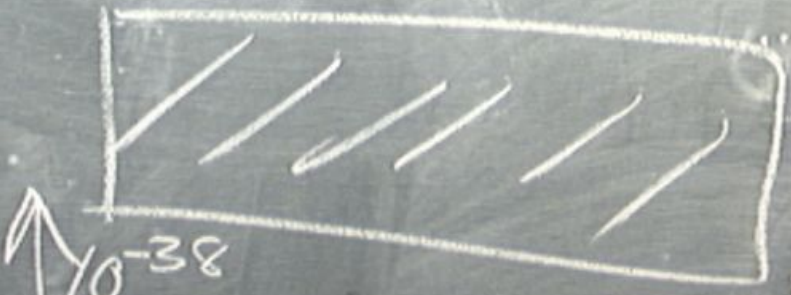
NaN

nan



-10^{38}

negative overflow
underflow



Underflow
Set to 0.0

NaN not a number

Double Precision 64-bit

Double Precision 64-bit

- 52 bits for the mantissa

- 11 bits for the exponent

-1022 — 1023

Double Precision 64-bit

- 52 bits for the mantissa

- 11 bits for the exponent

-1022 — 1023

$$\epsilon = 2.2 \times 10^{-16}$$

Double Precision 64-bit

- 52 bits for the mantissa

- 11 bits for the exponent

- 1022 — 1023

$$\epsilon = 2.2 \times 10^{-16}$$

Largest number

$$1.8 \times 10^{308}$$

Double Precision 64-bit

- 52 bits for the mantissa

- 11 bits for the exponent

-1022 — 1023

$$\epsilon = 2.2 \times 10^{-16}$$

Largest number
Smallest number

$$1.8 \times 10^{308}$$
$$10^{-308}$$

Toy Example

$$X = \pm (0, v_1, v_2, v_3)$$

Joy Example

$$X = \pm (0, l_1, l_2, l_3) \times 2^{\pm k}$$

show

Joy Example

$$x = \pm (0, b_1 b_2 b_3) \times 2^{\pm k} \quad b_i, k_i \in \{0, 1\}$$

0.000

Smaller

Joy Example

$$x = \pm (0, b_1 b_2 b_3) \times 2^{\pm k} \quad b_i, k_i \in \{0, 1\}$$

0.000

0.001

0.010

0.011

Joy Example

$$x = \pm (0, b_1 b_2 b_3) \times 2^{\pm k} \quad b_i, k_i \in \{0, 1\}$$

0:000

0,001

0,010

0,011

0,100

0,101

0,110

0,111

Smaller

Joy Example

$$x = \pm (0, b_1 b_2 b_3) \times 2^{\pm k} \quad b_i, k_i \in \{0, 1\}$$

0.000	0
0,001	1/8
0,010	2/8
0,011	3/8
0,100	4/8
0,101	5/8
0,110	6/8
0,111	7/8

5mmw

Joy Example

$$x = \pm (0, b_1 b_2 b_3) \times 2^{\pm k} \quad b_i, k_i \in \{0, 1\}$$

	$\times 2^0$	$\times 2^2$
0.000	0	0
0.001	1/8	2/8
0.010	2/8	4/8
0.011	3/8	6/8
0.100	4/8	8/8
0.101	5/8	10/8
0.110	6/8	12/8
0.111	7/8	14/8

Smaller ...

Joy Example

$$X = \pm (0, b_1, b_2, b_3) \times 2^{\pm k} \quad b_i, k_i \in \{0, 1\}$$

	$\times 2^0$	$\times 2^1$	$\times 2^{-1}$
0.000	0	0	
0.001	1/8	2/8	
0.010	2/8	4/8	
0.011	3/8	6/8	
0.100	4/8	8/8	
0.101	5/8	10/8	
0.110	6/8	12/8	
0.111	7/8	14/8	

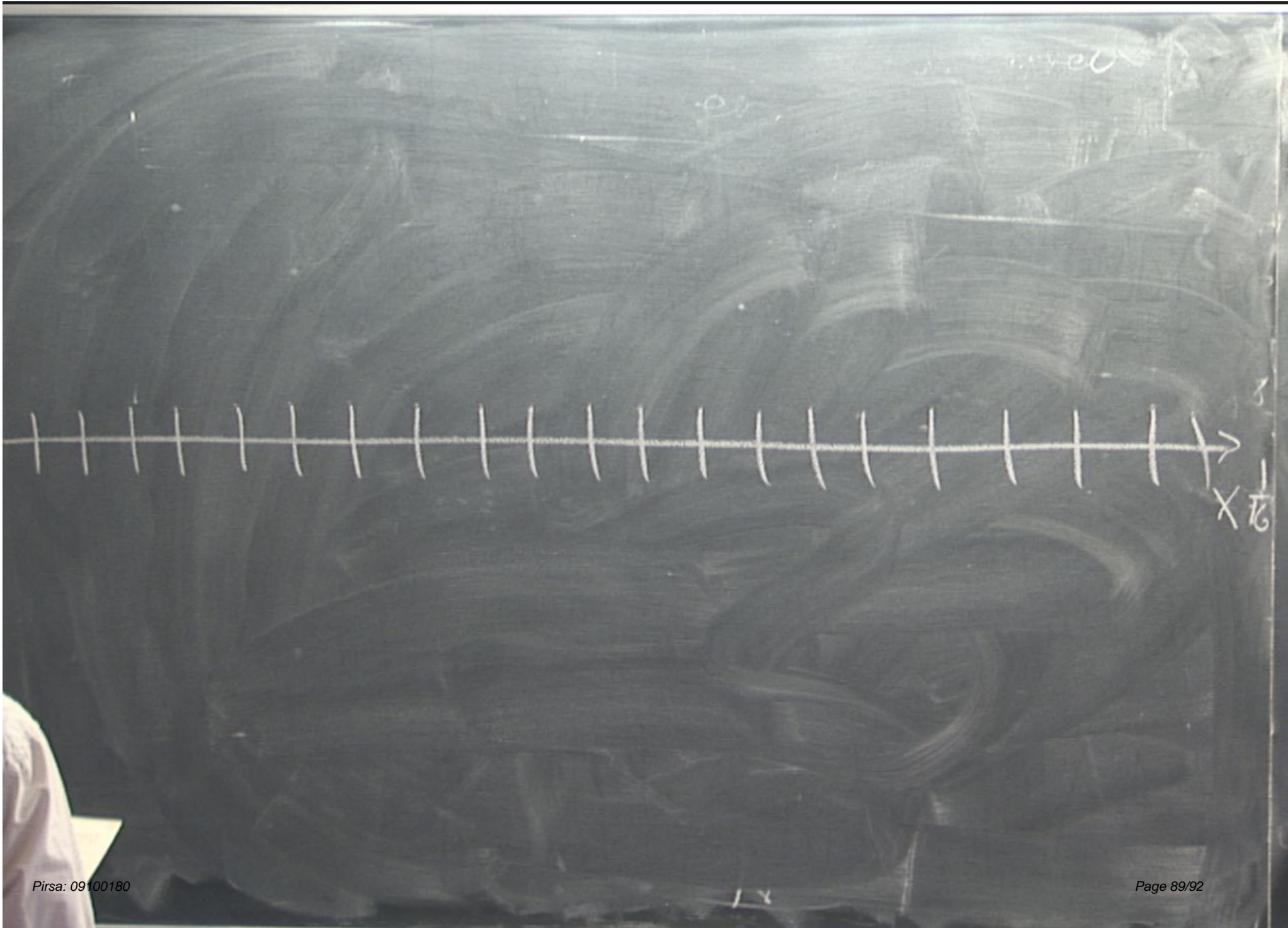
Smaller ...

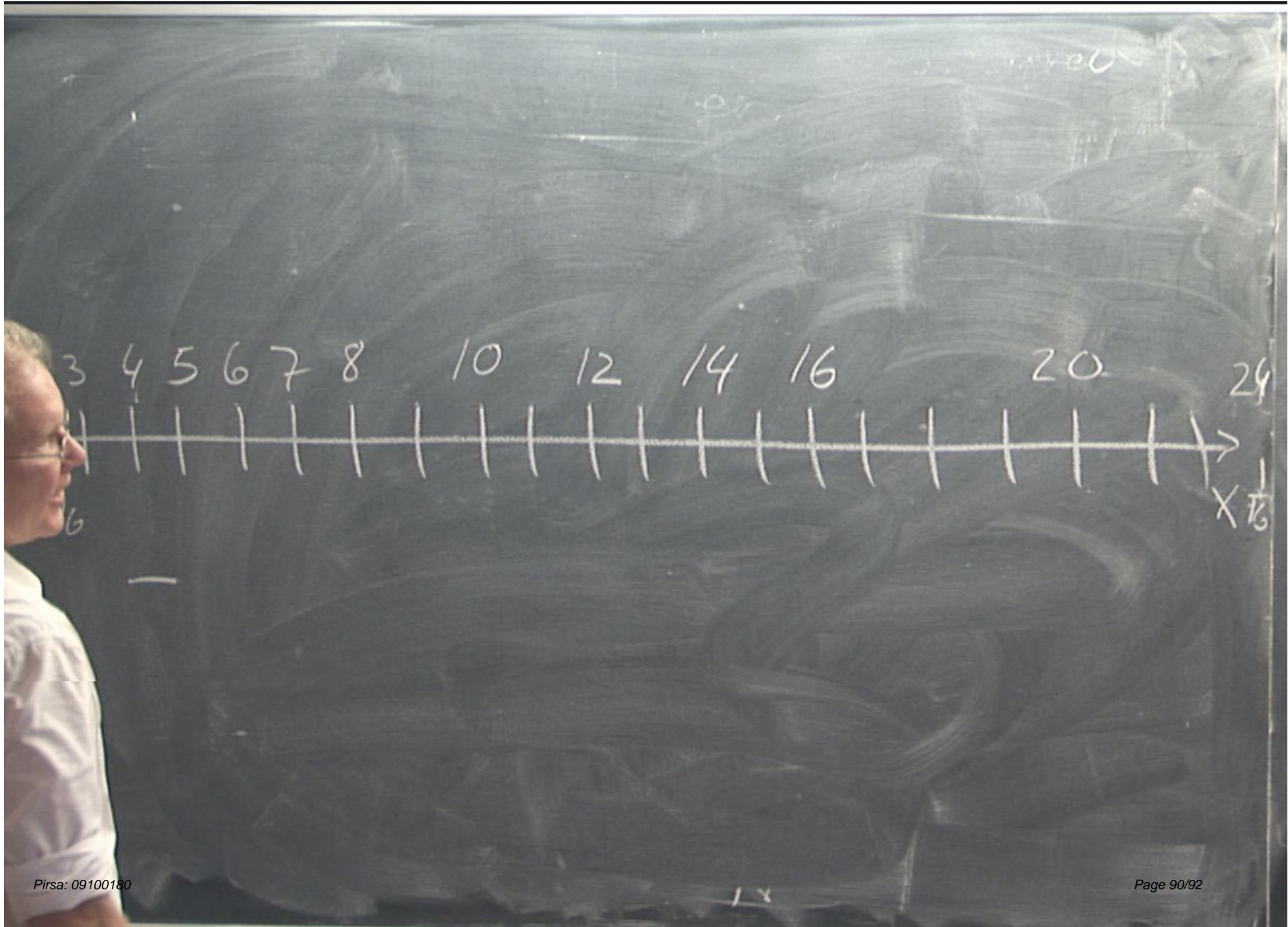
Joy Example

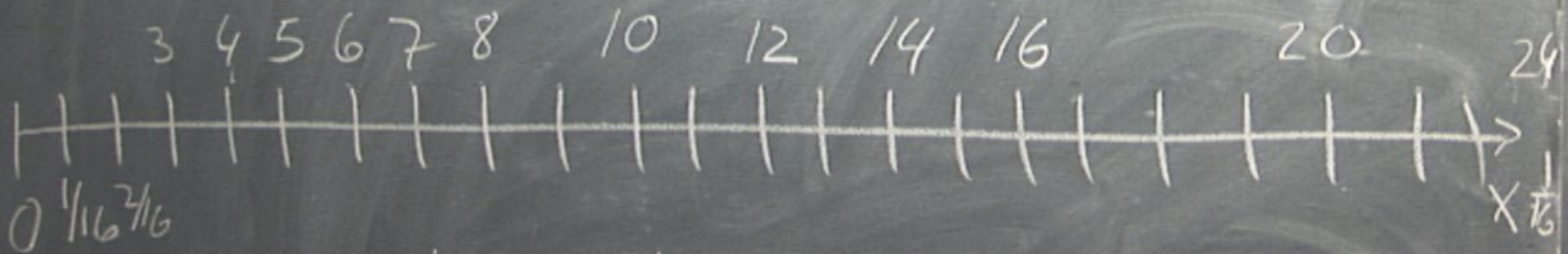
$$X = \pm (0, b_1, b_2, b_3) \times 2^{\pm k} \quad b_i, k_i \in \{0, 1\}$$

	$\times 2^0$	$\times 2^1$	$\times 2^{-1}$
0.000	0	0	0
0.001	1/8	2/8	1/16
0.010	2/8	4/8	2/16
0.011	3/8	6/8	3/16
0.100	4/8	8/8	4/16
0.101	5/8	10/8	5/16
0.110	6/8	12/8	6/16
0.111	7/8	14/8	7/16

Smaller ...







- Not evenly spaced

not normalized

Hole at zero

