

Class Progress

Basics of Linux, gnuplot, C

Visualization of numerical data

Roots of nonlinear equations

(Midterm 1)

Solutions of systems of linear equations

Solutions of systems of nonlinear equations

Monte Carlo simulation

Interpolation of sparse data points

Numerical integration

(Midterm 2)

Solutions of ordinary differential equations

Preserving Numerical Precision

Look at the rules of basic addition of numbers in base 10, for example, $3210 + 4.56$:

$$\begin{array}{r} 3210.00 \\ + 4.56 \\ \hline 3214.56 \end{array}$$

We shift numbers horizontally so that the decimal points are in vertical alignment

Now how about using “scientific notation”, for example, $3.21 \times 10^6 + 4.56 \times 10^{-4}$:

$$\begin{array}{r} 3.21000000000 \times 10^6 \\ 0.00000000456 \times 10^6 \\ \hline 3.21000000456 \times 10^6 \end{array}$$

Now the shifting also includes exponent addition or subtraction so that the decimal points are in vertical alignment **and** the exponents agree

What about an example of $3.21 \times 10^{18} + 4.56 \times 10^{-17}$? Do we have enough width on our computation space for the horizontal shifting?

Image for RAM Storage



Each box
has a
numerical
address or
pointer

What is the numerical capacity of each box?

Roundoff Error

Look at example $3.21 \times 10^6 + 2.46 \times 10^{-4}$:

$$3.21000000000 \times 10^6$$

$$0.00000000246 \times 10^6$$

$$3.21000000246 \times 10^6$$

Correct sum

$$3.21000000000 \times 10^6$$

$$0.00000000246 \times 10^6$$

$$3.210000002 \times 10^6$$

Width limit of arithmetic storage and processor

Rounded sum = correct sum - 4.6×10^{-5} roundoff error

$$3.21000000000 \times 10^6$$

$$0.00000000246 \times 10^6$$

$$3.21000000025 \times 10^6$$

Width limit of arithmetic storage and processor

Rounded sum = correct sum + 4.0×10^{-6} roundoff error

Roundoff Error

Look at example 1.0 / 3.0:

$$\begin{aligned} &3.333333333333\dots \times 10^{-1} \\ &3.3333333333 \times 10^{-1} \end{aligned}$$

Width limit of arithmetic storage and processor

Rounded result = correct result - $3.33\dots \times 10^{-13}$ roundoff error

Look at example 2.0 / 3.0:

$$\begin{aligned} &6.666666666666\dots \times 10^{-1} \\ &6.6666666667 \times 10^{-1} \end{aligned}$$

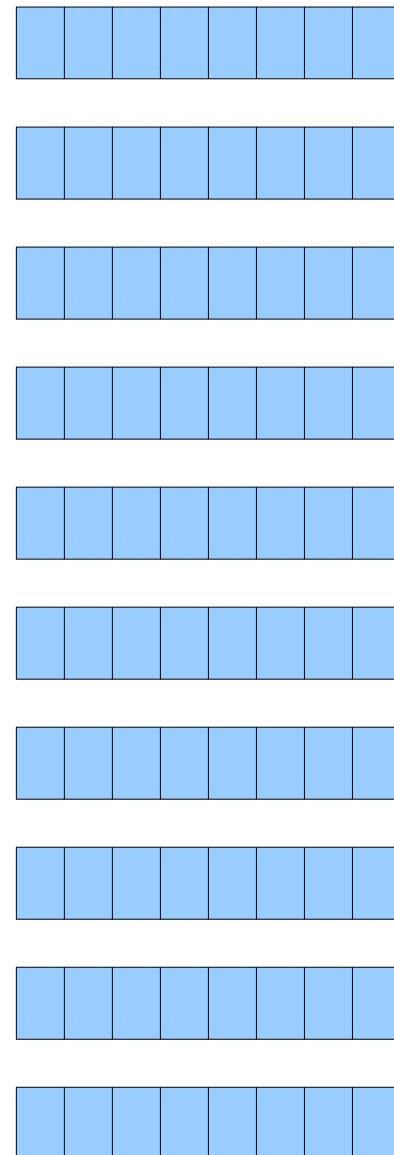
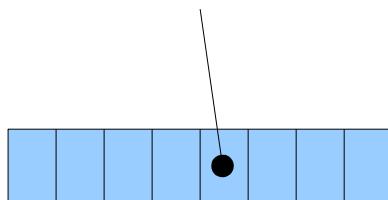
Width limit of arithmetic storage and processor

Rounded result = correct result + $3.33\dots \times 10^{-13}$ roundoff error

Understanding numerical limits and roundoff error is essential in scientific computing. Of course, computer hardware uses binary arithmetic, but the same principles apply.

Map of Bytes in RAM

One byte = 8 bits:
The minimum
addressable chunk
of data

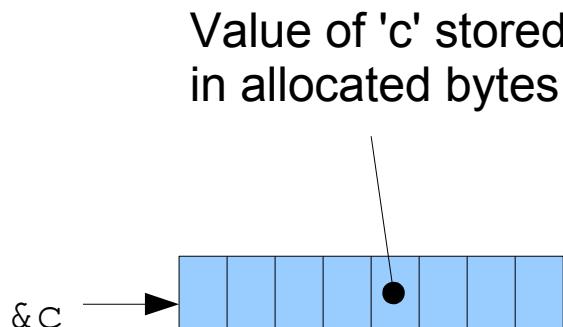


C Data Types for Numerical Programming

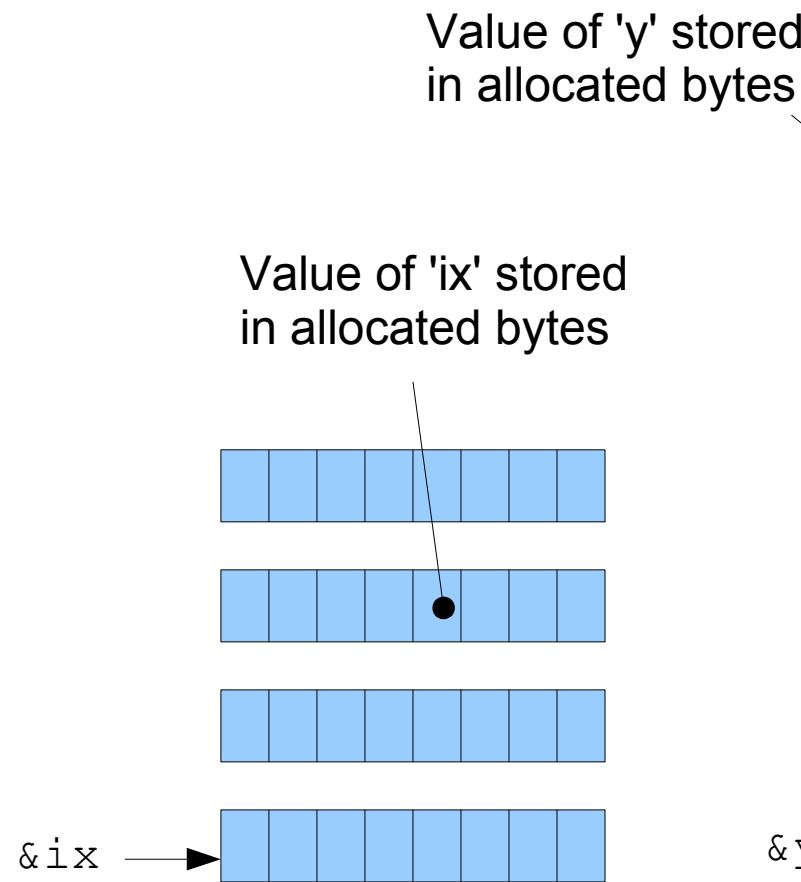
- Integer
 - Named “int”
 - Usually 32 bits, 4 bytes
 - Least significant byte at lowest memory address on Intel based processors
- Double precision floating point
 - Named “double”
 - Usually 64 bits, 8 bytes
- Character
 - Named “char”
 - Usually 8 bits, 1 byte

Allocation of C Data Types in RAM

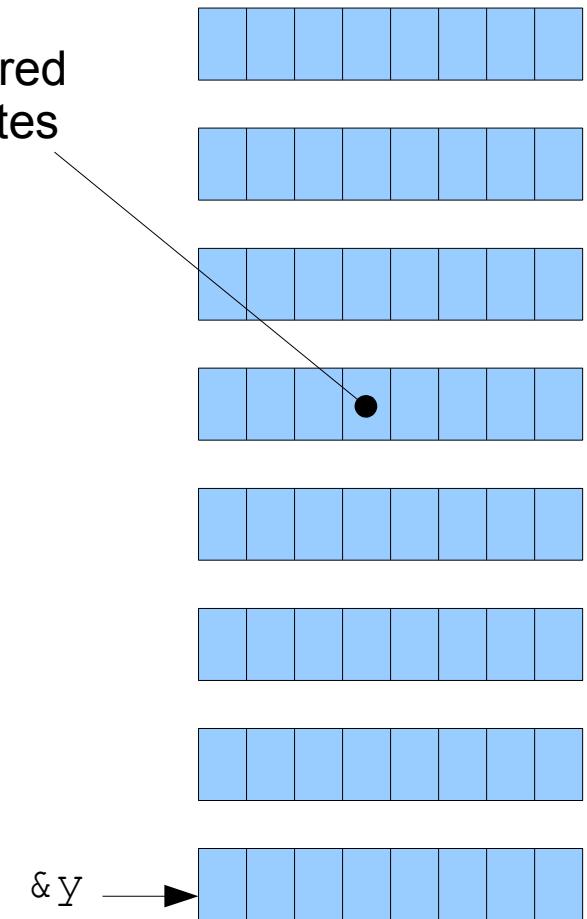
Generate address of an allocated variable with the construction "`&variable name`"



`char c;`



`int ix;`



`double y;`

Total Possible Bit Arrangements

For any collection of n items that can *independently* assume one of two states...



2^7 2^0 ←

1	0	0	1	1	0	0	1
---	---	---	---	---	---	---	---

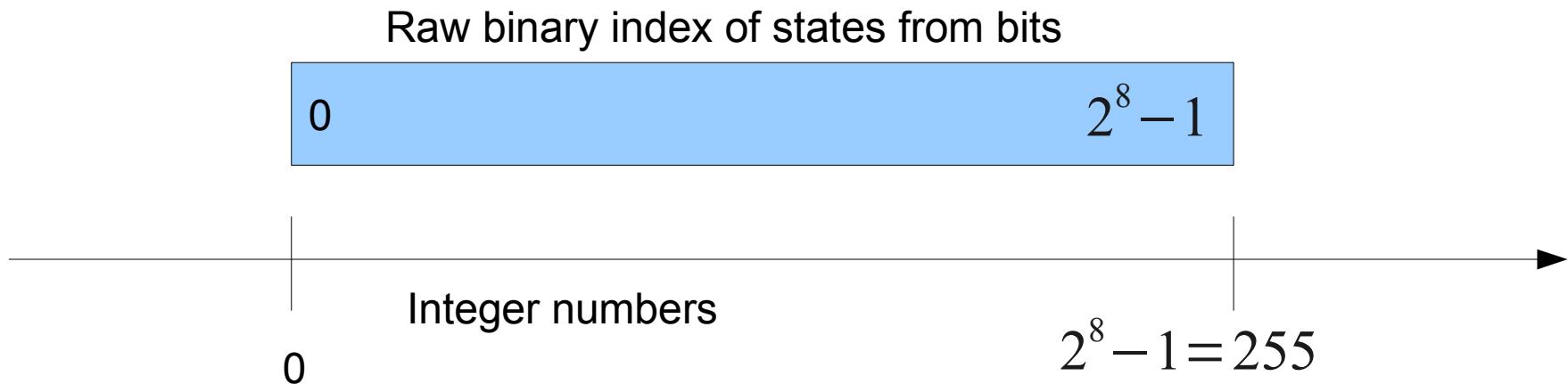
Binary weights to form
an index of all states

The total number of possible unique arrangements is ...

Unsigned 8-bit Character Values

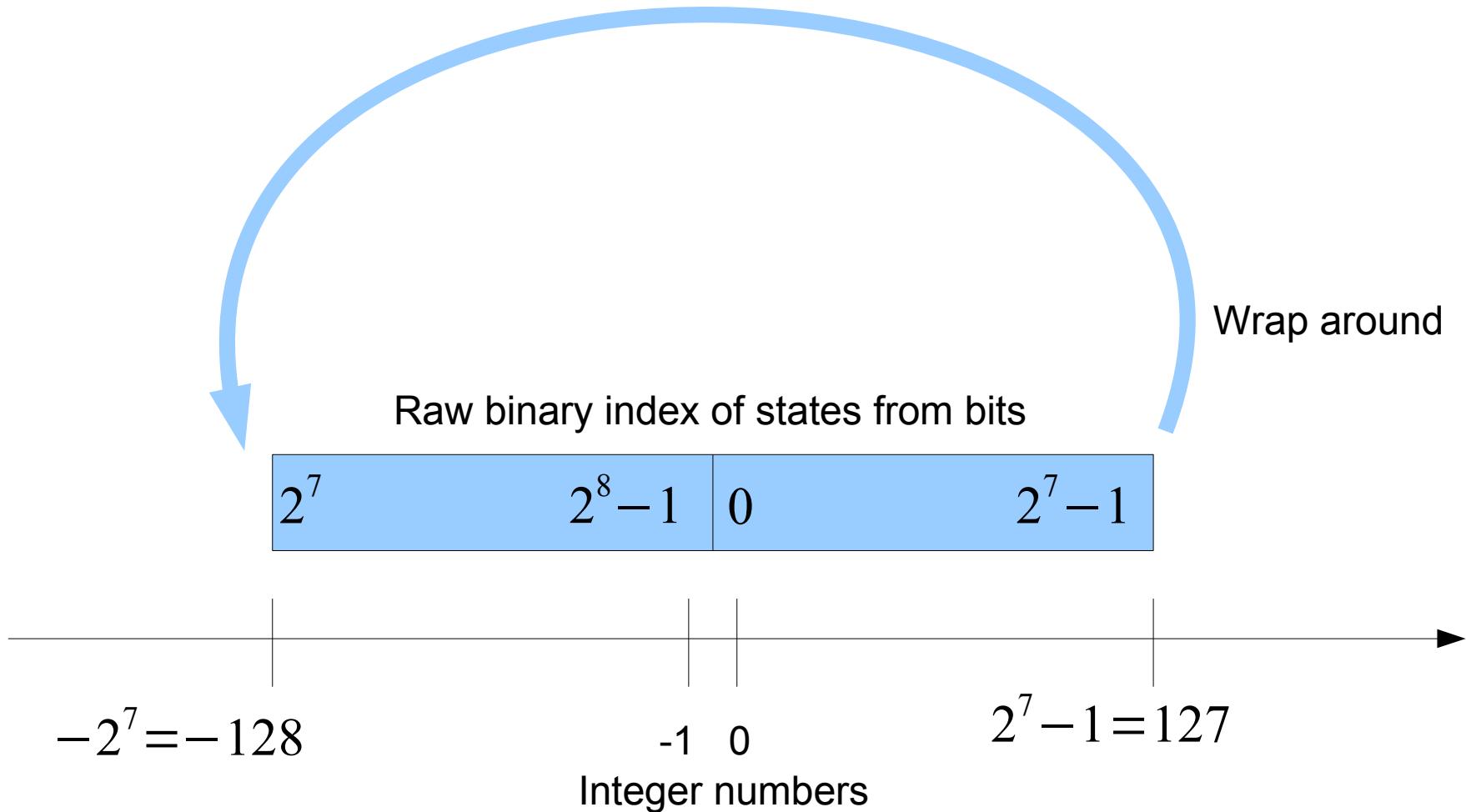
Total number of possible bit arrangements = $2^n = 2^8 = 256$

How do we assign numerical values to each of these arrangements?

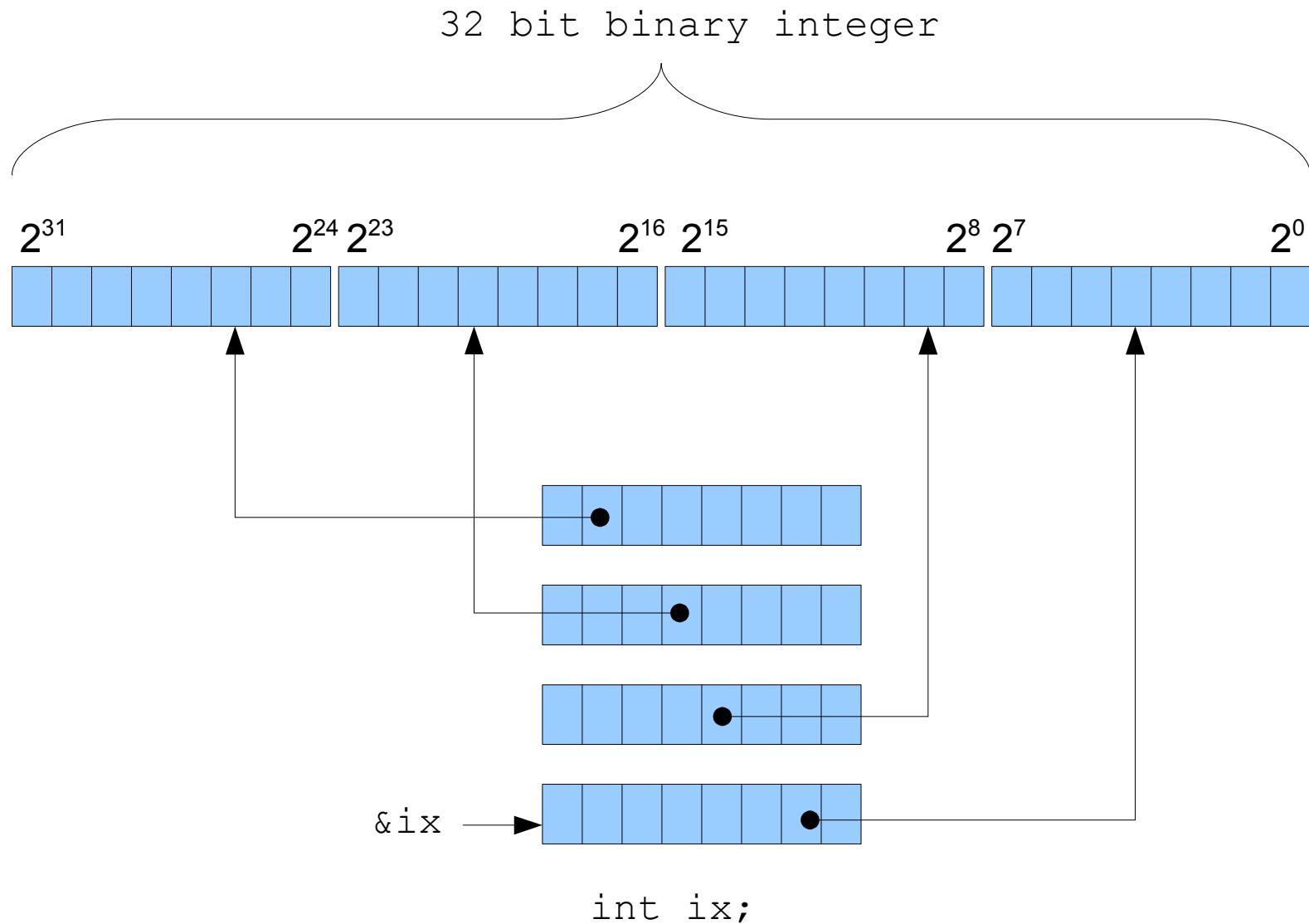


Signed 8-bit Character Values

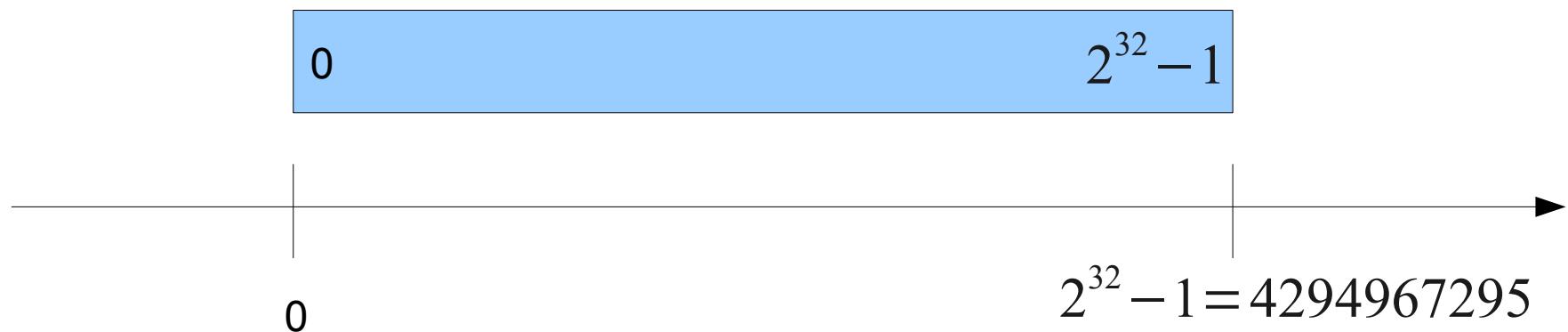
Same possible bit arrangements, but a different numerical mapping...



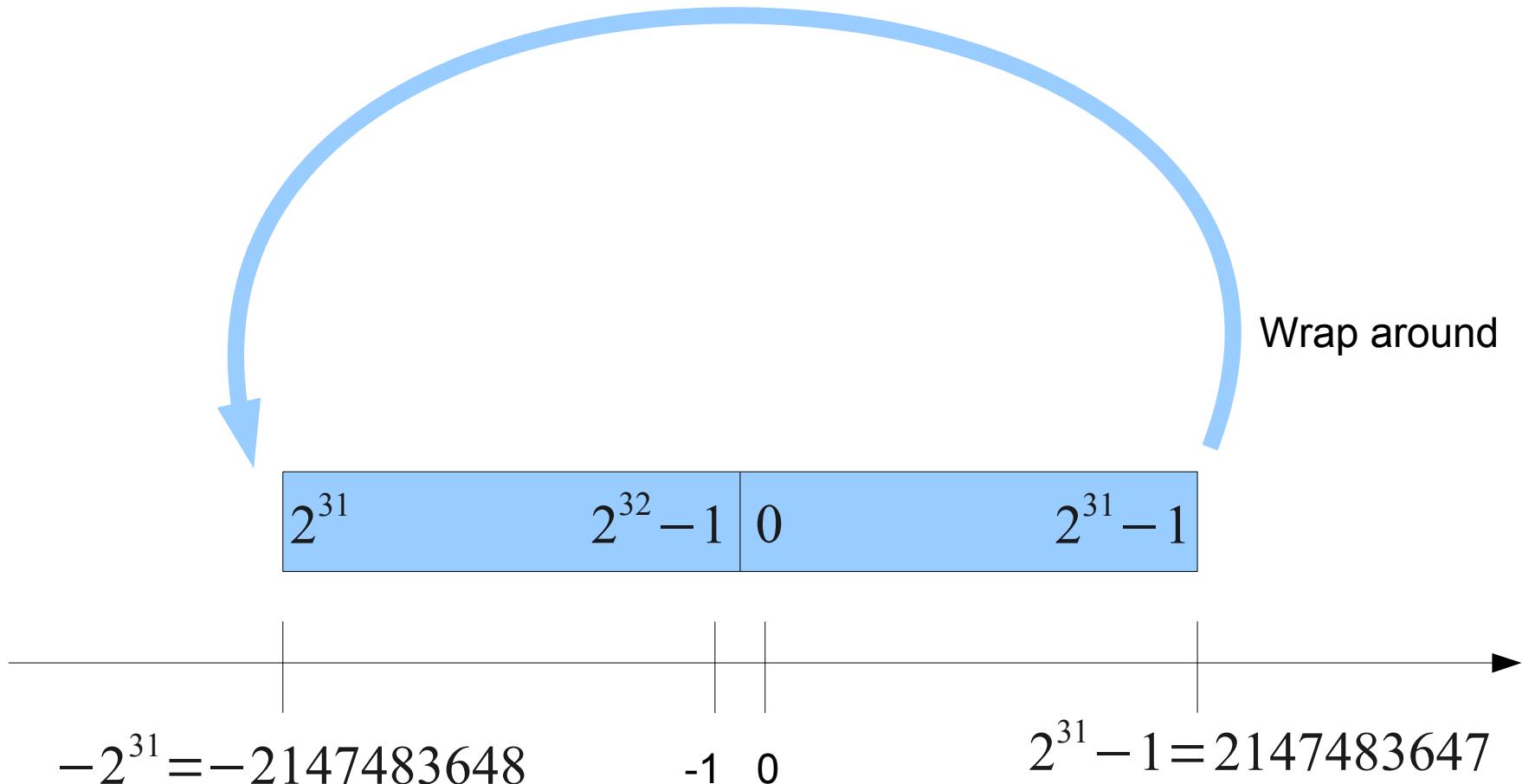
Integer Bit Arrangement for “Little Endian” Processors



Unsigned 32-bit Integer Values



Signed 32-bit Integer Values



Pseudocode for Seeing Integer Wraparound

```
integer i, s ;  
s  $\leftarrow$  0  
for i=1 to 50 do  
    s  $\leftarrow$  s+ 268435456  
    output s  
end for
```

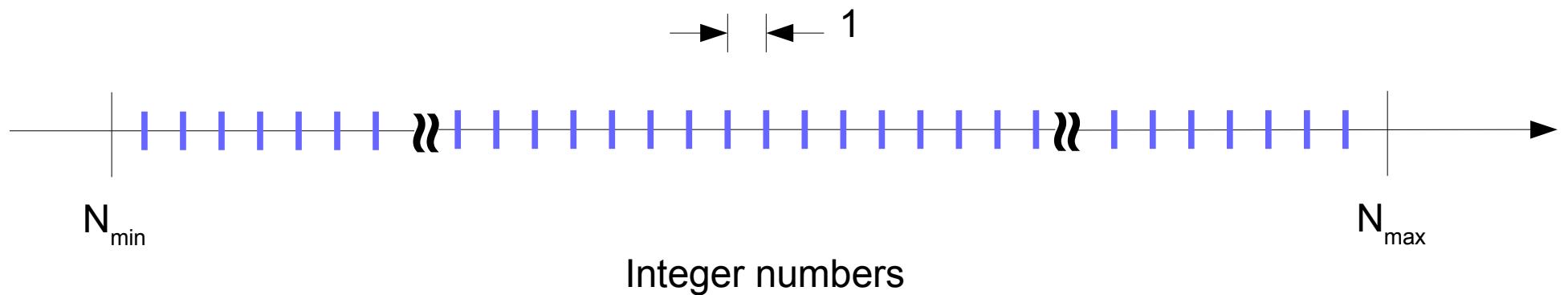
Note: $268435456 = 2^{28}$

Wraparound C Code Output

```
s=268435456  
s=536870912  
s=805306368  
s=1073741824  
s=1342177280  
s=1610612736  
s=1879048192  
s=-2147483648  
s=-1879048192  
s=-1610612736  
s=-1342177280  
s=-1073741824  
s=-805306368  
s=-536870912  
s=-268435456  
s=0  
s=268435456  
s=536870912  
s=805306368  
s=1073741824  
s=1342177280  
s=1610612736  
s=1879048192  
s=-2147483648
```

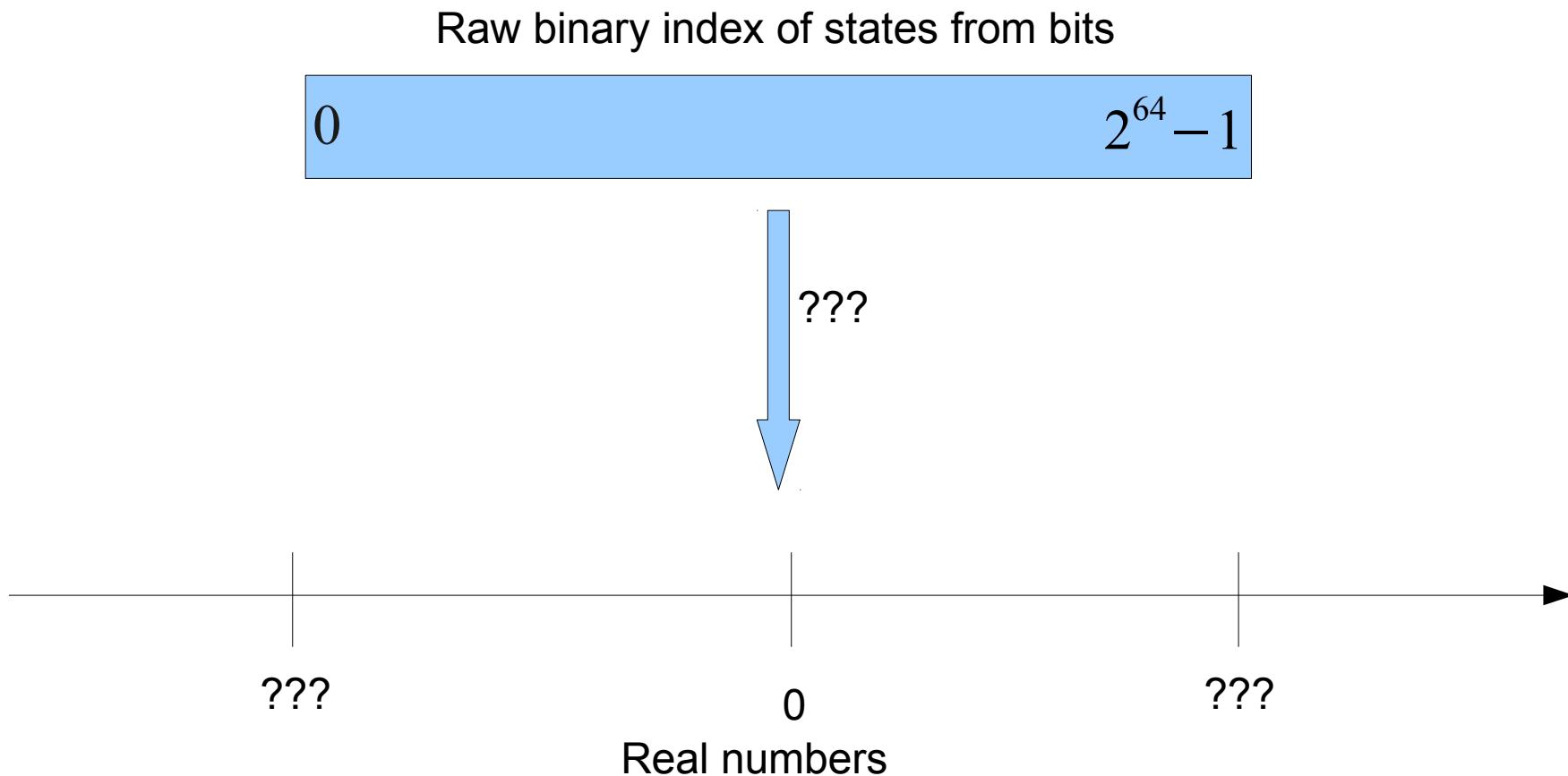
Density on Number Line

For all forms of integer variables, the density of numbers represented is uniform between the minimum and maximum numbers, with constant interval of 1



Representing “Real” Values

There are a finite number of bits available in RAM. We usually allot either 32 or 64 bits for each “real” variable in computer code. There are still either 2^{32} or 2^{64} possible bit combinations. How do we map values on an infinitely divisible line of real numbers?



Floating Point Normalization

When we use “Scientific Notation”, we naturally keep numbers in a normalized form. That is, we typically would write the number 123 as

$$1.23 \times 10^2$$

Or sometimes

$$0.123 \times 10^3$$

But not any of

$$12.3 \times 10^1$$

$$0.0123 \times 10^4$$

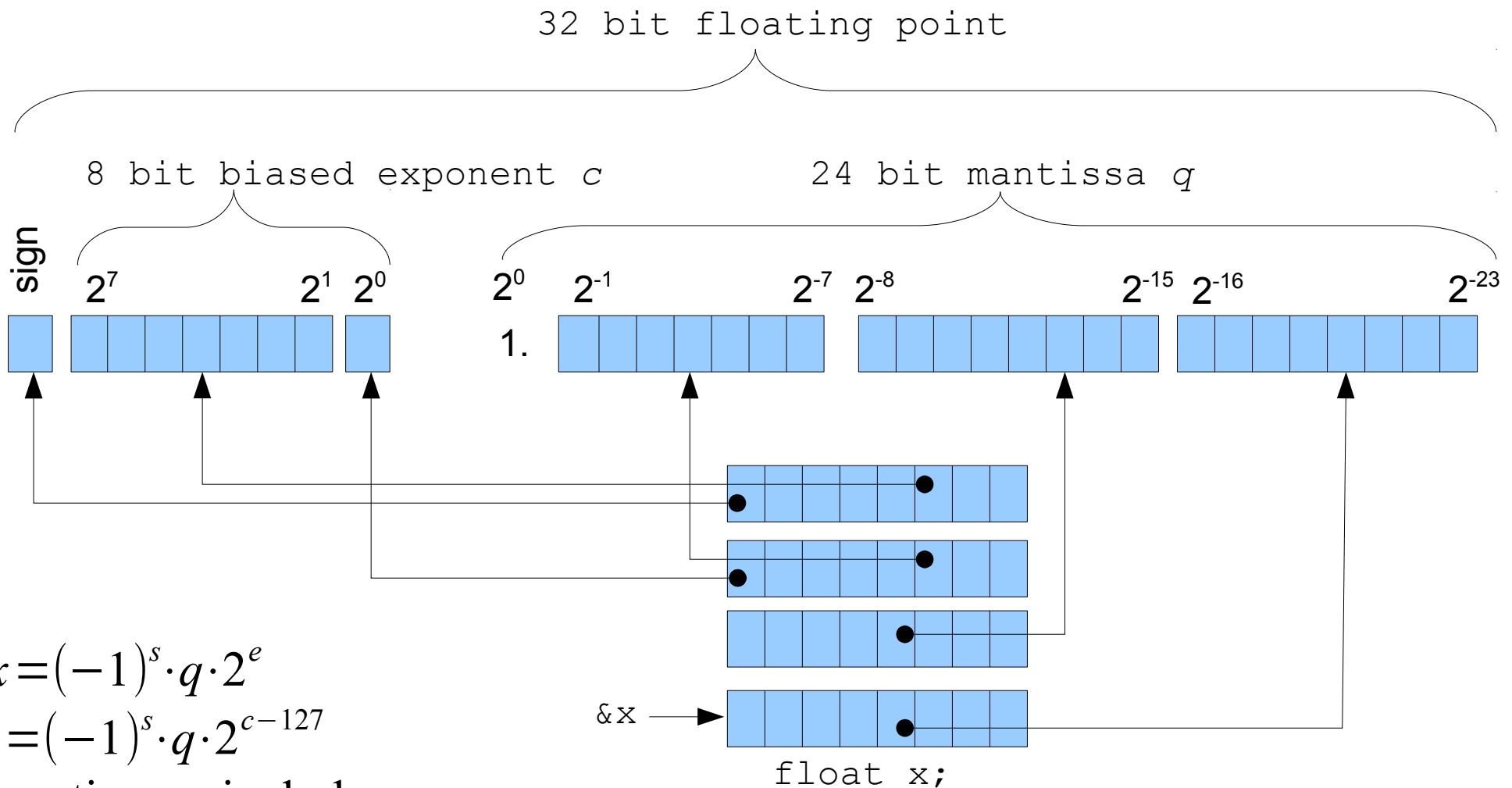
$$123.0 \times 10^0$$

$$0.00123 \times 10^5$$

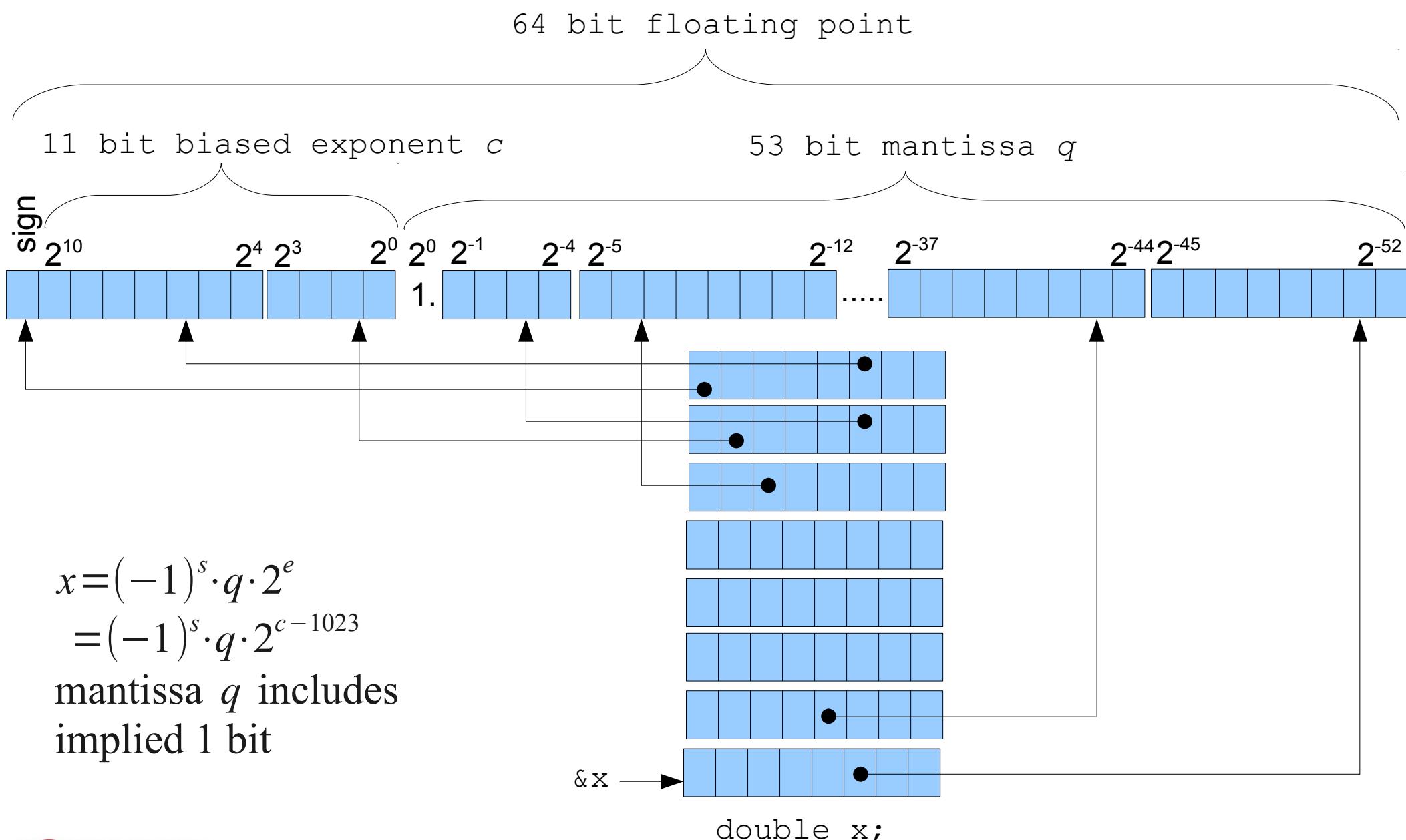
etc...

Similarly, machine hardware will keep floating point numbers normalized, so that the most significant '1' bit is always the hidden '1' bit.

Single Precision Floating Point Bit Arrangement



Double Precision Floating Point Bit Arrangement



Double Precision Examples

Actual bits shown in black, implied '1' bit shown in red

$$1.0 = 0 \ 0111111111 \ 1.0000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000000$$

$s = 0$ $e = 1023 - 1023 = 0$ $q = 1.0$ \rightarrow $1.0 = (-1)^0 \cdot 1.0 \cdot 2^0$

$$-1.0 = 1 \ 0111111111 \ 1.0000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000000$$

$s = 1$ $e = 1023 - 1023 = 0$ $q = 1.0$ \rightarrow $-1.0 = (-1)^1 \cdot 1.0 \cdot 2^0$

$$2.0 = 0 \ 1000000000 \ 1.0000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000000$$

$s = 0$ $e = 1024 - 1023 = 1$ $q = 1.0$ \rightarrow $2.0 = (-1)^0 \cdot 1.0 \cdot 2^1$

Double Precision Examples

Actual bits shown in black, implied '1' bit shown in red

$0.0 = 0 \ 0000000000 \ 1.0000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000000$

$s = 0$ all actual bits 0 \rightarrow 0.0 by convention

$10.0 = 0 \ 1000000010 \ 1.0100\ 00000000\ 00000000\ 00000000\ 00000000\ 00000000$

$s = 0$ $e = 1026 - 1023 = 3$ $q = 1.25$ \rightarrow $10.0 = (-1)^0 \cdot 1.25 \cdot 2^3$

$0.1 = 0 \ 0111111011 \ 1.1001\ 10011001\ 10011001\ 10011001\ 10011001\ 10011001\ 10011010$

$s = 0$ $e = 1019 - 1023 = -4$ $q \approx 1.6$ \rightarrow $0.1 \approx (-1)^0 \cdot 1.6 \cdot 2^{-4}$

Double Precision Examples

Actual bits shown in black, implied '1' bit shown in red

$$0.5 = 0 \ 0111111110 \ 1.000$$

$s = 0$ $e = 1022 - 1023 = -1$ $c = 1.0$ \rightarrow $0.5 = (-1)^0 \cdot 1.0 \cdot 2^{-1}$

$$0.25 = 0 \ 0111111101 \ 1.000$$

$s = 0$ $e = 1021 - 1023 = -2$ $c = 1.0$ \rightarrow $0.25 = (-1)^0 \cdot 1.0 \cdot 2^{-2}$

$$1.0000000000000002 =$$
$$0 \ 0111111111 \ 1.001$$

$s = 0$ $e = 1023 - 1023 = 0$ $c = 1.0000000000000002$ \rightarrow

$$1.0000000000000002 = (-1)^0 \cdot 1.0000000000000002 \cdot 2^0$$

Double Precision Examples

Actual bits shown in black, implied '1' bit shown in red

Largest possible number

$$= 0 \ 11111111110 \ 1.1111 \ 11111111 \ 11111111 \ 11111111 \ 11111111 \ 11111111$$

$s = 0 \quad e = 2046 - 1023 = 1023 \quad c \approx 2.0 \quad \rightarrow \quad (-1)^0 \cdot 2.0 \cdot 2^{1023} \approx 1.8 \cdot 10^{308}$

Example: The known size of the universe is 9.1×10^{10} light years = 8.6×10^{26} m = 8.6×10^{36} Å

Smallest possible number, excluding subnormal numbers

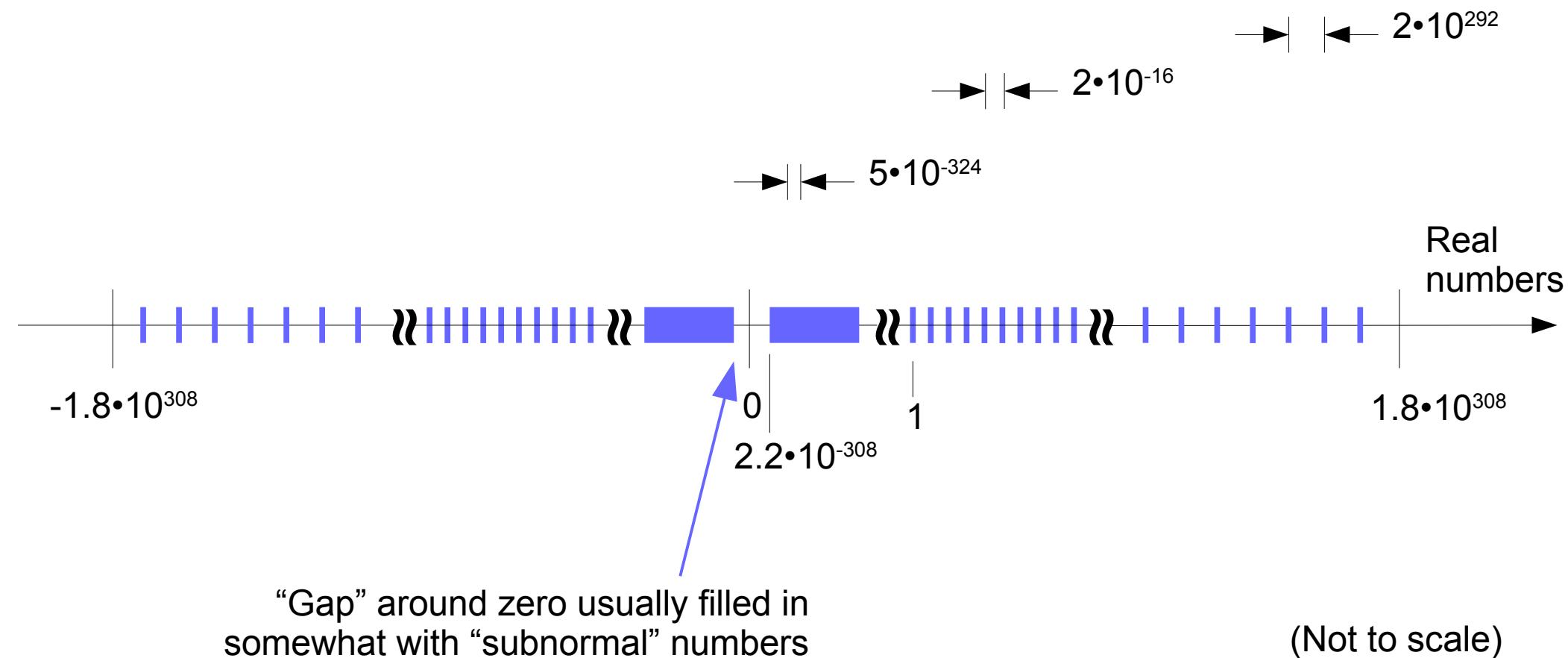
$$= 0 \quad 00000000001 \quad 1. \quad 0000 \ 0000000 \ 0000000 \ 0000000 \ 0000000 \ 0000000 \ 0000000$$

$s = 0 \quad e = 1 - 1023 = -1022 \quad c \approx 1.0 \quad \rightarrow \quad (-1)^0 \cdot 1.0 \cdot 2^{-1022} \approx 2.2 \cdot 10^{-308}$



Density on Number Line

For double precision floating point variables, the density of numbers represented is nonuniform between the minimum and maximum numbers



Pseudocode for Exploring Precision Limit

Mapping a finite number of bit states to an infinitely divisible real number line will clearly leave “gaps” in the line. This is an intrinsic limit of machine precision.

```
integer i; real x, y, z;  
x  $\leftarrow$  1.0  
y  $\leftarrow$  0.125  
for i=1 to 20 do  
    z  $\leftarrow$  x+ y  
    output x, y, z  
    y  $\leftarrow$  y·0.125  
end for
```

C Code Output Using Double Precision

```
x=1.0000000000000000 y=1.25000000000000e-01 z=1.1250000000000000
x=1.0000000000000000 y=1.56250000000000e-02 z=1.0156250000000000
x=1.0000000000000000 y=1.95312500000000e-03 z=1.0019531250000000
x=1.0000000000000000 y=2.44140625000000e-04 z=1.0002441406250000
x=1.0000000000000000 y=3.05175781250000e-05 z=1.0000305175781250
x=1.0000000000000000 y=3.814697265625000e-06 z=1.0000038146972656
x=1.0000000000000000 y=4.7683715820312500e-07 z=1.0000004768371582
x=1.0000000000000000 y=5.9604644775390625e-08 z=1.000000596046448
x=1.0000000000000000 y=7.4505805969238281e-09 z=1.000000074505806
x=1.0000000000000000 y=9.3132257461547852e-10 z=1.00000009313226
x=1.0000000000000000 y=1.1641532182693481e-10 z=1.000000001164153
x=1.0000000000000000 y=1.4551915228366852e-11 z=1.000000000145519
x=1.0000000000000000 y=1.8189894035458565e-12 z=1.000000000018190
x=1.0000000000000000 y=2.2737367544323206e-13 z=1.0000000000002274
x=1.0000000000000000 y=2.8421709430404007e-14 z=1.000000000000284
x=1.0000000000000000 y=3.5527136788005009e-15 z=1.000000000000036
x=1.0000000000000000 y=4.4408920985006262e-16 z=1.000000000000004
x=1.0000000000000000 y=5.5511151231257827e-17 z=1.000000000000000
x=1.0000000000000000 y=6.9388939039072284e-18 z=1.000000000000000
x=1.0000000000000000 y=8.6736173798840355e-19 z=1.000000000000000
```

Pseudocode for Exploring Precision Limit

```
integer i; real t,x,y,z;  
x  $\leftarrow$  1.0  
t  $\leftarrow$  0.1  
for i=1 to 20 do  
    y  $\leftarrow$  x+t  
    z  $\leftarrow$  y-x  
    output x,y,z  
    x  $\leftarrow$  x.8.0  
end for
```

C Code Output Using Double Precision

```
x=1.000000000000000e+00  y=1.100000000000001e+00  z=0.100000000000001
x=8.000000000000000e+00  y=8.099999999999996e+00  z=0.099999999999996
x=6.400000000000000e+01  y=6.409999999999994e+01  z=0.099999999999943
x=5.120000000000000e+02  y=5.121000000000002e+02  z=0.100000000000227
x=4.096000000000000e+03  y=4.096100000000004e+03  z=0.100000000003638
x=3.276800000000000e+04  y=3.276809999999999e+04  z=0.099999999985448
x=2.621440000000000e+05  y=2.621440999999998e+05  z=0.099999999767169
x=2.097152000000000e+06  y=2.097152100000001e+06  z=0.100000000931323
x=1.677721600000000e+07  y=1.677721610000001e+07  z=0.100000014901161
x=1.342177280000000e+08  y=1.342177280999999e+08  z=0.099999940395355
x=1.073741824000000e+09  y=1.073741824099999e+09  z=0.099999046325684
x=8.589934592000000e+09  y=8.589934592100004e+09  z=0.100003814697266
x=6.871947673600000e+10 y=6.871947673610006e+10  z=0.100061035156250
x=5.497558138880000e+11 y=5.4975581388809998e+11  z=0.099755859375000
x=4.398046511040000e+12 y=4.398046511040996e+12  z=0.099609375000000
x=3.518437208883200e+13 y=3.5184372088832102e+13 z=0.101562500000000
x=2.8147497671065600e+14 y=2.8147497671065612e+14 z=0.125000000000000
x=2.2517998136852480e+15 y=2.2517998136852480e+15 z=0.000000000000000
x=1.8014398509481984e+16 y=1.8014398509481984e+16 z=0.000000000000000
x=1.4411518807585587e+17 y=1.4411518807585587e+17 z=0.000000000000000
```

Machine ϵ

The “Machine ϵ ” is the smallest number ϵ such that $1 + \epsilon \neq 1$

For 32-bit float data type, $\epsilon = 2^{-24} \approx 5.96 \times 10^{-8}$

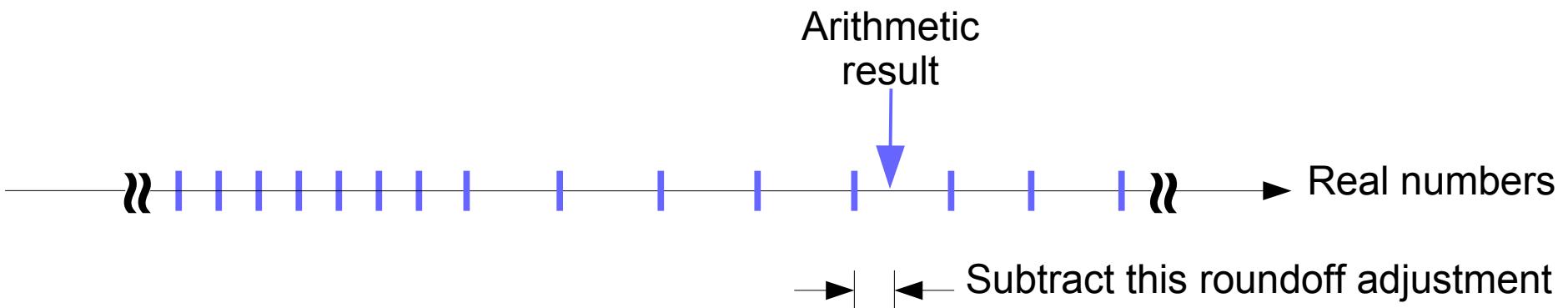
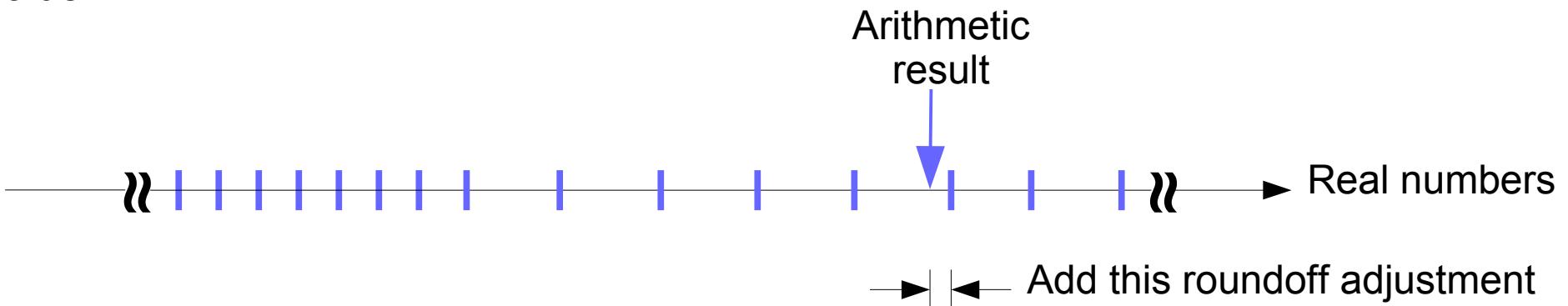
A number slightly larger than this ϵ will be rounded up to 2^{-23} when added to 1.0 and turn the least significant bit to '1'

For 64-bit double data type, $\epsilon = 2^{-53} \approx 1.11 \times 10^{-16}$

A number slightly larger than this ϵ will be rounded up to 2^{-52} when added to 1.0 and turn the least significant bit to '1'

Roundoff Error

When the result of an arithmetic operation lies between two allowed values on the real number line, “rounding off” to some allowed value before the value is stored in memory is equivalent to adding or subtracting some roundoff adjustment to end up at an allowed value



The accumulated roundoff error resulting at the end of a large number of arithmetic operations can be treated statistically

Pseudocode for Accumulating Roundoff Error

Use 64-bit doubles for these reals

Use 32-bit float for this real

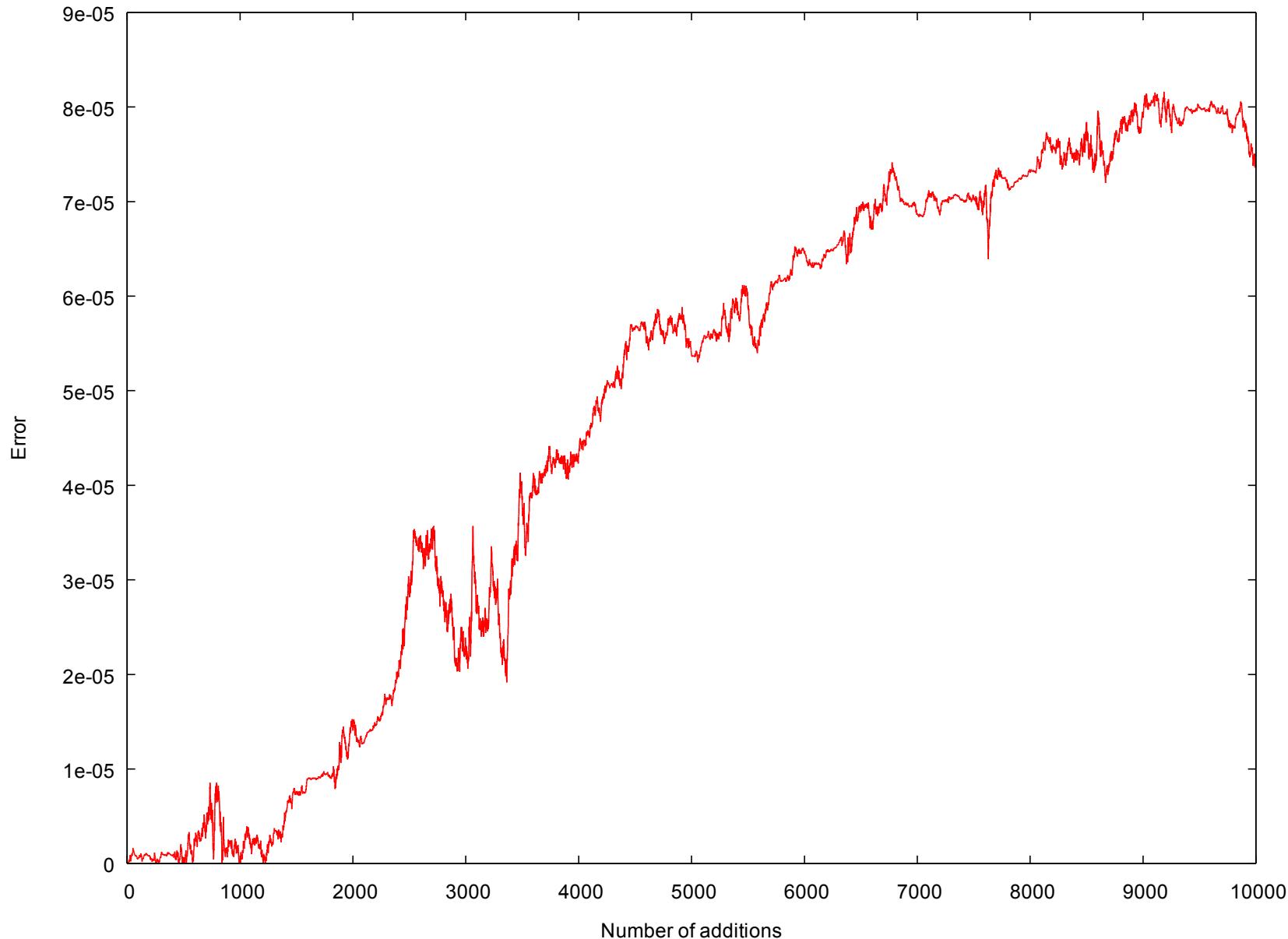
```
integer i ; real r , xd ; real xf ;
xd ← 0.0
xf ← 0.0
for i=1 to 10000 do
    r ← pseudorandom number between -1 and 1
    xf ← xf + r
    xd ← xd + r
    error ← xf - xd
    output i , error
end for
```

Accumulate sum of pseudorandom numbers in both float and double

Interpret difference as accumulated roundoff error for float

Accumulated roundoff error

Each roundoff event can be seen as an addition of a random number in the range $-\varepsilon$ to $+\varepsilon$



Evaluation of Polynomial with Multiple Roots

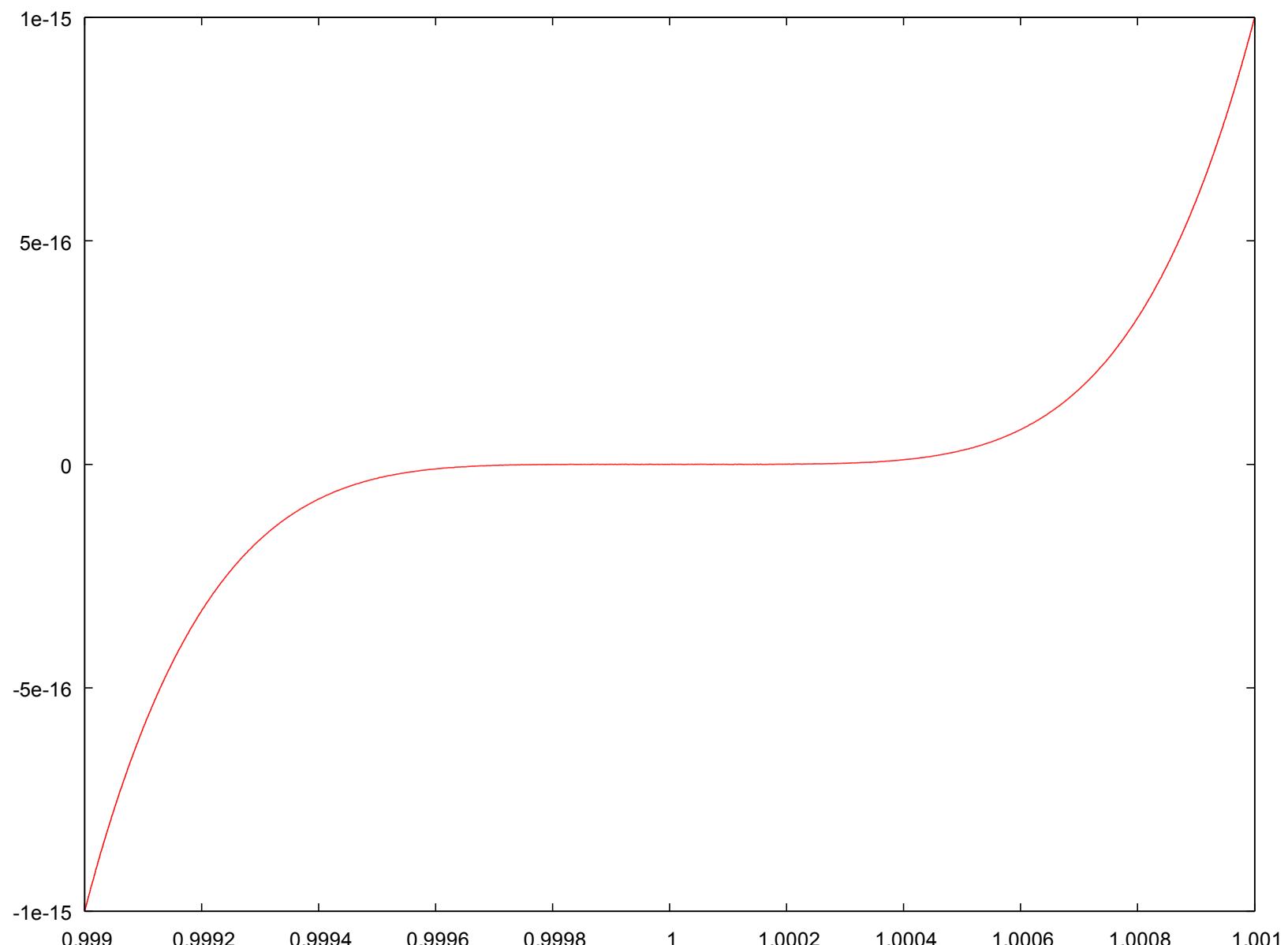
$$y = x^5 - 5x^4 + 10x^3 - 10x^2 + 5x - 1$$

Has root of multiplicity 5 at $x=1$

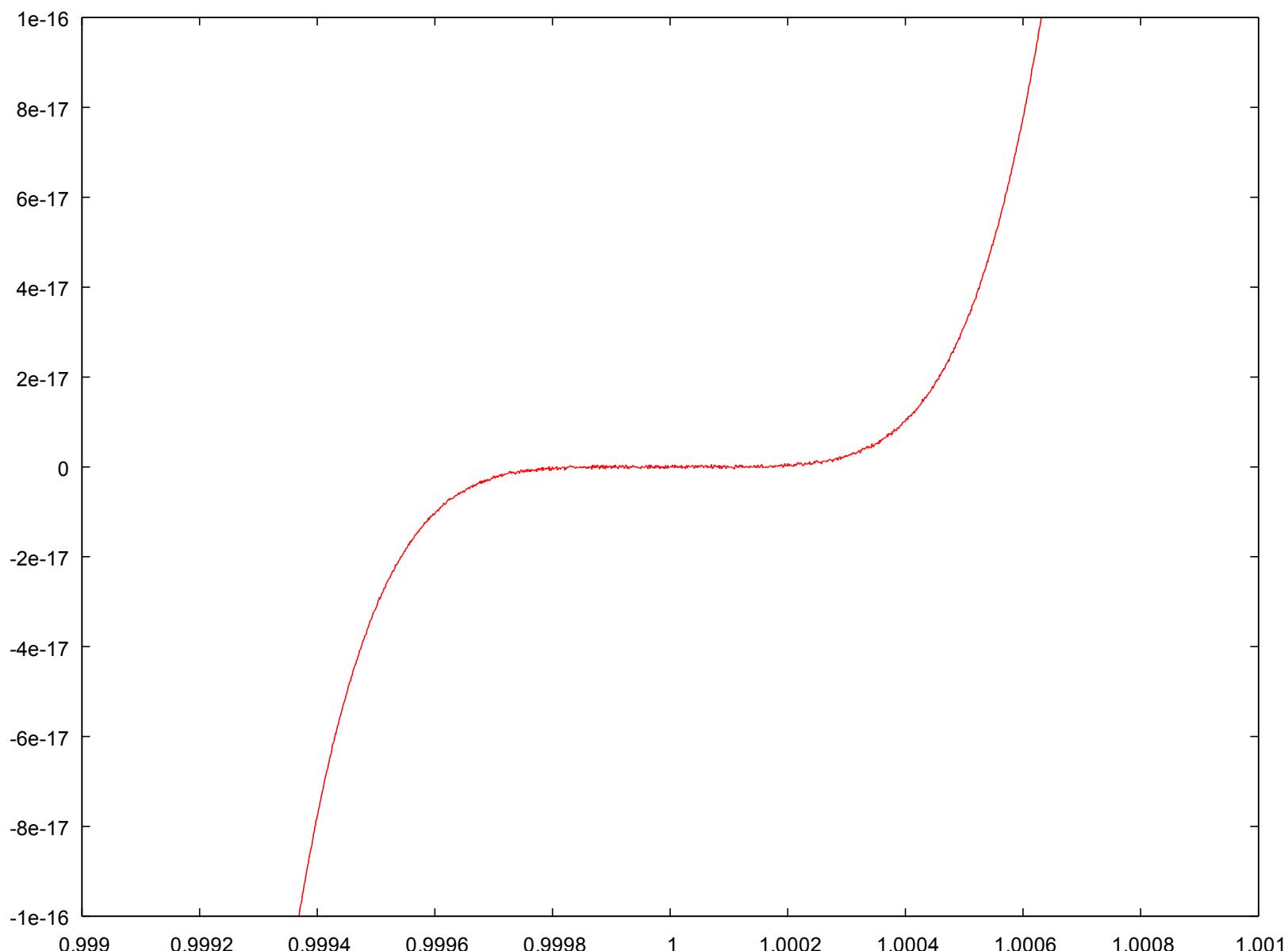
Of course, we actually evaluate this as

$$y = -1 + x \cdot (5 + x \cdot (-10 + x \cdot (10 + x \cdot (-5 + x))))$$

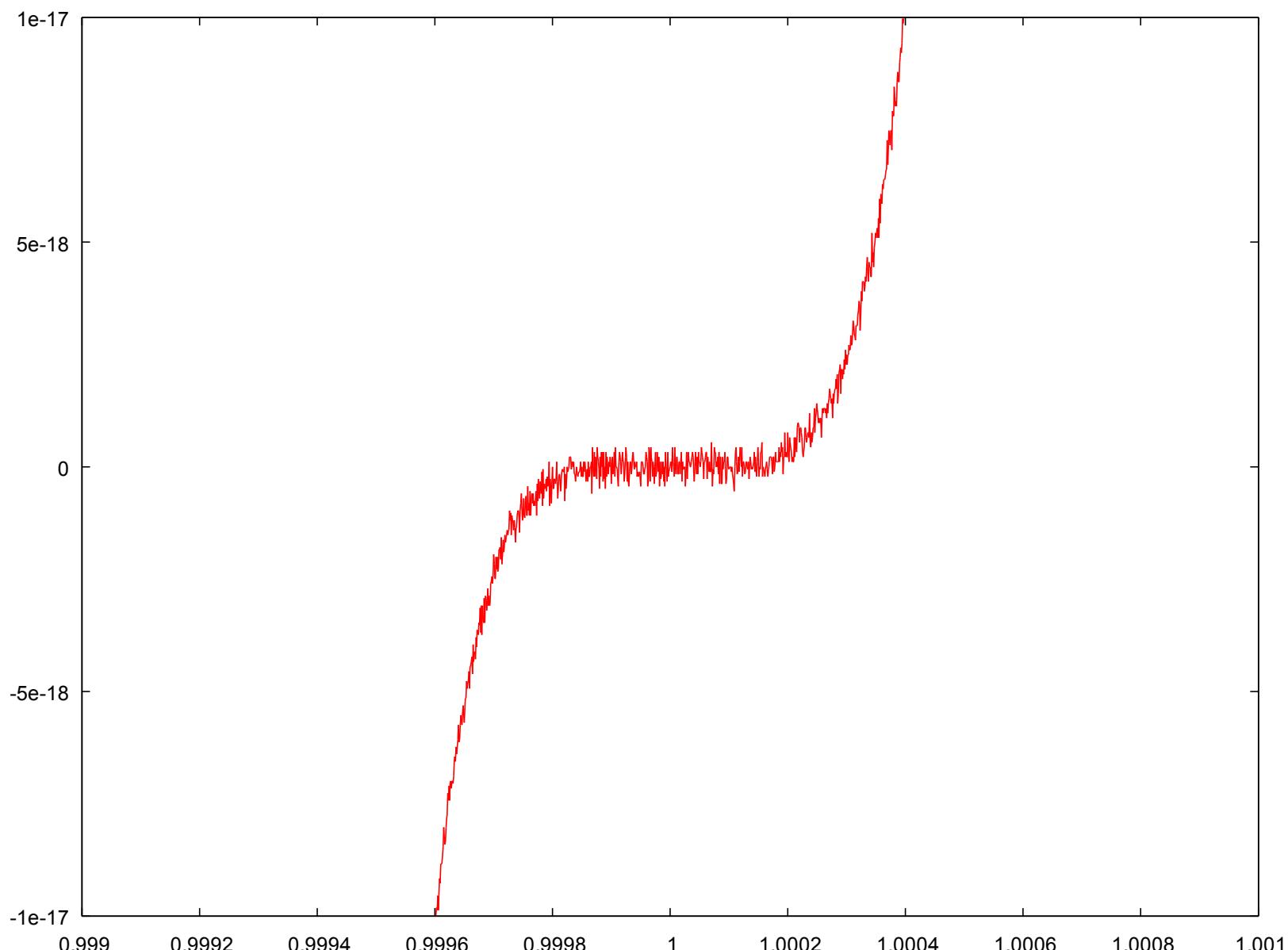
Region between $x=0.999$ and $x=1.001$



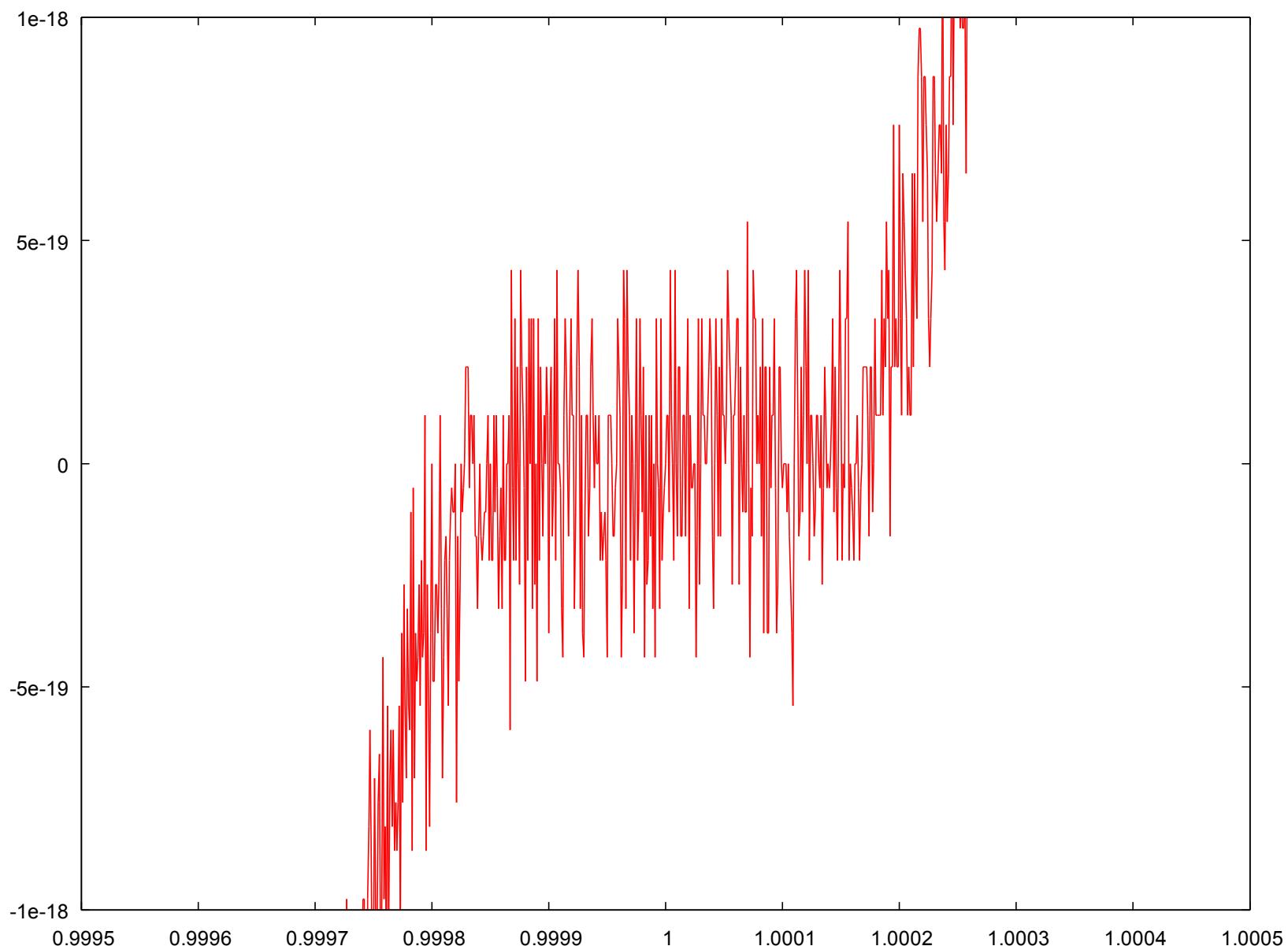
Zooming in on y-axis



Zooming in on y-axis



Zooming in on x-axis and y-axis



Raw Data from C Code

```
...
0.9999 -3.7947076e-19
0.999901 1.0842022e-19
0.999902 2.1684043e-19
0.999903 -1.6263033e-19
0.999904 -5.4210109e-20
0.999905 3.2526065e-19
0.999906 -2.1684043e-19
0.999907 4.3368087e-19
0.999908 0
0.999909 0
0.99991 -5.4210109e-20
0.999911 -3.2526065e-19
0.999912 -4.3368087e-19
0.999913 1.0842022e-19
0.999914 3.2526065e-19
0.999915 2.1684043e-19
0.999916 0
0.999917 -1.6263033e-19
0.999918 1.0842022e-19
0.999919 3.2526065e-19
0.99992 1.0842022e-19
...

```



Numerical “noise”
from accumulated
roundoff errors