# 64-bit representation of a real number (long real)

1985's IEEE (Institute for Electrical and Electronic Engineers) report specifies 64-bit (binary digit) representation for a real number.

The first bit is a sign indicator: s is 0 for a positive and 1 for a negative number.

This is followed by 11-bit exponent (called a characteristic): c, and 52-bit binary fraction,

$f$ , called the mantissa.

Here, the base of exponent is 2.

52 binary digits correspond to between 16 and 17 decimal digits; thus we can expect at least 16 decimal digits precision. The exponent of 11 binary digits corresponds to a range between

```
> 0*2^10+0*2^9+0*2^8+0*2^7+0*2^6+0*2^5+0*2^4+0*2^3+0*2^2+0*2^1+0*
  2^0;
```

$$0 \tag{1}$$

and

```
> 1*2^10+1*2^9+1*2^8+1*2^7+1*2^6+1*2^5+1*2^4+1*2^3+1*2^2+1*2^1+1*
  2^0;
```

$$2047 \tag{2}$$

Also note that

```
> 2^11-1;
```

$$2047 \tag{3}$$

However, we also want a good representation for numbers with small magnitudes. To insure that 1023 is substracted from the characteristic, so the range of exponents varies from -1023 to 1024. Using a normalization for a unique representation, the floating-point number is of the form

$$(-1)^s \, 2^{c-1023} \, (1+f)$$

For the machine number

0 10000000011
1011100100010000000000000000000000000000000000000000

```
> s:=0;
```
$$s := 0 \qquad\qquad (4)$$

so the number is positive. The characteristic

```
> c:=1*2^10+0*2^9+0*2^8+0*2^7+0*2^0+0*2^5+0*2^4+0*2^3+0*2^2+1*
  2^1+1*2^0;
```
$$c := 1027 \qquad\qquad (5)$$

The exponential part is

```
> 2^(c-1023);
```
$$16 \qquad\qquad (6)$$

and

```
> f:=1*(1/2)^(1)+1*(1/2)^3+1*(1/2)^4+1*(1/2)^5+1*(1/2)^8+1*(1/2)
  ^12;
```
$$f := \frac{2961}{4096} \qquad\qquad (7)$$

so the above machine number has the decimal representation

```
> (-1)^s*2^(c-1023)*(1+f);
```
$$\frac{7057}{256} \qquad\qquad (8)$$

```
> evalf(%);
```
$$27.56640625 \qquad\qquad (9)$$

The next smallest machine number, NS, is

         0 10000000011
10111001000011111111111111111111111111111111111111111

and the next largest machine number, NL, is

         0 10000000011
10111001000100000000000000000000000000000000000000001

The mantissa of the former is

```
> f1:=1*(1/2)^(1)+1*(1/2)^3+1*(1/2)^4+1*(1/2)^5+1*(1/2)^8+1*(1/2)
  ^13+1*(1/2)^14+1*(1/2)^15+1*(1/2)^16+1*(1/2)^17+1*(1/2)^18+1*
```

```
   (1/2)^19+1*(1/2)^20+1*(1/2)^21+1*(1/2)^22+1*(1/2)^23+1*(1/2)
   ^24+1*(1/2)^25+1*(1/2)^26+1*(1/2)^27+1*(1/2)^28+1*(1/2)^29+1*
   (1/2)^30+1*(1/2)^31+1*(1/2)^32+1*(1/2)^33+1*(1/2)^34+1*(1/2)
   ^35+1*(1/2)^36+1*(1/2)^37+1*(1/2)^38+1*(1/2)^39+1*(1/2)^40+1*
   (1/2)^41+1*(1/2)^42+1*(1/2)^43+1*(1/2)^44+1*(1/2)^45+1*(1/2)
   ^46+1*(1/2)^47+1*(1/2)^48+1*(1/2)^49+1*(1/2)^50+1*(1/2)^51+1*
   (1/2)^52;
```

$$f1 := \frac{3255653929844735}{4503599627370496} \tag{10}$$

while the mantissa of latter is

```
> f2:=1*(1/2)^(1)+1*(1/2)^3+1*(1/2)^4+1*(1/2)^5+1*(1/2)^8+1*(1/2)
  ^12+1*(1/2)^52;
```

$$f2 := \frac{3255653929844737}{4503599627370496} \tag{11}$$

The decimal representation of the next smallest machine number is

```
> NS:=(-1)^s*2^(c-1023)*(1+f1);
```

$$NS := \frac{7759253557215231}{281474976710656} \tag{12}$$

```
> Digits:=50;
```

$$Digits := 50 \tag{13}$$

```
> evalf(NS);
```

$$27.566406249999996447286321199499070644378662109375 \tag{14}$$

The decimal representation of the next largest machine number is

```
> NL:=(-1)^s*2^(c-1023)*(1+f2);
```

$$NL := \frac{7759253557215233}{281474976710656} \tag{15}$$

```
> evalf(%);
```

$$27.566406250000003552713678800500929355621337890625 \tag{16}$$

```
> NL-NS;
```

$$\frac{1}{140737488355328} \tag{17}$$

```
> evalf(%);
```

$$7.1054273576010018587112426757812500000000000000000 \times 10^{-15} \tag{18}$$

The smallest positive number is

0 00000000000
00000000000000000000000000000000000000000000000000001

```
> smallest:=2^(-1023)*(1+2^(-52));
```

$smallest := 4503599627370497 \Big/$                                                      **(19)**

    40480450661462123670499069343783461409911329952828423671380271605\
    48606791359906937839207674028742489903741557286336238227796174747\
    71586953734026799881477019843034885531327227289338315484186432682\
    47953535694549013712401496684938539723620671129831911268162011302\
    4717539104666829230461005064372655017292012526615415482186989568

```
> Digits:=800;
```

$$Digits := 800$$                                                      **(20)**

```
> evalf(smallest);
```

1.112536929253600938577939279289474120394004424341852097806565015605\ **(21)**
    98443019980034826489521461063144293195185068351409540085856480363\
    55955177563613706587576099527870021569402283901616688769940819588\
    69366447962303711463505653902690669854082668064822372476489476060\
    94895453919262916829509258948093999425752137573913180380108402181\
    10155646350227416036982422856551956112981967804161220311948453638\
    44509335852727463758699329466240520086911416412562289753282786909\
    5519004323455807914359948543236466107248984857733533601998959954\
    04580173812990192997712369923839430590047536255771881194801858107\
    58586490800577230217976564216270322096932266245269456889784045790\
    23962025496137069271374713102713202044199184595937090864938966701\
    39621383772282614543769341253209859132766723632812500000000000000\
    00000000000000000 $\times 10^{-308}$

which is approximately                                     $10^{-308}$

The largest positive machine number is
                    0 11111111111
1111111111111111111111111111111111111111111111111111

```
> largest:=2^1024*(1-2^(-52));
```

$largest :=$                                                      **(22)**
    17976931348623155085612432838450624023434343715745933592440487244\
    85818457545561143884706399431262203219608040271573715708098528849\
    64511743044087662767600909594331927728237078876188760579532563768\
    6986540648252621157710157914639830148577040812341945938624514172\
    3703148097529108423358883457665451722744025579520

```
> evalf(%);
```

1.797693134862315508561243283845062402343434371574593359244048724485\ **(23)**

$$81845754556114388470639943126220321960804027157371570809852884964\backslash$$
$$51174304408766276760009059433192772823707887618876057953256376869\backslash$$
$$86540648252621157710157914639830148577040081234194593862451417237\backslash$$
$$03148097529108423358883457665451722744025579520 \times 10^{308}$$

which is approximately

$$10^{308}$$

Numbers occurring in calculations that are smaller than

$$2^{-1023}\left(1+2^{-52}\right)$$

are treated as zero.

Floating-point representations in Maple are easy done. For example,

```
> Digits:=10;
```
$$\textit{Digits} := 10 \tag{24}$$

causes all arithmetic to be rounded to 100 digits. For instance, *fl(fl(x)+fl(y))* is performed using 100-digit rounding arithmetic by

```
> evalf(evalf(x)+evalf(y));
```
$$x + y \tag{25}$$

Implementing t-digit chopping aritmetic is slightly more complicated.

```
> chop:=proc(x,t)
  local e, x2;
  if x=0 then 0
  else
  e:= trunc(evalf(log10(abs(x))));
  if e>0 then e:=e+1 fi;
  x2:=evalf(trunc(x*10^(t-e))*10^(e-t))
  fi
  end:
> chop(12.226,4);
```
$$12.22000000 \tag{26}$$

```
> Digits:=10;
```
$$\textit{Digits} := 10 \tag{27}$$

Solving a quadratic equation

```
> solve({x^2+62.10*x+1},{x});
```
$$\{x = -0.01610723741\}, \ \{x = -62.08389276\} \tag{28}$$

```
> Digits:=4;
```
$$\tag{(29)}$$

$$Digits := 4 \tag{29}$$

```
> sqrt((62.10)^2-4*1.0*1.0);
```

$$62.06 \tag{30}$$

```
> floatx1:=(-62.10+62.06)/2.0;
```

$$floatx1 := -0.02000 \tag{31}$$

has large relative error: 2.4x10^(-1)

```
> floatx2:=(-62.10-62.06)/2.0;
```

$$floatx2 := -62.10 \tag{32}$$

while floatx2 has the small relative error: 3.2x10^(-4).  In order to obtain more accurate approximation for floatx1 one can note that

$$\left( x1 = \frac{-b + \sqrt{b^2 - 4\,a\,c}}{2\,a} \right) = -\frac{2\,c}{b + \sqrt{b^2 - 4\,a\,c}}$$

Using this formula floatx1 is given by

```
> floatx1:=-2.0/(62.10+62.06);
```

$$floatx1 := -0.01610 \tag{33}$$

which has the small relative error: 6.2x10^(-4).