

## IEEE 754: What it means for your computer, and humanity

The textbook<sup>1</sup> has an idealized view of floating point, which I think is wise. But, in this separate document, I think it is a good idea to lay out the basics of how the *real* standard is *actually* implemented on a computer.

- Computer memories are organized into *bytes*, that is, groups of 8 *bits*.<sup>2</sup> A bit is the irreducible atom of memory, which is in either of two states  $\{0, 1\}$ . You should know that *integers* are represented on computers using 1, 2, 4, or 8 bytes, thus 8, 16, 32, or 64 bits.
- The IEEE 754 standard is about how *real* numbers are approximately represented in memory, that is, how *floating point* numbers are represented. “Floating point” is essentially just scientific notation, but using only finitely-many bits and thus representing only a finite subset of real numbers.
- The two best-known floating point representations use 32 (“single”) and 64 (“double”) bits. In `single`, the number

$$x = (-1)^s \times (1.d_1d_2d_3 \dots d_{23})_2 \times 2^{(e_1 \dots e_8)_2 - 127}$$

is represented by 32 bits this way:

<i>s</i>	<i>e</i> <sub>1</sub>	<i>e</i> <sub>2</sub>	<i>e</i> <sub>3</sub>	<i>e</i> <sub>4</sub>	<i>e</i> <sub>5</sub>	<i>e</i> <sub>6</sub>	<i>e</i> <sub>7</sub>	<i>e</i> <sub>8</sub>	<i>d</i> <sub>1</sub>	<i>d</i> <sub>2</sub>	<i>d</i> <sub>3</sub>	<i>d</i> <sub>4</sub>	<i>d</i> <sub>5</sub>	<i>d</i> <sub>6</sub>	<i>d</i> <sub>7</sub>	<i>d</i> <sub>8</sub>	<i>d</i> <sub>9</sub>	<i>d</i> <sub>10</sub>	<i>d</i> <sub>11</sub>	<i>d</i> <sub>12</sub>	<i>d</i> <sub>13</sub>	<i>d</i> <sub>14</sub>	<i>d</i> <sub>15</sub>	<i>d</i> <sub>16</sub>	<i>d</i> <sub>17</sub>	<i>d</i> <sub>18</sub>	<i>d</i> <sub>19</sub>	<i>d</i> <sub>20</sub>	<i>d</i> <sub>21</sub>	<i>d</i> <sub>22</sub>	<i>d</i> <sub>23</sub>
----------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	------------------------	------------------------	------------------------	------------------------	------------------------	------------------------	------------------------	------------------------	------------------------	------------------------	------------------------	------------------------	------------------------	------------------------

In `double`, the number

$$x = (-1)^s \times (1.d_1d_2d_3 \dots d_{52})_2 \times 2^{(e_1 \dots e_{11})_2 - 1023}$$

is represented by 64 bits this way:

<i>s</i>	<i>e</i> <sub>1</sub>	<i>e</i> <sub>2</sub>	<i>e</i> <sub>3</sub>	<i>e</i> <sub>4</sub>	<i>e</i> <sub>5</sub>	<i>e</i> <sub>6</sub>	<i>e</i> <sub>7</sub>	<i>e</i> <sub>8</sub>	<i>e</i> <sub>9</sub>	<i>e</i> <sub>10</sub>	<i>e</i> <sub>11</sub>	<i>d</i> <sub>1</sub>	<i>d</i> <sub>2</sub>	<i>d</i> <sub>3</sub>	<i>d</i> <sub>4</sub>	<i>d</i> <sub>5</sub>	<i>d</i> <sub>6</sub>	<i>d</i> <sub>7</sub>	<i>d</i> <sub>8</sub>	<i>d</i> <sub>9</sub>	<i>d</i> <sub>10</sub>	...	<i>d</i> <sub>51</sub>	<i>d</i> <sub>52</sub>
----------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	------------------------	------------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	-----------------------	------------------------	-----	------------------------	------------------------

- Note that the “1.” in the above `single` and `double` representations, which appears before the *d<sub>i</sub>* bits, is always there. Thus it is not represented! It is called the “implicit leading bit”.
- The IEEE 754 standard is slightly more abstract than the concrete way the bits are arranged above. The standard says that every representable *nonzero* number is of the form

$$(1) \quad x = (-1)^s \times \frac{m}{\beta^{t-1}} \times \beta^e$$

for fixed positive integers  $\beta$  (the *base*) and  $t$  (the *precision*). The other numbers, namely  $s \in \{0, 1\}$  (the *sign*), the integer  $m$  (the *mantissa*), and the integer  $e$  (the *exponent*), depend on, and determine,  $x$ . These satisfy

$$(2) \quad \begin{aligned} \beta^{t-1} &\leq m \leq \beta^t - 1, \\ e_{\min} &\leq e \leq e_{\max}. \end{aligned}$$

<sup>1</sup>L. Trefethen and D. Bau, *Numerical Linear Algebra*, SIAM Press, 1997.

<sup>2</sup>“bit” = binary digit

- Unlike the system **F** in the textbook, in IEEE and other actual floating-point representations there are only finitely-many allowed values of the exponent  $e$ , and thus only finitely-many representable floating point numbers.
- In IEEE 754 there are five basic formats, but two of these are (oddly enough) *decimal* standards with  $\beta = 10$ , and rarely used.
- The three IEEE 754 basic formats that matter most are *binary*, that is, they all have base  $\beta = 2$ . They use 32, 64, or 128 bits, respectively. We have already shown how the first two appear in memory. In terms of form (1) and constraints (2), they follow this table:

name	common name	precision $t$	exponent bits	exponent bias	$e_{min}$	$e_{max}$
binary32	single	24	8	$2^7 - 1 = 127$	-126	+127
binary64	double	53	11	$2^{10} - 1 = 1023$	-1022	+1023
binary128	quadruple	113	15	$2^{14} - 1 = 16383$	-16382	+16383

- If you convert the above binary description to decimal you get these heuristic values:

name	decimal precision	decimal $e_{max}$	decimal $e_{min}$
binary32	7.22	38.23	-37.93
binary64	15.95	307.95	-307.65
binary128	34.02	4931.77	-4931.47

- Regarding the exponent, if all bits  $e_i$  are zero or all are one then the number has special meaning. That is, for normal numbers in `single` the standard requires

$$(e_1 \dots e_8)_2 \in \{1, 2, \dots, 254\}$$

and in `double` the standard requires

$$(e_1 \dots e_{11})_2 \in \{1, 2, \dots, 2046\}.$$

- Representing the number zero, which is *not* in form (1), is an example of “special meaning.” It is done by setting all bits other than  $s$  to zero. Oddly enough, this means “+0” and “-0” both exist on our computers as separate representations.
- Also there are representations of  $+\infty$  and  $-\infty$ , of things that are “not a number” (“NaN”), and of “subnormal” numbers. For subnormal numbers, in the `single` representation  $(e_1 \dots e_8)_2 = 0$  but not all bits  $d_i$  are zero.
- The IEEE 754 standard also addresses the rounding errors which occur with operations (addition, multiplication, etc.). In essence, the goal is that axiom (13.7) applies. But this is beyond my scope.
- Note “IEEE” stands for “Institute of Electrical and Electronics Engineers”. For more on IEEE 754 see the wikipedia page

[en.wikipedia.org/wiki/IEEE\\_floating\\_point](https://en.wikipedia.org/wiki/IEEE_floating_point)