

5 Number Representations

There are two separate concepts that you need to know to understand how numbers are represented on a computer. First, the binary system. And second, the finite precision/range.

5.1 Binary system

Information on a computer is stored in two-state devices such as a flip-flop made out of two transistors. One of these devices can either be turned on or off. This is the smallest amount of information a computer can store. We call this a *bit*, which can be either be 0 or 1. As humans, we normally represent numbers by writing their digits in base 10 (we have 10 fingers). However, because computers work with transistors that are either on or off (but do not have ten different levels of *on*) it makes sense to describe numbers in base 2 on a computer. There is nothing special about writing a number in either base 2 or base 10.

Let's go over a few examples of converting numbers from decimal to binary and back. You don't need to become an expert in this, but you should be familiar with the issues this conversion can create.

Decimal	Binary
0	0
1	1
2	10
3	11
4	100
10	1010
100	1100100
0.5	0.1
0.25	0.01
0.75	0.11
0.1	0.0001001100110011...
0.33333...	0.101010101...

Code 25: Example of binary/decimal conversion.

Let's discuss the last two examples in more detail, which brings us to the issue of finite precision.

5.2 Finite precision

On a computer (also on paper) there is only a certain amount of space available. We want to make best use of the available space which is why there are different ways of storing a number on a computer (and on paper).

As you can see in the above example, some of the number cannot be written exactly on any finite piece of memory (or paper). One example is the number $1/3$. We had to come up with a completely new way of writing number, namely as a fraction, to be able to write it down. Some numbers are representable exactly in base 10, but not in base 2, one example is the number $1/10$. Note that computers do (in general) not understand what a fraction is. They use a representation called floating point (see below) as the representation for everything. There are some exceptions to this rule. First of all, you could write a program that understands fractions. Second, there are mathematical packages such as Maple and Mathematica that already come with this functionality. We will not cover those in this lecture.

Now, all types of information on a computer such as numbers, texts, sounds and videos are stored as a collection of bits. Often 8 bits are grouped together and called a byte. The following figure shows

other commonly used names for a certain number of bits. You’ve probably heard *bit* and *byte* before, but maybe not *word* or *nibble*.

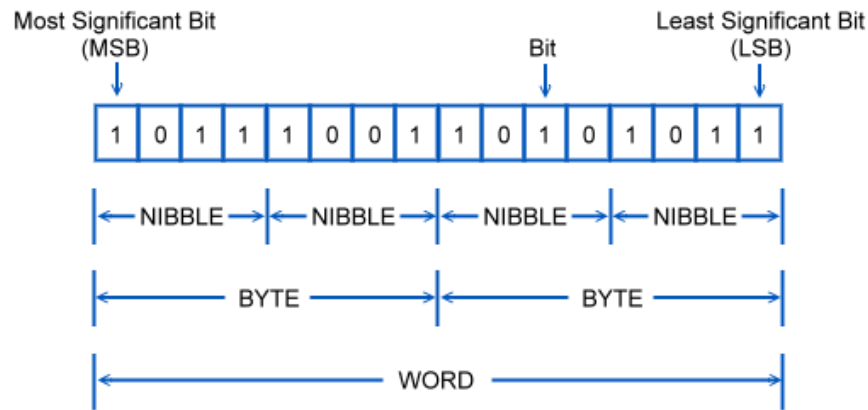


Figure 5: Bits and bytes. Source: <http://www.pcldev.com>.

Most importantly for us, each number on a computer is represented by a finite number of bits (or bytes). This has dramatic consequences for the algorithms we will implement. For that reason, we will discuss in quite some detail how numbers are stored.

5.3 Integers

Integers are one of the most fundamental datatypes in computer science. Typically an integer consists of at least two bytes on a modern computer. With two bytes (16 bits) we can store numbers from 0 (including) up to

$$2^{16} - 1 = 65535.$$

However, often we also want to be able to store negative integers. We need to reserve one bit for the sign of the integer, thus leaving us 15 bits for the magnitude. The range of a signed two-byte integer is therefore

$$-2^{15} = -32768 \dots 32767 = 2^{15} - 1.$$

Let’s look at what this finite range of possible integers means in practice. If we start with an integer set to 0 and add 1 many times, eventually we will reach the maximum. When this happens, we will get a negative number! In python you can find out what the maximum integer is that it can represent by using the `sys` module. However, python is pretty smart. It automatically increases the number of bytes for the integer if you try to add 1 to the maximum allowed value for the integer. In that case, the integers becomes what is called a *long integer*, or just *long*. In many other programming languages, this is not done automatically for you and you need to be extremely careful about having integers which are larger then the maximum allowed size.

An example code for python is given in Code 26. Remember, python does increase the number of bytes automatically. The same code in C, would give you a negative number.

```
>>> import sys
>>> sys.maxint
9223372036854775807
>>> sys.maxint + 1
9223372036854775808L
```

Code 26: Interactive python session testing the maximum size of int.

5.4 Floats

Most numbers in scientific calculations will be represented by so something called a *float* or a *floating point number* on the computers that we use. The name comes from the fact that the decimal point will shift (float) and is not fixed. You can consider the integer datatype as a fixed precision datatype (i.e. no information about the digits after the decimal point is stored). In python and almost all other programming languages, the industry standard of representing such a number is IEEE-754. We use the *double precision* version of this standard. It used 8 bytes, 64 bits. On some special hardware, for example GPUs, a datatype with less precision, *single precision* is used. So how do we represent a number x where the decimal point can shift? We describe the number as a combination of three other numbers:

$$x = \text{sign} \times \text{mantissa} \times 2^{\text{exponent}}$$

The sign is just a single bit, either 0 or 1, describing whether the number is positive or negative. The mantissa (sometimes called fraction) consists of 52 bits. These are the digits after the decimal point. The exponent contains 11 bits, which is just an integer. Putting this all together we have

$$x = (-1)^{\text{sign}} \left(1 + \sum_{i=1}^{52} b_{52-i} 2^{-i} \right) 2^{e-1023}$$

This is crucial to understand how a floating point number is constructed. It has severe limitations for both the precision and the range of possible numbers.

The following figure illustrates how a double floating point number is stored in memory.

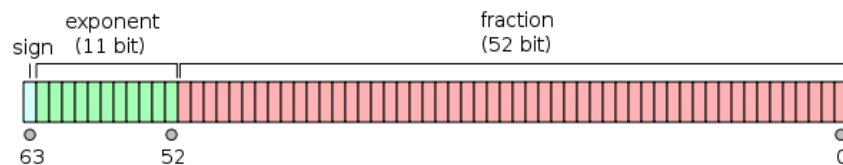


Figure 6: Double floating point number in memory. Source: <http://en.wikipedia.org>.

Let's look at the precision first. There are only 52 bits available. Thus the binary representation of a number is truncated after 52 binary digits. This corresponds to approximately 16 decimal digits. Any digits after that cannot be represented and are simply ignored. Let's look at $1/10$ and $2/10$. Both of these numbers cannot be represented exactly as a binary fraction (see above). For that reason, the addition of the two numbers will also not exactly give us $3/10$ (which can not be represented exactly in the first place). Have a look at the python snippet in Code 27. The first line loads an external module that allows you to print the decimal number very accurately on the screen.

```

>>> from decimal import Decimal
>>> Decimal(0.1)
Decimal('0.1000000000000000000555111512312578270211815')
>>> Decimal(0.2)
Decimal('0.20000000000000000001110223024625156540423631')
>>> Decimal(0.2+0.1)
Decimal('0.30000000000000000004440892098500626161694526')
>>> Decimal(0.3)
Decimal('0.29999999999999999998889776975374843459576368')

```

Code 27: Python listing testing the floating point precision.

Now let's look at the finite range. Let's start with the number $x = 1$ and multiply it with the number 2 multiple times. Eventually we will run over the maximum range of what double precision floating point numbers can represent. This happens when we run out of bits in the exponent. Let's see when this occurs in python. A simple test is implemented in Code 29. The code outputs the largest number before we reach infinity. Here, infinite means that the computer has noticed that we ran above the maximum range. It's not actually infinity, just a number larger than what the computer can represent. In python (using the default double precision) we get a value around 10^{307} . Note that python outputs this in *scientific notation* with a syntax involving `..e+...`. Can you explain why we get the number 307 (roughly)?

```

>>> x = 1.
>>> while x < x*2.:
...     x = 2.*x
...     print x
...
2.0
4.0
8.0
16.0
[...]
4.49423283716e+307
8.98846567431e+307
inf

```

Code 28: Python listing testing the floating point maximum.

5.5 Compensated summation

As we found out earlier, one problem with floating point numbers is the limited precision. One could simply use a data format with more bits. The buzzword for this is extended (or even arbitrary) precision. The downside to this is that one needs to either build a new computer, which is expensive, or implement it in software, which is slow.

There is another way to solve the problem, at least in precisions. This solution only uses floating point numbers. It will be our first algorithm that we discuss. It is known as compensated summation, or Kahan summation.

Let's look at the problem first. Adding two numbers works well, only if the numbers are of similar size. Otherwise we get large errors.

```

>>> print "%.20f" % (1. + .1)
1.100000000000000000008882
>>> print "%.20f" % (1e15 + 0.1)
1000000000000000.12500000000000000000
>>> print "%.20f" % (1e15 + 0.1 - 1e15)
0.125000000000000000000000

```

Code 29: Python listing showing the problem using non-compensated summation.

Don't worry about the complicated looking print statement. It's only purpose is to show you the number with 20 decimal digits (instead of 5 which is the default). What you see is that a huge problem exists when we add large numbers to small numbers. The equation

$$(10^{15} + 0.1) - 10^{15} = 0.1$$

is only correct to within 25%. Suppose you have 10 cents in your bank account. If the bank accidentally deposits 10^{16} dollars into your account, notices the change and then withdraws the money again, you would end up with 12.5 cents. So you would have made a 2.5 cent profit. With compensated summation we can avoid this issue. Here is how the algorithm works.

Let y_0 be the initial number and δ_n a series of number that we want to add to y_0 . We need one additional variable e , also a floating point number. Initially we set $e = 0$. Then, in pseudo code:

```

for n=0,1,2,... do
     $a = y_n$ 
     $e = e + \delta_n$ 
     $y_{n+1} = a + e$ 
     $e = e + (a - y_{n+1})$ 
end for

```

If these were all mathematical statements as you know them, we could simplify all the above to simply get

```

for n=0,1,2,... do
     $y_{n+1} = y_n + \delta_n$ 
end for

```

which is of course just the simple summation. Let's look at what happens to the variables in the compensated summation algorithm in detail.

$a = y_n$	a'	a''		
e			e'	0
δ_n		δ'	δ''	
$e = e + \delta_n$		δ'	$e' + \delta''$	
$y_{n+1} = a + e$	a'	$a'' + \delta'$		
$e = e + (a - y_{n+1})$			$e' + \delta''$	0

Using only standard summation, the last digits (δ'') would be lost. Here, we carry them along!

You will have to implement this algorithm for the next assignment. To make things worse, you'll be asked to implement this in assembler. Let's start with an algorithm for non-compensated summation.

You can use three new assembler commands. The following table lists them.

Command	Argument 1	Argument 2	Argument 3	Result
INPUT	r0			Input value from user and store it in r0
SUB	r0	r1	r2	$r2 := r0 - r1$
COPY	r0	r1		Copy the value of r0 into r1

```
SET 0 r0
SET -3 r9
INPUT r1
ADD r0 r1 r0
PRINT r0
JUMP r9
```

Code 30: Non-compensated summation in assembler.

To give you some help, you may use of the registers as follows:

```
y = r0
d = r1
a = r2
e = r3
temporary value = r4
```

End of lecture 3.
