

CHAPTER 1

INTRODUCTION TO ROUND-OFF ERRORS

Round-off errors can lower the stability and accuracy of a computer simulation. This lecture aims to show where the round-off errors come from and how to minimize the impact of the round-off errors in the numerical simulation. To understand the source of the round-off errors, we need to know the computer architecture and the computer representations of the integer numbers and the real numbers.

1.1. Review of Word, Byte, and Bit in the Computer Architecture

In the modern computer processor, a byte consists of 8 bits. A word consists of 4 bytes for a 32-bit computer or 8 bytes for a 64-bit computer. Thus, the word's size is 32 bits in the 32-bit computer, but 64 bits for a 64-bit computer.

Exercise 1.1.

To learn more about the historical development of computer architecture, please read the information on “word in computer architecture” in Wikipedia,

[http://en.wikipedia.org/wiki/Word_\(computer_architecture\)](http://en.wikipedia.org/wiki/Word_(computer_architecture))

From Exercise 1.1, we can see that the most popular sizes of a word found in different processors are 64, 32, 16, 8, and 4 bits. But one can also find processors with the size of a word equal to 60, 50, 48, 40, 39, 36, 34, 27, 26, 25, 22, 18, 15, 12, or 9 bits. The size of a byte is equal to the size of a character. The size of a byte is 8 bits in most of the processors. However, there were processors in the early days with a byte equal to 5, 6, or 9 bits.

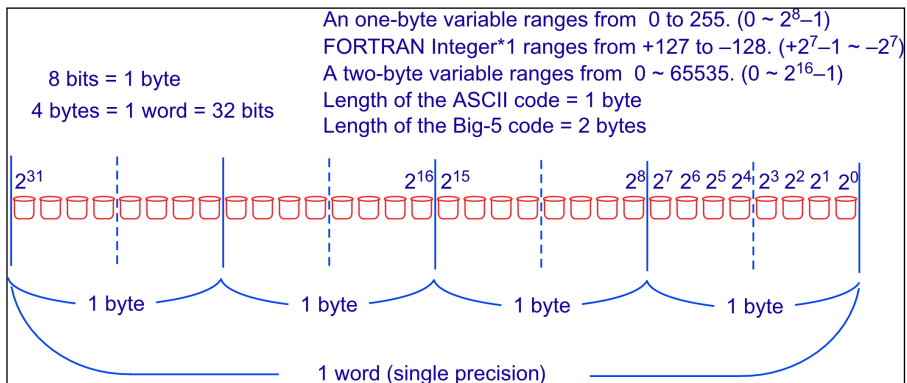


Figure 1. The structure of word, bytes, and bits in a 32-bit computer

1.2. Computer Representations of Integer Numbers

Figure 1 shows the relationship among a single precision word, byte, bit in a 32-bit computer. For the modern 32-bit computer, a single precision word consists of 4 bytes. A byte consists of 8 bits. Thus, a single precision word consists of 32 bits.

From Figure 1, we can also conclude that, for the 32-bit computer, the maximum integer register is $2^{32} = 4294967296 \sim 4G$. But the range of integer depends on how the processor treats the negative integer. For the modern computer processor, the computer representation of integer is given in the following way:

- When the first bit is 0, the integer is positive. The absolute value of the positive integer is determined by the binary representation of the rest 31 bits.
- When the first bit is 1, the integer is negative. The absolute value of the negative integer is determined by the binary representation of the complement of the rest 31 bits plus 1.

Exercise 1.2.

Read the information on “signed number representations” in Wikipedia, http://en.wikipedia.org/wiki/Signed_number_representations

Exercise 1.3.

Read the information on “integer in computer science” in Wikipedia, [http://en.wikipedia.org/wiki/Integer_\(computer_science\)](http://en.wikipedia.org/wiki/Integer_(computer_science))

Q: What is $(10011000)_2$?

The first bit is 1. It means it is a negative integer. Changing the rest bits from 0011000 to 1100111, it yields $(1100111)_2 = 64 + 32 + 0 + 0 + 4 + 2 + 1 = 103$.

The absolute value of the negative integer is $103 + 1 = 104$. Thus, $(10011000)_2 = -104$ (When the author was a graduate student in NCU, the university has purchased a 60-bits CDC-Cyber computer. In order to speed up the calculation, this computer taking the complement of the rest 59 bits without adding 1 when it evaluates the negative integers. As a result, both $(000 \dots 000)_2$ and $(111 \dots 111)_2$ equal to 0 in the 60-bits CDC-Cyber computer.)

1.3. Computer Representations of Integers and Characters

The integers and characters are connected through a given code table. Examples of such code include the ASCII code for English and the Big-5 code for Traditional Chinese. One may find intrinsic functions to make converting between the integers and ASCII code. Or, one can use “A format” in the FORTRAN language to convert integers and “characteristics” (or sometimes called “strings”). However, it is even more important to know that we can output our integer or floating-point data in “A format” as binary data to save a lot of disk space. Note that after the year 2000, the formatted binary data has become a machine-independent data structure. Thus, most computer graphic software, such as IDL or MATLAB, can read cross-platform binary data without a problem.

Exercise 1.4

Write a test program to determine the binary structures of the ASCII code, the integers, and the floating numbers used in your computer. (e.g., Write 31~127 to a file in A1 format or a binary format. Then, open the file and find out what you can see in that file.)

The following are two FORTRAN programs for ASCII code and Chinese Big-5 code

```
C==GETASCII.f=====
```

```
C This program shows the binary structures of the ASCII code
```

```

Program GETASCII
DO I=31, 127
WRITE(3,20) I, I
ENDDO
20 FORMAT(1X,I3,1X,A1)
STOP
END

```

```
C==GETBIG5.f=====
```

C This program outputs the Chinese characters in Big-5 code.
C To view the Chinese characters, you can open the output file by
a web browser
C and choose viewing text by Big-5 encoding. Then, you can make
C a copy of the Chinese characteristics and past them to a regular
word document.

```
C
```

```

PROGRAM GETBIG5
C HIGH: A1-F9 (161-249)
C LOW: 40-7E (64-126), A1-FE (161-254)
INTEGER*1 IA(20000),IB(20000)
JJ=0
DO I=161,249
DO II=64,126
JJ=JJ+1
IB(JJ)=II
IA(JJ)=I
ENDDO
DO II=161,254
JJ=JJ+1
IB(JJ)=II
IA(JJ)=I
ENDDO
ENDDO
JJ0=JJ
WRITE(11,1) (IA(K),IB(K),K=1,JJ0)
1 FORMAT(100A1)
STOP
END

```

Additional example of FORTRAN program to determine the ASCII code of a given character and vice versa.

```
C==ASCII_TEST.f=====
C This program determines the ASCII code of a given character and
vice versa
  program ascii_test
  character*1 a
  byte i
  id=1 !id can be any integer between 1 and 99 except 5 and 6
      !id=5 is reserved for the system_input, such as the
terminal
      !id=6 is reserved for the system_output, such as the
terminal
  10 continue
  print *, 'enter one character'
  read(5,*) a
  write(id,1)a
  1 format(A1)
  rewind id
  read(id,1)I
  print *, 'I=',I
  rewind id
  write(id,2)a
  2 format(A2)
  rewind id
  read(id,2)I
  print *, 'I=',I
C
  print *, 'enter an integer, or enter 0 to stop'
  read(5,*) I
  if(i.eq.0) go to 99
  rewind id
  write(id,1)I
  rewind id
  read(id,1)a
  print *, 'a=',a
  rewind id
  write(id,2)I
  rewind id
  read(id,2)a
  print *, 'a=',a
  go to 10
  99 continue
```

```

stop
end

```

1.4. Computer Representations of Integers and Real Numbers

Table 1.1 shows the lower and upper limits of the integers and byte(s) at different length.

Table 1.1. The lower and upper limits of the integers and byte(s) at different length

Fortran Data Type	lower limit of the data	upper limit of the data
Byte	0	+ 255 (= $2^8 - 1$)
Integer*1	- 128 (= $- 2^7$)	+ 127 (= $2^7 - 1$)
Bytes	0	+ 65535 (= $2^{16} - 1$)
Integer*2	-32768 (= $- 2^{15}$)	+ 32767 (= $2^{15} - 1$)
Integer*4	- 2147483648 (= $- 2^{31}$)	~ 2147483647 (= $2^{31} - 1$)

Exercise 1.5.

To learn more about the historical development of the floating point in the computer architecture, please read the information on “floating point” on Wikipedia,

http://en.wikipedia.org/wiki/Floating_point

Exercise 1.6

Write a program to verify the results shown in Table 1.1.

Table 1.2 shows the computer representation of floating-point data at different lengths. The significant digits listed in Table 1.2 will give rise to round-off error. Thus, we must use double precision in our numerical simulation studies.

Note that there is no round-off error in the integer expression and calculation. However, the maximum and minimum of an integer expression are much less than the extrema of a floating-point expression at the same word length. Both real numbers and complex numbers are floating-point expressions with finite significant digits. Since a complex number consists of two real numbers, the size of a complex number is twice that of a real number.

Table 1.2. Computer representations of floating-point numbers at different lengths

Type	Sign	Exponent	Significand	Total bits	Exponent upper limit	significant digits
Half (IEEE 754-2008)	1	$5(= 1 + 4)$	10	16	$15(= 2^4 - 1)$	~3.3
Single	1	$8(= 1 + 7)$	23	32	$127(= 2^7 - 1)$	~7.2
Double	1	$11(= 1 + 10)$	52	64	$1023(= 2^{10} - 1)$	~15.9
Double extended (80-bit)	1	$15(= 1 + 14)$	64	80	$16383(= 2^{14} - 1)$	~19.2
Quad	1	$15(= 1 + 14)$	112	128	$16383(= 2^{14} - 1)$	~34.0

Exercise 1.7

- (a) Write a program to verify the last column shown in Table 1.2.
 (b) Write a program to check the value in the first three columns shown in Table 1.2 in your computer.

List below is an example of Fortran program, which determines the extrema of the real number and the integer number that can be resolved by the current computing system.

```

PROGRAM MAXIMUM_TEST
REAL*8 AA
1 CONTINUE
PRINT *, 'ENTER AA, A, I'
READ(5,*) AA, A, I
B=EXP(A)
IF(I.EQ.0) GO TO 99
PRINT *, 'AA, A, EXP(A), I ='
PRINT *, AA, A, B, I
GO TO 1
99 CONTINUE
STOP
END

```

1.5. How to Determine the Relative Error of a Floating-Point Expression

If U is the relative error of 1, then an iteration scheme is convergent when $|^{k+1}y^{n+1} - ^ky^{n+1}| < U * |^ky^{n+1}|$.

where $^ky^{n+1}$ is the k th iteration result of y^{n+1} .

The relative error U can be obtained from the program GETU.f as given below. (e.g., Shampine and Gordon, 1975). The relative error of a number A is $U * |A|$. The relative error is a machine-dependent error before year 2000. The reason that different floating-point processor has different relative error can be understood by the historical review given in Exercise 1.5.

```
C== GETU.f =====
C This subroutine determines machine-dependent relative error
C relative to 1.
  Subroutine GETU(U)
  Implicit double precision (a-h,o-z)
  A1=1.d0      !for double precision
  AH=0.5d0     !for double precision
C   A1=1.      !for single precision program
C   AH=0.5     !for single precision program
  U=A1
  UU=U
1 CONTINUE
  UU=UU*AH
  UT=U+UU
  IF(UT.GT.U) GO TO 1
  U=UU*2
  RETURN
  END
```

1.6. How to Minimized the Numerical Errors due to Round-off Error

We can minimize the numerical errors due to round-off errors by the following way (e.g., Tsai et al., 2009).

If $|A - B| < (\text{the relative error of } A \text{ and } B)$, then set $A - B = 0$

Examples of the function programs and main program that can demonstrate the round-off errors in the derivatives of $\tanh(x)$ are given below.

```
C==DIFAB.f=====
C This function determines A-B with minimized round-off error
```



```

FUNCTION DIFAB(A,B)
Implicit double precision (a-h,o-z)
common RELERR
TEMP=A-B
IF(DABS(TEMP).LT.MAX(DABS(A),DABS(B))*RELERR) TEMP=0
DIFAB=TEMP
RETURN
END

```

```

C==ADDAB.f=====
C This function determines A+B with minimized round-off error
FUNCTION ADDAB(A,B)
Implicit double precision (a-h,o-z)
common RELERR
TEMP=A+B
IF(DABS(TEMP).LT.MAX(DABS(A),DABS(B))*RELERR) TEMP=0
ADDAB=TEMP
RETURN
END

```

```

C==test_DIFAB.f=====
C This program determines the derivatives of the tanh(x)
program test_DIFAB
Implicit double precision (a-h,o-z)
common RELERR
call DGETU(RELERR)
RELERR=RELERR*10
A01=0.1d0
A0H=A01/2.d0
FL=1 ! exp(0)=1
DO I=1,180
X=I*A01
temp=dexp(x)
temp1=1.d0/temp
FR=(temp-temp1)/(temp+temp1)
PF1=(FR-FL)/A01
FR=DIFAB(temp,temp1)/ADDAB(temp,temp1)
PF2=DIFAB(FR,FL)/A01
IF(I.GT.160) THEN
WRITE(2,*) , 'X1, PF1, PF2='
WRITE(2,*) , X1, PF1, PF2

```

```

ENDIF
FL=FR
ENDDO
STOP
END
C=====
INCLUDE 'DIFAB.f'
INCLUDE 'ADDAB.f'
INCLUDE 'GETU.f'

```

1.7. Summary

Computer simulation is a special type of numerical method, which means to solve a system of ordinary differential equations (ODEs) or partial differential equations (PDEs) with a time-derivative term presented in each of the ODEs or PDEs. Simulation results can provide information on how and why the given system will evolve with time. The simulation results will depend on the governing equations, the simulation scheme, the grid size, the simulation domain's size, the time steps, the initial conditions, and the boundary conditions. Since most numerical simulations are time-consuming and memory-consuming, we need to do our best to save the real-time and find a balance between saving the CPU time and saving the memory. For better diagnostics of the simulation results, we need a lot of storage space to save the simulation results. Thus, we need to do our best to save disk space.

Exercise 1.7. Please discuss

- How to save the CPU time?
- How to save the real execution time?
- How to save the disk space?

The following examples are FORTRAN programs in which the first one will take much longer time to complete the execution in comparing with the second program.

```

C==TEST_IO_SLOW.f=====
C THIS PROGRAM SHOWS A BAD EXAMPLE OF I/O
PROGRAM TEST_IO_SLOW
PARAMETER(NCX=1000000)
DO I=1,NCX

```

```
    ARRAY=I*0.1D0
    WRITE(8,8) I, ARRAY
8  FORMAT(1X, I6, 1X, F15.7)
    ENDDO
    STOP
    END
```

```
C==TEST_IO_FAST.f=====
C THIS PROGRAM SHOWS A GOOD EXAMPLE OF I/O
PROGRAM TEST_IO_SLOW
PARAMETER(NCX=1000000)
DIMENSION ARRAY(NCX)
DO I=1,NCX
ARRAY(I)=I*0.1
ENDDO
WRITE(9,8) (I, ARRAY(I), I=1,NCX)
8  FORMAT(200A4)
STOP
END
```

Note that one can use the following command

```
time ./a.out
```

to find the real execution time in a Linux operating system (OS), where a.out is the execution file. An example of execution results is given below.

```
C==Real_CPU_Time.txt=====
Real Execution time and CPU time
```

```
$ gfortran TEST_IO_FAST.f
$ time ./a.out
```

```
real    0m0.197s
user    0m0.119s
sys     0m0.017s
```

```
$ gfortran TEST_IO_SLOW.f
$ time ./a.out
```

```
real    0m1.335s
user    0m1.257s
```

```
sys    0m0.037s
```

```
$ ls -l fort.*
```

```
-rw-r--r--  1 lyuling-hsiao  staff  24000000 Mar 25 19:52 fort.8  
-rw-r--r--  1 lyuling-hsiao  staff   8010000 Mar 25 19:52 fort.9
```

REFERENCES

Shampine, L. F., and M. K. Gordon (1975), *Computer Solution of Ordinary Differential Equation: the Initial Value Problem*, W. H. Freeman and Company, San Francisco.

Tsai, T. C., L. H. Lyu, J. K. Chao, M. Q. Chen, and W. H. Tsai (2009), A theoretical and simulation study of the contact discontinuities based on a Vlasov simulation code, *J. Geophys. Res.*, 114, A12103, doi:10.1029/2009JA014121.