

Lab 5

1. int

```
#include <stdio.h>

//unsigned int a = 4294967295;
int a = -2147483648;
int b[32];
int c;

void bit()
{
    printf("%d \n", a);
    for(int i = 0; i < 32; i ++)
    {
        asm(
            "movq $1, %rax \n"
            "andq a(%rip), %rax \n"
            "movq %rax, c(%rip) \n"
            "shrq $1, a(%rip) \n"
        );
        b[31 - i] = c;
    }
    for(int i = 0; i < 32; i ++)
    {
        printf("%d", b[i]);
    }
}
```

Examples

Number	Two's complement
-2147483648	10000000000000000000000000000000
-2147483646	100000000000000000000000000000010
-3	1111111111111111111111111111111101
2147483648(int)	100000000000000000000000000000000 (equals to -2147483648)
4294967295(unsigned int)	111111111111111111111111111111111
2147483648(unsigned int)	100000000000000000000000000000000

2. float

```
#include <stdio.h>

float a = 1.5;
//double a = 1.5;
int c;
int b[sizeof(a) * 8];

void bit()
{
    printf("%f \n", a);
    for(int i = 0; i < sizeof(a) * 8; i++)
    {
        asm(
            "movq $1, %rax \n"
            "andq a(%rip), %rax \n"
            "movq %rax, c(%rip) \n"
            "shrq $1, a(%rip) \n"
            );
        b[sizeof(a) * 8 - i - 1] = c;
    }
    for(int i = 0; i < sizeof(a) * 8; i++)
    {
        if(sizeof(a) == 4 && (i == 0 || i == 8)) printf("%d ", b[i]);
        else if(sizeof(a) == 8 && (i == 0 || i == 11)) printf("%d ", b[i]);
        else printf("%d", b[i]);
    }
}
```

Examples

Number	Two's complement
1.5(float)	0 01111111 100000000000000000000000
1.5(double)	0 0111111111 1000 0000000000000000
3.14(float)	0 10000000 10010001111010111000011 (=3.1400001049041748046875)
+(-)0 (float)	0(1) 00000000 000000000000000000000000
+(-) ∞	0(1) 11111111 000000000000000000000000
NaN	0 11111111 100000000000000000000000

3. Mantissa overflow

```
float d = 1073741824;  
d += 1.0;  
printf("%f", d);
```

output: 1073741824.0

4. Non-associativity

```
a*(b-c)!=a*b-a*c: -0.12999999803 != -0.12999999654  
a/e*b*c)!=a*b*c/e: 0.05199999973 != 0.05200000011
```

5. Assembly listing

```
float a = 3.1;
double b = 3.4;

int main()
{
    double c = b + 1.2;
    float d = a + 4.3;
}
```

- movsd - moves scalar double-precision floating-point value(b to %xmm1(first 64 bits are rewritten))
- movss – same as movsd, but performs with float values
- addsd – adds two double values
- cvtss2sd - converts one single-precision floating-point value to one double-precision floating-point value(%xmm0 to %xmm1 - first 32 bits of %xmm0 are rewritten as double in first 64 bits of %xmm1)

Interestingly, float and double numbers are written as huge numbers, e.g. float a is written as 10783555588. That is because floating-point number representation(in binary) is converted to decimal. As for doubles, it's value is divided in 2 rows: first row is for last 32 bits of double(mantissa), second row is for first 32 bits(sign + exp + mantissa)

```
a:
    .long    1078355558
    .globl   b
    .align   8
    .type    b, @object
    .size    b, 8

b:
    .long    858993459
    .long    -1073007821
    .text
    .globl   main
    .type    main, @function

main:
.LFB0:
    .cfi_startproc
    endbr64
    pushq    %rbp
    .cfi_def_cfa_offset 16
    .cfi_offset 6, -16
    movq     %rsp, %rbp
    .cfi_def_cfa_register 6
    movsd    b(%rip), %xmm1
    movsd    .LC0(%rip), %xmm0
    addsd    %xmm1, %xmm0
    movsd    %xmm0, -8(%rbp)
    movss    a(%rip), %xmm0
    cvtss2sd    %xmm0, %xmm1
    movsd    .LC1(%rip), %xmm0
    addsd    %xmm1, %xmm0
    cvtsd2ss    %xmm0, %xmm0
    movss    %xmm0, -12(%rbp)
    movl     $0, %eax
    popq     %rbp
    .cfi_def_cfa 7, 8
    ret
    .cfi_endproc

.LFE0:
    .size    main, .-main
    .section .rodata
    .align   8

.LC0:
    .long    858993459
    .long    1072902963
    .align   8

.LC1:
    .long    858993459
    .long    1074869043
```

```

#include <stdio.h>

float a = 3.1;
double b = -3.4;

int main()
{
    double c = b * 1.2;
    float d = a * 4.3f;
}

```

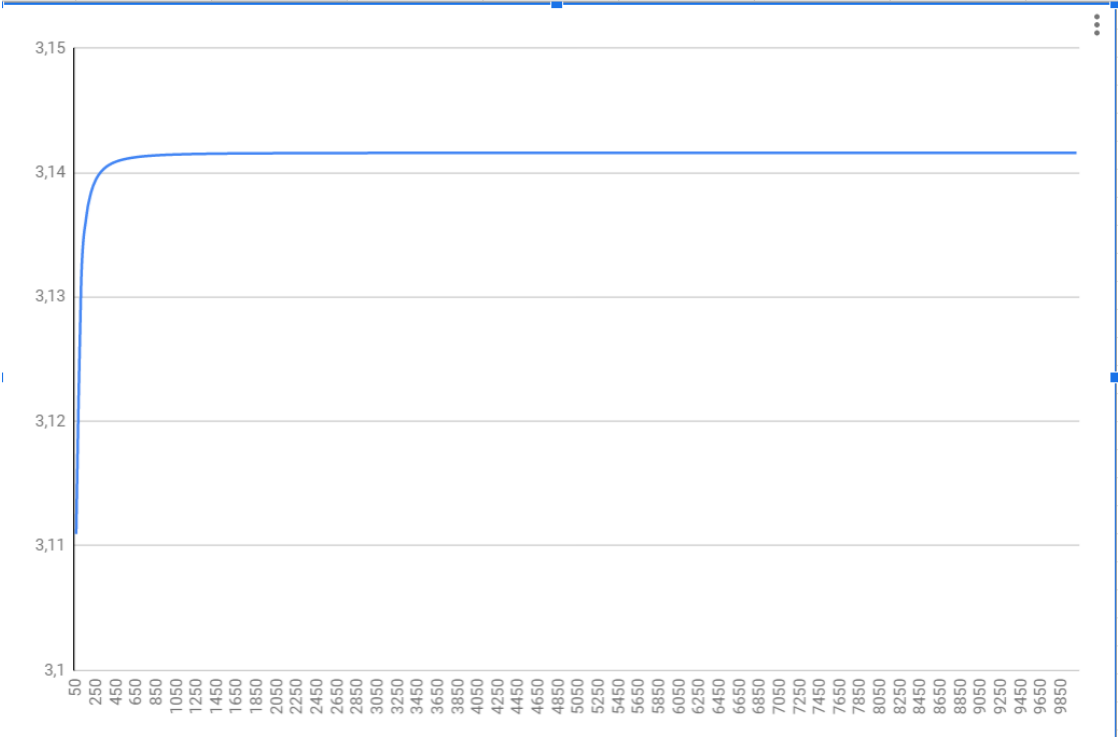
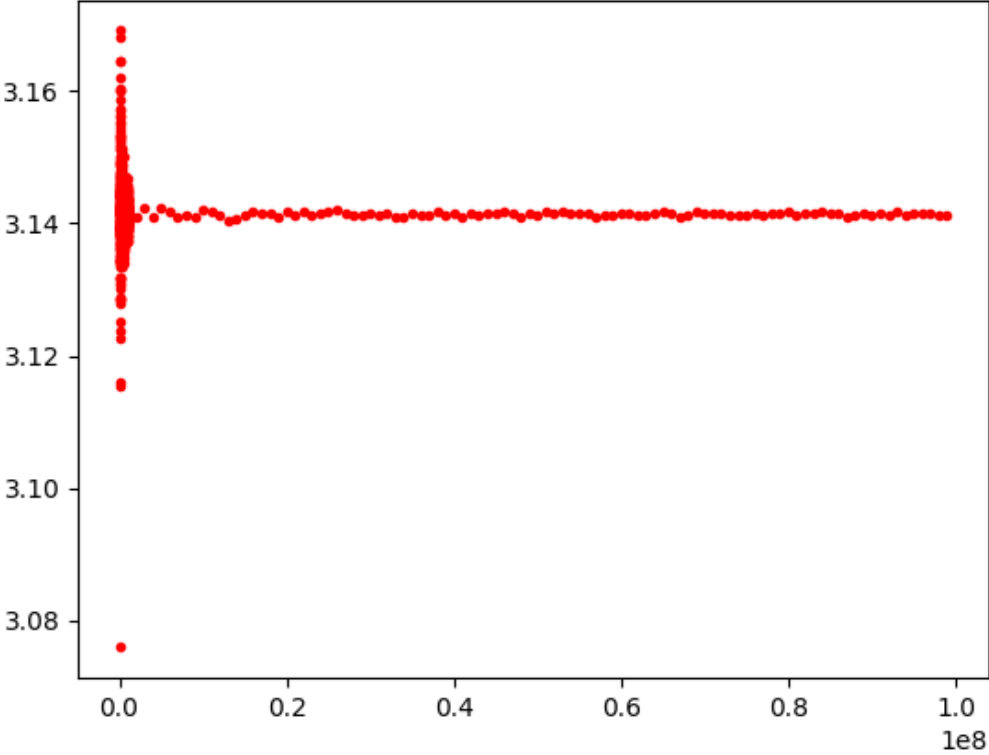
- mulsd – multiplies 2 double values
- mulss – multiplies 2 float values

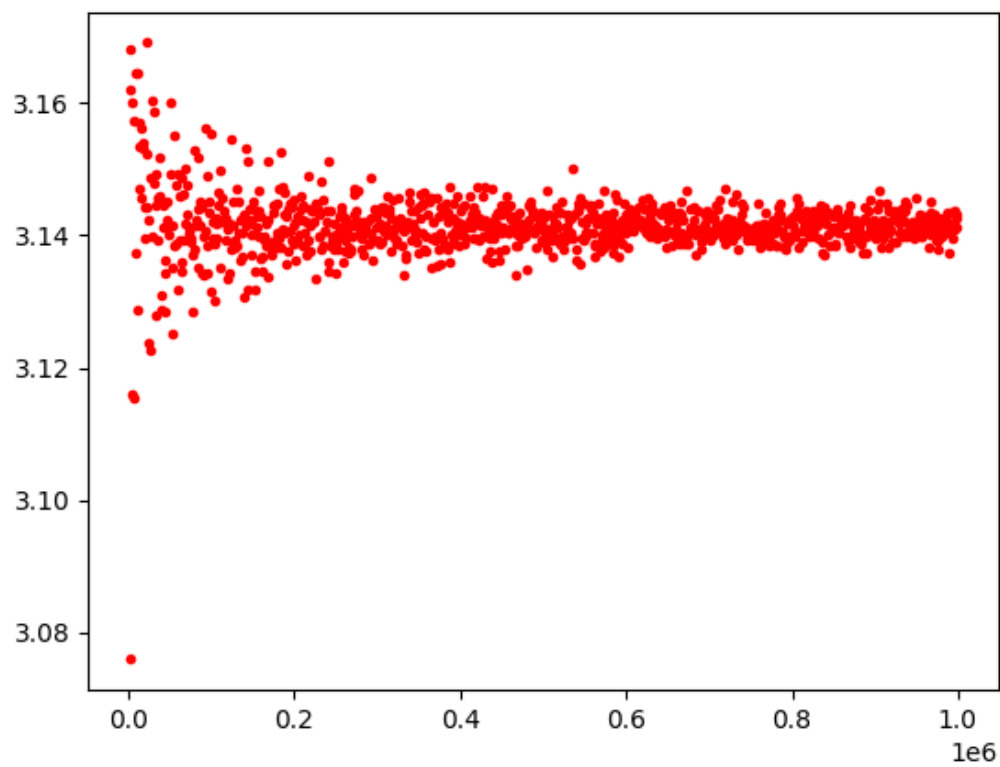
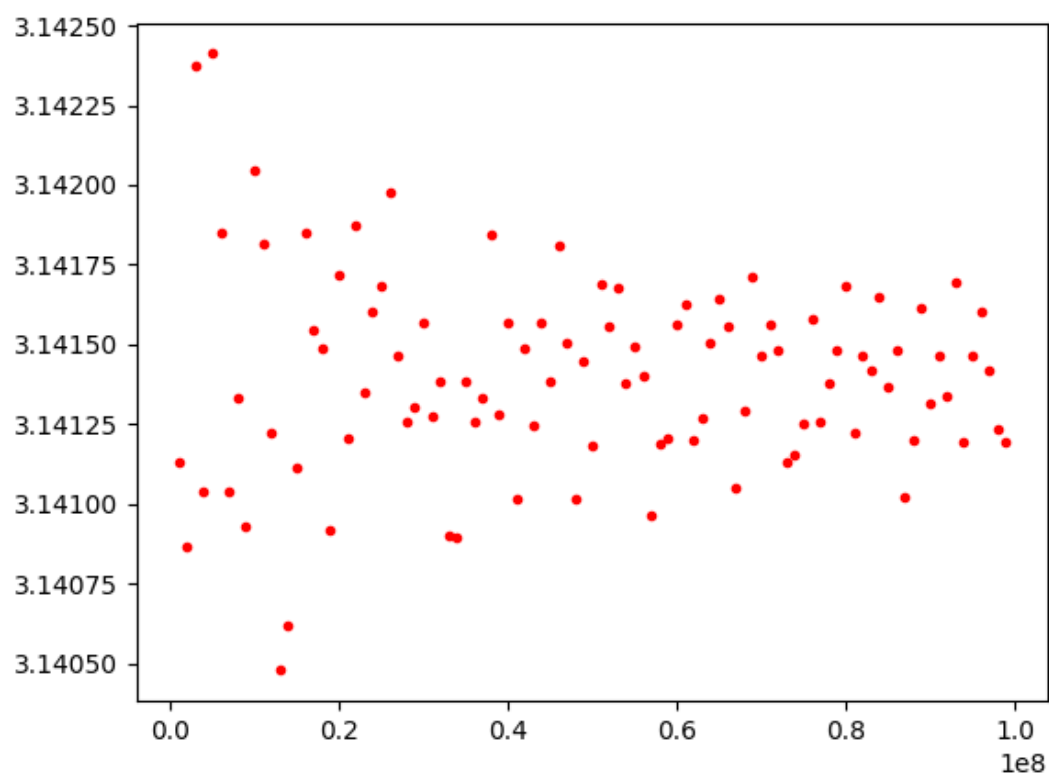
```

main:
.LFB0:
    .cfi_startproc
    endbr64
    pushq    %rbp
    .cfi_def_cfa_offset 16
    .cfi_offset 6, -16
    movq     %rsp, %rbp
    .cfi_def_cfa_register 6
    movsd    b(%rip), %xmm1
    movsd    .LC0(%rip), %xmm0
    mulsd    %xmm1, %xmm0
    movsd    %xmm0, -8(%rbp)
    movss    a(%rip), %xmm1
    movss    .LC1(%rip), %xmm0
    mulss    %xmm1, %xmm0
    movss    %xmm0, -12(%rbp)
    movl     $0, %eax
    popq     %rbp

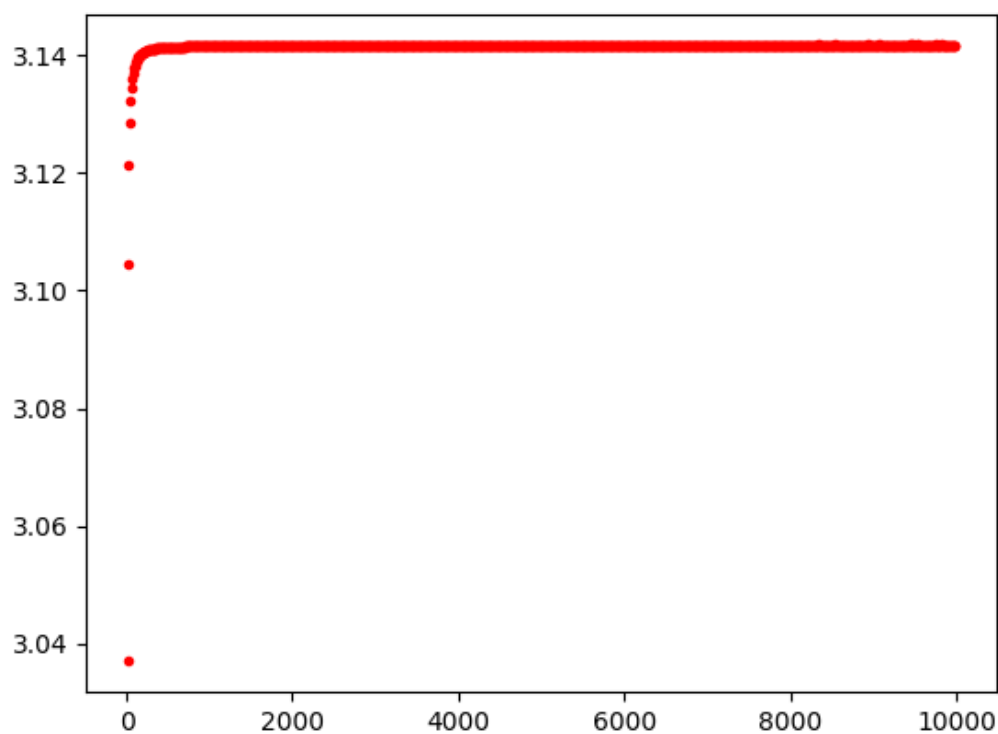
```

6. $\pi(N)$

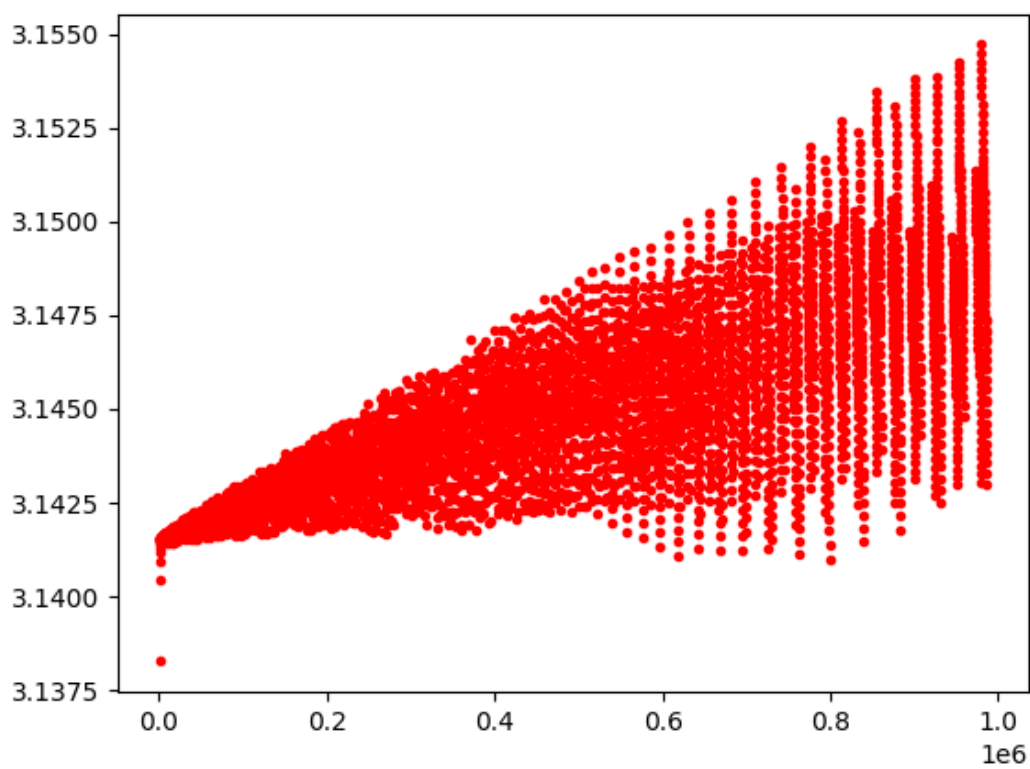
Method	Plot
Wallis product	 <p>The plot shows the convergence of the Wallis product. The x-axis represents the number of terms, ranging from 50 to 9850 in increments of 250. The y-axis represents the value of the product, ranging from 3,1 to 3,15 in increments of 0,01. A blue line starts at approximately 3,11 for 50 terms and rises sharply, reaching a plateau of about 3,14159 by 1000 terms, which it maintains for the remainder of the iterations.</p>
Monte-Carlo method	 <p>The plot shows the convergence of the Monte-Carlo method. The x-axis represents the number of iterations, ranging from 0.0 to 1.0 with a multiplier of $1e8$ at the end, in increments of 0.2. The y-axis represents the value of the product, ranging from 3,08 to 3,16 in increments of 0,02. Red dots represent individual iterations. There is a high initial variance, with points scattered between 3,075 and 3,165 for the first 0,05 $1e8$ iterations. After this initial phase, the points converge to a stable horizontal line at approximately 3,14159.</p>



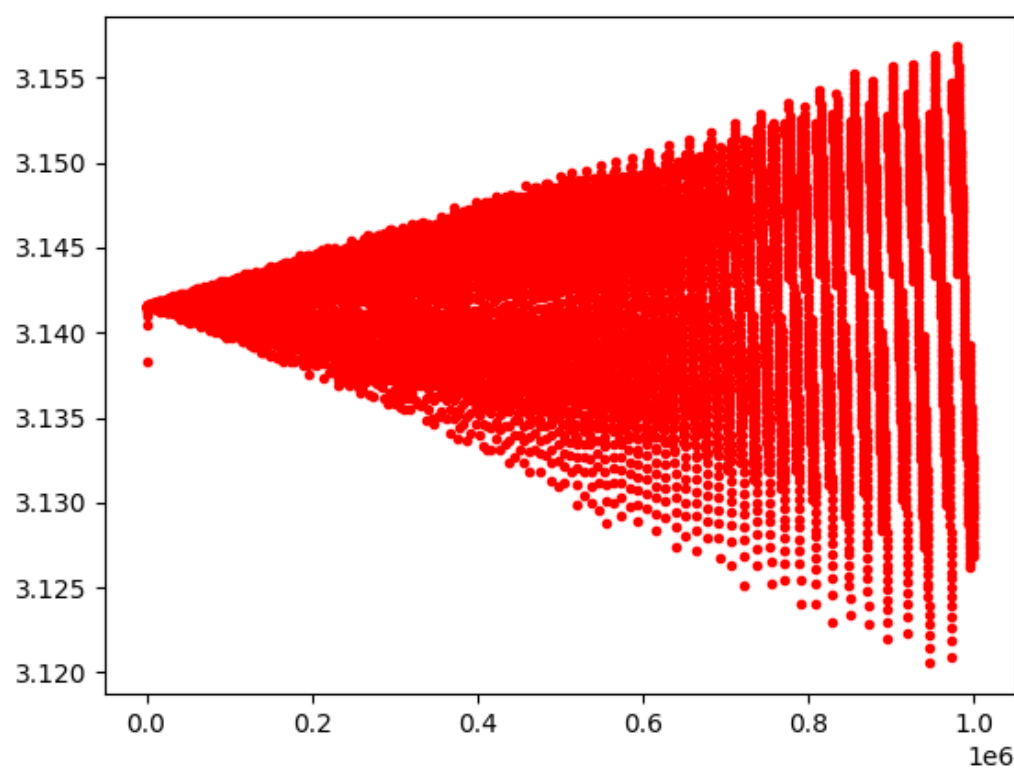
$$\int_{-1}^1 \sqrt{1-x^2}$$



$dx > 2 * 10^{-4}$



$dx > 2 * 10^{-4}$



$dx < 2 * 10^{-4}$ (more strict integration boundaries)

7. Denormalized numbers

```
#include <iostream>
#include <xmmintrin.h>

int main()
{
    //_mm_setcsr(_mm_getcsr() | 0x8040);
    float a = 1;
    int i = 0;
    while(a > 0)
    {
        a /= 2;
        i += 1;
        std::cout << a << "\n";
    }
    printf("%d", i);
}
```

	float	double
Normalized	1.17549e-38 0 127sonamkr1m@s	2.22507e-308 0 1023sonamkr1m
Denormalized	1.4013e-45 0 150sonamkr1	4.94066e-324 0 1075sonamkr1

First value is the minimum floating-point number; third value is the abs of maximum power of 2, which represents 0.

8. FTZ and DAZ

FTZ - sets denormal results from floating-point calculations to zero.

DAZ - creates denormal values used as input to floating-point instructions as zero.