

Лабораторная работа 7 (0111 = 7)

Вычисления с плавающей запятой. Скалярные команды AVX/SSE

Цель работы: исследовать особенности арифметики с плавающей запятой. Научиться использовать скалярные команды расширений AVX/SSE.

Задание Л7.№1

Разработайте программу на языке C/C++, выполняющую вычисления над числами с плавающей запятой одинарной точности (*float*). Проверьте, что программа действительно работает с операндами одинарной точности, а не приводит к типу *float* окончательный результат.

Для частичной суммы гармонического ряда $S(N) = \sum_{i=1}^N \frac{1}{i} \in \mathbb{R}$ найдите две её оценки:

$S_d(N)$ — последовательно складывая члены, начиная от $i = 1$ и заканчивая $i = N$ («наивный» порядок),

$S_a(N)$ — от $i = N$ к $i = 1$.

Сравните $S_d(N)$ и $S_a(N)$ для различных значений N : 10^3 , 10^6 , 10^9 . Объясните результат.

Измените тип операндов на *double*. Объясните результат.

Все N печатаются в экспоненциальной форме, а не с хвостом из нулей.

Все $S_d(N)$ и $S_a(N)$ печатаются *print32()* и *print64()* из Л2

```
Задание №1
=====
SystemInfo
=====
OS: Linux
Архитектура процессора: x86_64 (64-бит)
Compiler: GCC
Version: 13.2.1
=====
N = 1.000000e+03
float sum forward: 40EF890A 100000011101111100010010001010 1089440010 +1089440010 +0x1.df1214p+2 +7.485478e+00 +7.49
float sum backward: 40EF88FC 1000000111011111000100011111100 1089439996 +1089439996 +0x1.df11f8p+2 +7.485472e+00 +7.49
double sum forward: 401DF11F45F4E618 1000101111101001110011000011000 4620113909371954712 +4620113909371954712 +0x1.df11f45f4e618p+2 +7.485471e+00 +7.49
double sum backward: 401DF11F45F4E615 1000101111101001110011000010101 4620113909371954709 +4620113909371954709 +0x1.df11f45f4e615p+2 +7.485471e+00 +7.49
N = 1.000000e+06
float sum forward: 416587BD 1000001011001011011101111011101 1097185213 +1097185213 +0x1.cb6f7ap+3 +1.435736e+01 +14.36
float sum backward: 4166484D 1000001011001100100010001001101 1097222221 +1097222221 +0x1.cc909ap+3 +1.439265e+01 +14.39
double sum forward: 402CC9137A1DF006 1111010000111011111000011010110 4624292002893000918 +4624292002893000918 +0x1.cc9137a1df00d6p+3 +1.439273e+01 +14.39
double sum backward: 402CC9137A1DF28F 1111010000111011111001010001111 4624292002893001359 +4624292002893001359 +0x1.cc9137a1df28fp+3 +1.439273e+01 +14.39
N = 1.000000e+09
float sum forward: 4176757C 1000001011101100111010101111100 1098282364 +1098282364 +0x1.eceaf8p+3 +1.540368e+01 +15.40
float sum backward: 4196769E 1000001100101100111011010011110 1100379806 +1100379806 +0x1.2ced3cp+4 +1.880792e+01 +18.81
double sum forward: 40354CEC5B11A3D5 1011011000100011010001111010101 4626688770216928213 +4626688770216928213 +0x1.54cec5b11a3d5p+4 +2.130048e+01 +21.30
double sum backward: 40354CEC5B11A131 1011011000100011010000100110001 4626688770216927537 +4626688770216927537 +0x1.54cec5b11a131p+4 +2.130048e+01 +21.30
=====
```

Рис. 1: Результаты вычисления гармонического ряда

Объяснение результатов

- Для *float* при больших N наивный порядок даёт меньшую точность, так как маленькие слагаемые теряются при сложении с большими.
- Обратный порядок уменьшает ошибку, т.к. сначала суммируются маленькие слагаемые.
- Для *double* точность намного выше, расхождение между двумя порядками меньше.
- При очень больших N (например, 10^9) вычисления на *float* сильно теряют точность.

Листинг:

Файл task7_1.c:

```
18. void run_task7_1()
19. {
20.     printf("\nЗадание №1\n");
21.     printf("=====");
22.     printSystemInfo();
23.     printf("=====\n");
24.
25.     int Ns[] = {1000, 1000000, 1000000000};
26.     int count = sizeof(Ns)/sizeof(Ns[0]);
27.
28.     for (int idx = 0; idx < count; idx++) {
29.         int N = Ns[idx];
30.         printf("N = %e\n", (double)N);
31.
32.         float Sf = harmonic_sum_float_forward(N);
33.         float Sb = harmonic_sum_float_backward(N);
34.
35.         double Sd = harmonic_sum_double_forward(N);
36.         double Sb_d = harmonic_sum_double_backward(N);
37.
38.         printf("float sum forward: ");
39.         print32(&Sf);
40.         printf("float sum backward: ");
41.         print32(&Sb);
42.
43.         printf("double sum forward: ");
44.         print64(&Sd);
45.         printf("double sum backward: ");
46.         print64(&Sb_d);
47.     }
48.
49.     printf("=====\n");
50. }
51.
52. // Сумма float от 1 до N (наивный порядок)
53. float harmonic_sum_float_forward(int N) {
54.     float sum = 0.0f;
55.     for (int i = 1; i <= N; i++) {
56.         float term = 1.0f / (float)i;
57.         sum += term;
58.     }
59.     return sum;
60. }
61.
62. // Сумма float от N до 1 (обратный порядок)
63. float harmonic_sum_float_backward(int N) {
64.     float sum = 0.0f;
65.     for (int i = N; i >= 1; i--) {
66.         float term = 1.0f / (float)i;
67.         sum += term;
68.     }
69.     return sum;
70. }
71.
72. // Сумма double от 1 до N (наивный порядок)
73. double harmonic_sum_double_forward(int N) {
74.     double sum = 0.0;
75.     for (int i = 1; i <= N; i++) {
76.         double term = 1.0 / (double)i;
77.         sum += term;
78.     }
79.     return sum;
80. }
81.
82. // Сумма double от N до 1 (обратный порядок)
83. double harmonic_sum_double_backward(int N) {
84.     double sum = 0.0;
85.     for (int i = N; i >= 1; i--) {
86.         double term = 1.0 / (double)i;
87.         sum += term;
88.     }
89.     return sum;
90. }
```

Задание Л7.№2

Рассчитайте на языке C/C++ для заданного 32-битного значения с плавающей запятой x его модуль $|x|$, используя только битовые (целочисленные) операции и преобразование указателей.

Исходное значение и результат печатаются `print32()` из Л2.

```
Задание №2
=====
SystemInfo
=====
ОС: Linux
Архитектура процессора: x86_64 (64-бит)
Compiler: GCC
Version: 13.2.1
=====
x = C2F6E979 1100001011110110110100101111001 3270961529 -1024005767 -0x1.edd2f2p+6 -1.234560e+02 -123.46
|x| = 42F6E979 100001011110110110100101111001 1123477881 +1123477881 +0x1.edd2f2p+6 +1.234560e+02 +123.46
=====
```

Рис. 2: результат выполнения `float_abs`

Листинг:

Файл `task7_2.c`:

```
1. float float_abs(float x);
2.
3. void run_task7_2()
4. {
5.     printf("\nЗадание №2\n");
6.     printf("=====");
7.     printSystemInfo();
8.     printf("=====\\n");
9.
10.    float x = -123.456f;
11.    printf("x = ");
12.    print32(&x);
13.
14.    float abs_x = float_abs(x);
15.    printf("|x| = ");
16.    print32(&abs_x);
17.
18.    printf("=====\\n");
19. }
20.
21. float float_abs(float x) {
22.     uint32_t *p = (uint32_t *)&x;
23.     uint32_t bits = *p;
24.
25.     // Сброс знакового бита (старший бит)
26.     bits &= 0x7FFFFFFF;
27.
28.     float *res = (float *)&bits;
29.     return *res;
30. }
31.
```

Задание Л7.№3

Разработайте на ассемблере функцию `inc32_asm(void *p)`, которая принимает адрес переменной `p` и выполняет целочисленный инкремент (команда `inc`) 32-битного (суффикс `l`) значения по этому адресу.

Разработайте на C/C++ программу, которая применяет `inc32_asm(void *p)` к 32-битным переменным с плавающей запятой типа `float`. Значения переменных `a`, `b`, `c`, `d` соответствуют варианту (таблица Л2.2 из Л2). Каждое исходное значение и каждый результат печатаются `print32()` из Л2.

Как можно реализовать аналогичную функцию `inc32_c(void *p)` на C/C++?

```
Задание №6
=====
SystemInfo
=====
ОС: Linux
Архитектура процессора: x86_64 (64-бит)
Compiler: GCC
Version: 13.2.1
=====
Исходные значения:
x = 00000005 00000101 5 +5 +0x1.4p-147 +7.006492e-45 +0.00
y = FFFFFFFB 111111111111111111111111111111011 4294967291 -5 -nan -nan -nan
a = 3F800000 1111111000000000000000000000000000 1065353216 +1065353216 +0x1p+0 +1.000000e+00 +1.00
b = 40000000 1000000000000000000000000000000000 1073741824 +1073741824 +0x1p+1 +2.000000e+00 +2.00
c = 483C6159 1001011001111000110000101011001 1262248281 +1262248281 +0x1.78c2b2p+23 +1.234569e+07 +12345689.00
d = 4CEB79AF 1001100111010110111100110101111 1290500527 +1290500527 +0x1.d6f35ep+26 +1.234569e+08 +123456888.00

После inc32_asm:
a = 3F800001 1111111000000000000000000000000000 1065353217 +1065353217 +0x1.000002p+0 +1.000000e+00 +1.00
b = 40000001 1000000000000000000000000000000000 1073741825 +1073741825 +0x1.000002p+1 +2.000000e+00 +2.00
c = 483C615A 1001011001111000110000101011001 1262248282 +1262248282 +0x1.78c2b4p+23 +1.234569e+07 +12345690.00
d = 4CEB79B0 1001100111010110111100110101000 1290500528 +1290500528 +0x1.d6f36p+26 +1.234569e+08 +123456896.00

После inc32_c:
x = 00000006 00000110 6 +6 +0x1.8p-147 +8.407791e-45 +0.00
y = FFFFFFFC 111111111111111111111111111111100 4294967292 -4 -nan -nan -nan
=====
```

Рис. 3: результат выполнения `inc32_c` и `inc32_asm`

Листинг:

Файл task7_3.c:

```
1. #define showABCD printf("a = "); print32(&a); \
2.     printf("b = "); print32(&b); \
3.     printf("c = "); print32(&c); \
4.     printf("d = "); print32(&d);
5.
6. #define showXY printf("x = "); \
7.     print32(&x); \
8.     printf("y = "); \
9.     print32(&y);
10.
11. void inc32_asm(void *p);
12. void inc32_c(void *p);
13.
14. void run_task7_3()
15. {
16.     printf("\nЗадание №6\n");
17.     printf("=====");
18.     printSystemInfo();
19.     printf("=====\n");
20.
21.     // Инициализация переменных
22.     float a = 1.0f;
23.     float b = 2.0f;
24.     float c = 12345689.0f;
25.     float d = 123456891.0f;
26.
27.     // Значения x=5, y=-5, a,b,c,d
28.     int x = 5;
29.     int y = -5;
30.
31.     // Печать исходных значений
32.     printf("Исходные значения:\n");
33.     showXY
34.     showABCD
35.
36.     // Применяем inc32_asm к float-переменным
37.     inc32_asm(&a);
38.     inc32_asm(&b);
39.     inc32_asm(&c);
40.     inc32_asm(&d);
41.
42.     printf("\nПосле inc32_asm:\n");
43.     showABCD
44.
45.     // Для сравнения применяем inc32_c к x и y
46.     inc32_c(&x);
47.     inc32_c(&y);
48.
49.     printf("\nПосле inc32_c:\n");
50.     showXY
51.
52.     printf("=====\n");
53. }
54.
55. void inc32_c(void *p) {
56.     int *ip = (int *)p;
57.     (*ip)++;
58. }
59.
60. void inc32_asm(void *p) {
61.     __asm__ volatile (
62.         "incl (%0)"
63.         :
64.         : "r"(p)
65.         : "memory"
66.     );
67. }
```

Задание Л7.№4.

Вычислите для заданных x и y с плавающей запятой двойной точности выражение из таблицы Л7.1, используя скалярные AVX-команды либо их SSE-аналоги

Вариант 2: $z = 1 - 5/x - y^2/7$

```
Задание №4
=====
SystemInfo
=====
ОС: Linux
Архитектура процессора: x86_64 (64-бит)
Compiler: GCC
Version: 13.2.1
SSE поддерживается
AVX поддерживается
=====
z = 1 - 5/2.000 - 3.000^2/7 = -2.785714e+00
=====
```

Рис. 4: выполнение avx_compute

Листинг:

Файл task7_4.c:

```
1. double avx_compute(double x, double y);
2.
3. void run_task7_4()
4. {
5.     printf("\nЗадание №4\n");
6.     printf("=====");
7.     printSystemInfo();
8.     checkAVXorSSE();
9.     printf("=====\n");
10.
11.     double x = 2.0;
12.     double y = 3.0;
13.
14.     double z = avx_compute(x, y);
15.     printf("z = 1 - 5/%.3f - %.3f^2/7 = %e\n", x, y, z);
16.
17.     printf("=====\n");
18. }
19.
20. double avx_compute(double x, double y) {
21.     __m128d vx = _mm_set_sd(x);
22.     __m128d vy = _mm_set_sd(y);
23.
24.     __m128d five = _mm_set_sd(5.0);
25.     __m128d one = _mm_set_sd(1.0);
26.     __m128d seven = _mm_set_sd(7.0);
27.
28.     __m128d div1 = _mm_div_sd(five, vx);
29.     __m128d y_sq = _mm_mul_sd(vy, vy);
30.     __m128d div2 = _mm_div_sd(y_sq, seven);
31.     __m128d res = _mm_sub_sd(one, div1);
32.     res = _mm_sub_sd(res, div2);
33.
34.     double result;
35.     _mm_store_sd(&result, res);
36.     return result;
37. }
```

Задание Л7.№5.

Разработайте программу, вычисляющую по введенным значениям x и y с плавающей запятой двойной точности значение z (таблица Л7.2), вызывая функции `libm pow()/atan2()`.

Если программа не собирается из-за отсутствия ссылок на `pow()/atan2()`, добавьте к команде сборки ключ `-lm` (указание компоновщику использовать `libm`).

Вариант 2: $z = \text{atan2}(x, y)$, угол между вектором (x, y) и осью абсцисс.

```
Задание №5
=====
Введите x y: 3 9
atan2(3.000, 9.000) = 3.217506e-01
=====
```

Рис. 5: выполнение `atan2` из `libm`

```
Задание №5
=====
SystemInfo
=====
ОС: Linux
Архитектура процессора: x86_64 (64-бит)
Compiler: GCC
Version: 13.2.1
=====
Введите x y: 3 9
atan2(3.000000, 9.000000) = 3.217506e-01
=====
```

Рис. 6: выполнение `my_atan2` на ассемблерной вставке

Листинг:

Файл `lr7_5.c`:

```
1. int main()
2. {
3.     printf("\nЗадание №5\n");
4.     printf("=====\n");
5.     double x, y;
6.     printf("Введите x y: ");
7.     if (scanf("%lf %lf", &x, &y) != 2) {
8.         printf("Ошибка ввода\n");
9.         return 1;
10.    }
11.
12.    double z = atan2(x, y);
13.    printf("atan2(%.3f, %.3f) = %e\n", x, y, z);
14.    printf("=====\n");
15. }
```

Файл `task7_6.c`:

```
1. double my_atan2(double y, double x);
2. void run_task7_5()
3. {
4.     double x, y;
5.     printf("Введите x y: ");
6.     if (scanf("%lf %lf", &x, &y) != 2) {
7.         printf("Ошибка ввода\n");
8.         return;
9.     }
10.
11.    double z = my_atan2(x, y);
12.    printf("atan2(%.3f, %.3f) = %e\n", x, y, z);
13. }
14.
15. double my_atan2(double y, double x) {
16.     double result;
17.     __asm__ (
18.         "fpatan"           // Инструкция процессора для atan2
19.         : "=t" (result)    // Результат в ST(0)
20.         : "0" (x), "u" (y) // Входные значения в ST(0) и ST(1)
21.         : "st(1)"          // Регистр FPU, который будет изменён
22.     );
23.     return result;
24. }
```

Задание Л7.№6.

Разработайте функцию `double mce_sd(void * p, size_t N)`, которая аналогично Л6.№6 обрабатывает массив из `double` (таблица вариантов Л6.5).

```
Задание №6
=====
SystemInfo
=====
ОС: Linux
Архитектура процессора: x86_64 (64-бит)
Compiler: GCC
Version: 13.2.1
=====
Array:
3FF8000000000000 0 4609434218613702656 +4609434218613702656 +0x1.8p+0 +1.500000e+00 +1.50
4000000000000000 0 4611686018427387904 +4611686018427387904 +0x1p+1 +2.000000e+00 +2.00
4008000000000000 0 4613937818241073152 +4613937818241073152 +0x1.8p+1 +3.000000e+00 +3.00
Product of array elements =
4022000000000000 0 4621256167635550208 +4621256167635550208 +0x1.2p+3 +9.000000e+00 +9.00
=====
```

Рис. 7: выполнение `mce_sd`

Листинг:

Файл `task7_6.c`:

```
1. void run_task7_6()
2. {
3.     printf("\nЗадание №6\n");
4.     printf("=====");
5.     printSystemInfo();
6.     printf("=====\n");
7.
8.     double arr[] = {1.5, 2.0, 3.0};
9.     size_t N = sizeof(arr) / sizeof(arr[0]);
10.
11.    double prod = mce_sd(arr, N);
12.    printf("Array:\n");
13.    PRINT_ARRAY(arr, print64);
14.    printf("Product of array elements = \n");
15.    print64(&prod);
16.
17.    printf("=====\n");
18. }
19.
20. double mce_sd(void *p, size_t N) {
21.     double *arr = (double *)p;
22.     double result = 1.0;
23.     for (size_t i = 0; i < N; i++) {
24.         result *= arr[i];
25.     }
26.     return result;
27. }
```