

Московский Физико-Технический Институт

---

Лаба по флоатам

---

Выполнил:  
Студент 1 курса ФАКТ  
Группа Б03-504  
Подмосковнов Лев

## 0. unsigned int -> binary

```
1 void ui_to_bin(unsigned int a) {
2     for (int i = 32 - 1; i >= 0; i--) {
3         cout << ((a >> i) & 1);
4         if (i != 0 && i % 8 == 0) {
5             cout << " ";
6         }
7     }
8 }
```

Функция для перевода и печати в двоичную систему переменную типа unsigned int.

## 1. float -> binary

```
1 union fu {
2     float f;
3     unsigned int u;
4 };
5
6 void ui_to_bin(unsigned int a) {
7     for (int i = 32 - 1; i >= 0; i--) {
8         cout << ((a >> i) & 1);
9         if (i == 31 || i == 23) {
10             cout << "|";
11         }
12         if (i != 0 && i % 4 == 0) {
13             cout << " ";
14         }
15     }
16 }
```

В структуру fu записываем в f наше число. Теперь нужно посмотреть на биты float через unsigned int.

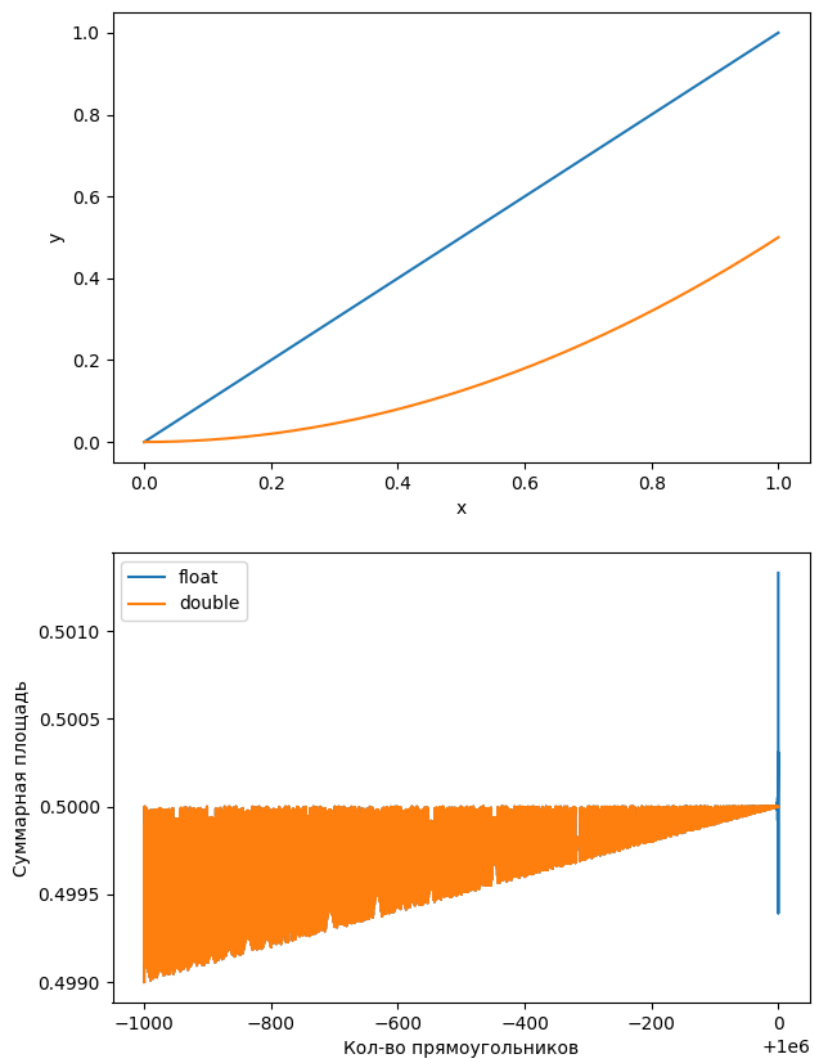
## 2. Переполнение мантисы

В ходе эксперимента найдено, что последняя степень 10, которая не переполняет мантису 11. Когда хотим положить в float  $2^{12}$  получается число 99999997952.00.

## 3. Бесконечный цикл

В ходе эксперимента было замечено, что к число 16777216.0 в float нельзя добавить единицу.

## 4. Численное интегрирование



При увеличении количества столбцов результат становится точнее и в float, и в double. Float ломается быстрее чем double.

## 5\*. Антипереполнение

Subnormal / denormal — это числа с очень маленьким модулем, которые находятся между нулём и самым маленьким нормальным числом IEEE-754. Они имеют нулевой (biased) экспонент-поле и ненулевую мантиссу; обеспечивают плавное (gradual) антипереполнение. Без них некоторые разности могли бы сразу дать 0.

Underflow (антипереполнение) — ситуация, когда результат арифметики меньше минимального нормального значения. С появлением subnormals результат может быть subnormal (потеря точности), а без них — прямо 0.

DAZ (Denormals-Are-Zero) — если входной операнд (операнды) является денормалом, он заменяется на ноль до выполнения операции (то есть входные денормалы “приравняются к 0”).

FTZ (Flush-To-Zero) — если результат операции получился денормалом, вместо него возвращается 0.

При DAZ/FTZ OFF вы часто увидите, что операция с денормалами значительно медленнее