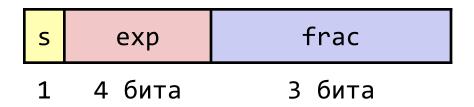
Лекция 16

4 апреля

## Пример



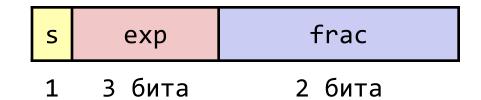
- 8-разрядные числа с плавающей точкой
  - Знаковый бит старший бит
  - Следующие четыре бита порядок, смещение 7
  - Последние три бита дробная часть (мантисса)
- Выполнены все требования стандарта IEEE 754 к формату числа
  - Реализованы нормализованные и денормализованные числа
  - Представлены значения 0, NaN, бесконечность

# Диапазоны значений (только для положительных чисел)

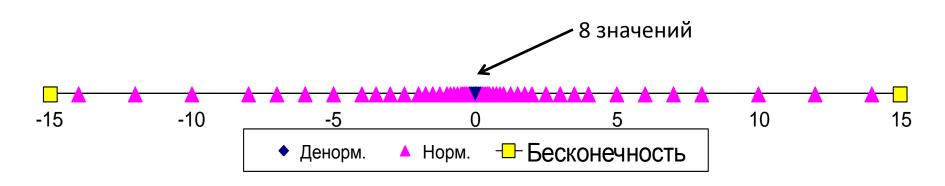
|                 | s   | exp  | frac | E   | Значения                               |  |
|-----------------|-----|------|------|-----|--|--|
|                 | 0   | 0000 | 000  | -6  | 0                                      |  |
|                 | 0   | 0000 | 001  | -6  | 1/8*1/64 = 1/512 Ближайшее к 0         |  |
| _               | 0   | 0000 | 010  | -6  | 2/8*1/64 = 2/512                       |  |
| Денормализов    | анн | ые   |      |     |  |  |
| числа           | 0   | 0000 | 110  | -6  | 6/8*1/64 = 6/512                       |  |
|                 | 0   | 0000 | 111  | -6  | 7/8*1/64 = 7/512 Наибольшее денорм.    |  |
|                 | 0   | 0001 | 000  | -6  | 8/8*1/64 = 8/512<br>Наименьшее норм.   |  |
|                 | 0   | 0001 | 001  | -6  | 9/8*1/64 = 9/512                       |  |
|                 |     |      |      |     |  |  |
|                 | 0   | 0110 | 110  | -1  | 14/8*1/2 = 14/16                       |  |
|                 | 0   | 0110 | 111  | -1  | 15/8*1/2 = 15/16 Ближайшее к 1 «снизу» |  |
|                 | 0   | 0111 | 000  | 0   | 8/8*1 = 1                              |  |
|                 | 0   | 0111 | 001  | 0   | 9/8*1 = 9/8 Ближайшее к 1 «сверху»     |  |
|                 | _   | 0111 | 010  | 0   | 10/8*1 = 10/8                          |  |
| Нормализованные |     |      |      |     |  |  |
| числа           | 0   | 1110 | 110  | 7   | 14/8*128 = 224                         |  |
|                 | 0   | 1110 | 111  | 7   | 15/8*128 = 240 Наибольшее норм.        |  |
|                 | 0   | 1111 | 000  | n/a | inf                                    |  |
|                 |     |      |      |     | 2                                      |  |

## Распределение значений по числовой прямой

- 6-разрядный формат
  - е = 3 бита порядка
  - f = 2 бита мантиссы
  - Смещение 2<sup>3-1</sup>-1 = 3

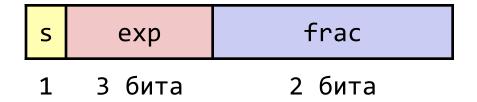


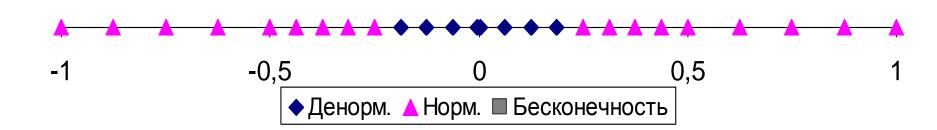
• Распределение сильно «сгущается» в окрестности 0



# Распределение значений по числовой прямой (вид вблизи)

- 6-разрядный формат
  - е = 3 бита порядка
  - f = 2 бита мантиссы
  - Смещение 3





## Некоторые числа

| Описание   | exp  | frac | Численное значение                             |  |  |  |  |
|--|--|------|--|--|--|--|--|
| • Ноль   | 0000   | 0000 | 0.0  |  |  |  |  |
| • Наименьшее «+» денорм.                                     | 0000   | 0001 | $2^{-\{23,52\}} \times 2^{-\{126,1022\}}$      |  |  |  |  |
| <ul> <li>Одинарная точность ≈ 1.4 x 1</li> </ul>             | 0 <sup>-45</sup>   |      |  |  |  |  |  |
| <ul> <li>Двойная точность ≈ 4.9 х 10<sup>-3</sup></li> </ul> | 24   |      |  |  |  |  |  |
| • Наибольшее денорм.   | 0000   | 1111 | $(1.0 - \varepsilon) \times 2^{-\{126,1022\}}$ |  |  |  |  |
| <ul> <li>Одинарная точность ≈ 1.18 x</li> </ul>              | <ul> <li>Одинарная точность ≈ 1.18 x 10<sup>-38</sup></li> </ul> |      |  |  |  |  |  |
| <ul> <li>Двойная точность ≈ 2.2 x 10<sup>-3</sup></li> </ul> | 808  |      |  |  |  |  |  |
| • Наименьшее «+» норм.                                       | 0001   | 0000 | $1.0 \times 2^{-\{126,1022\}}$                 |  |  |  |  |
| – Немногим больше чем наибольшее денормализованное           |  |      |  |  |  |  |  |
| • Единица  | 0111   | 0000 | 1.0  |  |  |  |  |
| • Наибольшее норм.   | 1110   | 1111 | $(2.0 - \varepsilon) \times 2^{\{127,1023\}}$  |  |  |  |  |
| <ul> <li>Одинарная точность ≈ 3.4 x 1</li> </ul>             | $0^{38}$   |      |  |  |  |  |  |
| <ul> <li>Двойная точность ≈ 1.8 х 10<sup>30</sup></li> </ul> | 8  |      |  |  |  |  |  |

Точность {одинарная,двойная}

## Особенности кодировки

- FP ноль совпадает с целочисленным нулем
  - Все биты = 0
- Допустимо (в большинстве случаев) использовать беззнаковое целочисленное сравнение
  - Сперва сравниваем знаковые биты
  - Необходимо рассматривать –0 = 0
  - NaNs
    - В целочисленной интерпретации больше, чем любые другие числа
    - Что необходимо выдавать в качестве результата сравнения?
  - В противном случае ...
    - Денормализованные vs. Нормализованные
    - Нормализованные vs. Бесконечность

## Операции над числами с плавающей точкой

• 
$$x +_f y = Round(x + y)$$

• 
$$x \times_f y = Round(x \times y)$$

- Основная идея
  - Сперва вычислить точный результат
  - Поместить результат в требуемый размер точности
    - Переполнение, если порядок слишком большой
    - Возможно придется округлять поле frac

## Округление

## • Способы округления

| •                               | 1.40 | 1.60 | 1.50 | 2.50 | -1.50 |
|---------------------------------|------|------|------|------|-------|
| – К нулю                        | 1    | 1    | 1    | 2    | -1    |
| $-$ К наименьшему ( $-\infty$ ) | 1    | 1    | 1    | 2    | -2    |
| $-$ К наибольшему (+ $\infty$ ) | 2    | 2    | 2    | 3    | -1    |
| – К ближайшему (✔)              | 1    | 2    | 2    | 2    | -2    |

## Округление к ближайшему целому числу

- Основной способ округления
  - Все остальные способы дают статистическое смещение
    - Пример: суммирование положительных чисел будет давать устойчивую недо- или пере- оценку результата
- Применимо при округлении в произвольной позиции дроби
  - Когда число расположено точно посредине двух значений к которым можно округлить
    - Округляют к тому числу, у которого наименьшая значащая цифра четная
  - Например, округление до ближайших сотых

```
1.2349999
1.23
1.2350001
1.24
1.2350000
1.24 (середина — округляем к большему)
1.2450000
1.24 (середина — округляем к меньшему)
```

## Округление двоичных чисел

#### • Двоичные дробные числа

- "Четные" числа у которых младший значащий бит 0
- "Середина" когда биты справа от позиции к которой происходит округление = 100...<sub>2</sub>

#### • Примеры

– Округление до ближайшей 1/4 (2 бита справа от бинарной точки)

| Число  | Двоичное                 | Окр.   | Действие    | Окр. число |
|--------|--------------------------|--------|-------------|------------|
| 2 3/32 | 10.000112                | 10.002 | (<1/2—down) | 2          |
| 2 3/16 | 10.00 <mark>110</mark> 2 | 10.012 | (>1/2—up)   | 2 1/4      |
| 2 7/8  | 10.11 <mark>100</mark> 2 | 11.002 | ( 1/2—up)   | 3          |
| 2 5/8  | $10.10100_2$             | 10.102 | ( 1/2—down) | 2 1/2      |

#### Умножение

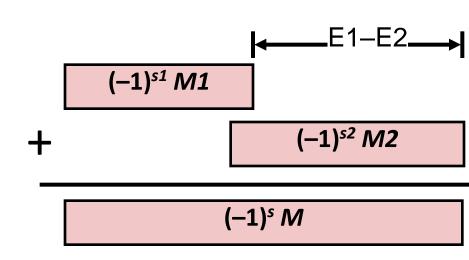
- $(-1)^{s1}$  M1  $2^{E1}$  x  $(-1)^{s2}$  M2  $2^{E2}$
- Точный результат: (-1)<sup>s</sup> M 2<sup>E</sup>
  - Знаковый бит s: s1 ^ s2
  - Мантисса М: М1 х М2
  - Порядок Е: E1 + E2

#### • Исправление

- Если  $M \ge 2$ , сдвигаем M вправо (делим на 2), увеличивая E
- Если Е выходит за пределы, переполнение
- Округляем *М* до соответствующего размера поля frac

#### Сложение

- $(-1)^{s1} M1 2^{E1} + (-1)^{s2} M2 2^{E2}$ - $\Pi$ yctb E1 > E2
- Точный результат: (-1)<sup>s</sup> M 2<sup>E</sup>
  - -Знаковый бит s, мантисса M:
    - Результат выравнивания и сложения
  - -Порядок E: E1



#### • Исправление

- –Если М ≥ 2, сдвигаем М вправо, увеличивая Е
- –Если M < 1, сдвигаем M влево на k позиций, уменьшая E на k
- -Переполнение если Е выходит за пределы
- –Округляем M до соответствующего размера поля frac

#### Математические свойства сложения

- Выполняются ли свойства Абелевых групп
  - Замкнутость?
    - Результатом может быть бесконечность или NaN
  - Коммутативность?
  - Ассоциативность?
    - Переполнения и изменение результата при округлении
  - 0.0 нейтральный элемент?
  - Каждый элемент имеет обратный
    - За исключением бесконечности и NaN
- Монотонность
  - $-a \ge b \Rightarrow a+c \ge b+c$ ?
    - За исключением бесконечности и NaN

## Математические свойства умножения

- Выполняются ли свойства коммутативных колец
  - Замкнуто ли относительно умножения?
    - Результат может быть бесконечность или NaN
  - Умножение коммутативно?
  - Умножение ассоциативно?
    - Возможность переполнения, неточности округления
  - 1.0 мультипликативная единица?
  - Умножение дистрибутивно над сложением?
    - Возможность переполнения, неточности округления
- Монотонность
  - $-a \ge b \& c \ge 0 \Rightarrow a * c \ge b *c?$ 
    - Исключение бесконечность и NaN

## Числа с плавающей точкой в языке Си

- Язык Си вводит два уровня точности
  - -float одинарная точность
  - -double двойная точность
- Приведение типа
  - —Приведение типа между int, float, и double включает изменение битового представления
  - double/float  $\rightarrow$  int
    - Отбрасывается дробная часть (аналогично округлению к нулю)
    - Поведение не определено, когда значение вне допустимого диапазона или NaN: как правило устанавливается TMin
  - int  $\rightarrow$  double
    - Точное приведение, поскольку long и int 32 бита ≤ 53 бита
  - int  $\rightarrow$  float
    - Будет округляться согласно принятым соглашениям

## Задачи

- Для каждого Си-выражения объяснить:
  - почему оно верно для любого значения переменных, ...
  - ... либо почему ложно

```
int x = ...;
float f = ...;
double d = ...;
```

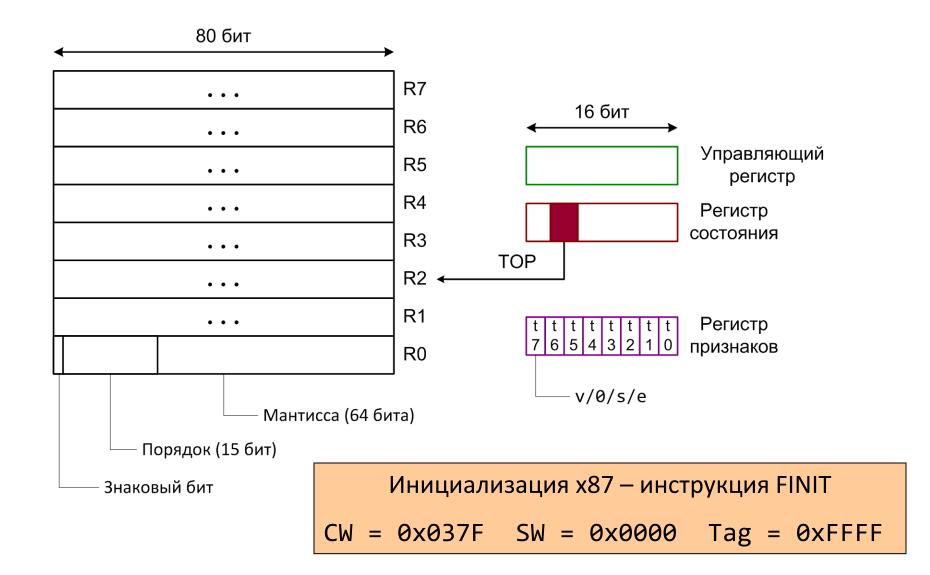
Предполагается, что d и f не являются NaN

```
• x == (int)(double) x
• x == (int)(float) x
• f == (float)(double) f
• d == (float) d
• f == -(-f);
• 2/3 == 2/3.0
• d < 0.0 \Rightarrow ((d*2) < 0.0)
• d > f \Rightarrow -f > -d
• d * d >= 0.0
• (d+f)-d == f
```

## Промежуточные итоги

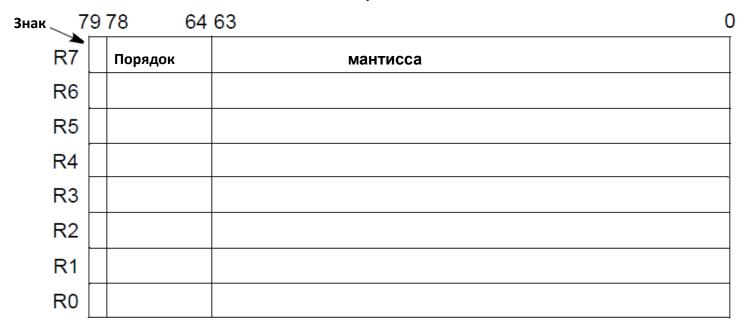
- IEEE754 четкое определение математических свойств
- Представляются числа вида М х 2<sup>E</sup>
- Семантика операций не зависит от особенностей аппаратуры
  - Сперва точное вычисление, затем округление
- Отличия от «настоящей» арифметики
  - Нарушаются свойства ассоциативности и дистрибутивности
  - Создаются сложности для компилятора и серьезных математических вычислений

## Упрощенная схема х87



## Размер чисел с плавающей точкой

#### Регистры данных



#### Обмен данными только с памятью.

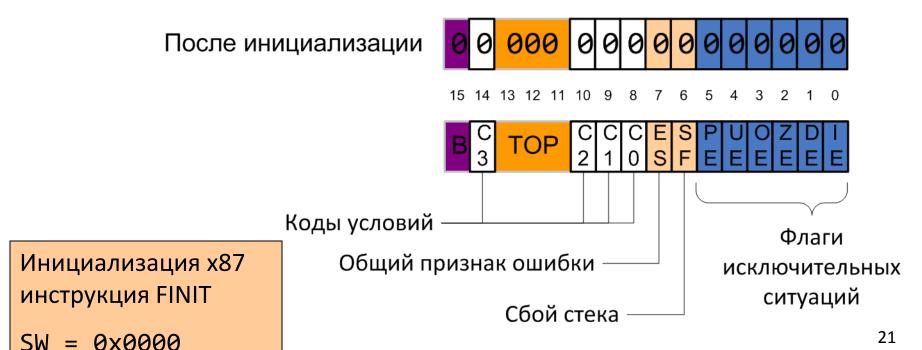
dd 1.234567e20 ; Константы одинарной точности

dq 1.234567e20 ; Двойной точности

dt 1.234567e20 ; Расширенной точности

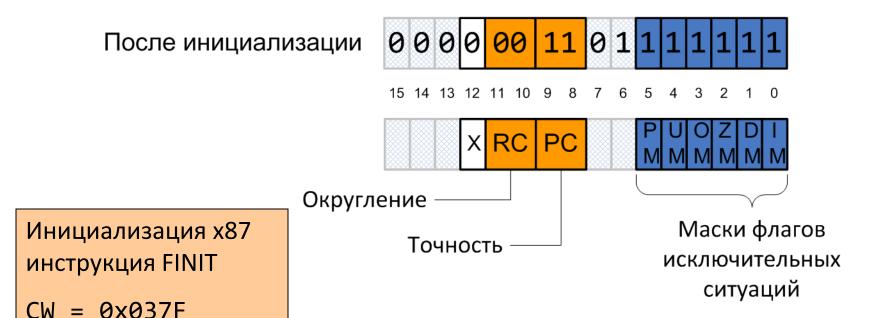
## Слово (регистр) состояния

- SF переполнение стека (С1 показывает направление)
- Исключительные ситуации: точность, переполнение, деление на ноль, денормализованный операнд, «неправильные» данные



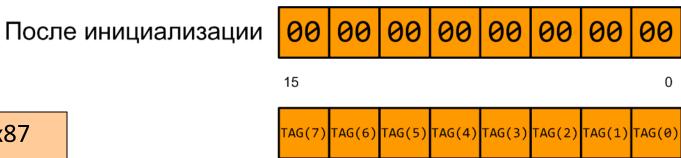
## Управляющий регистр

- Точность: одинарная, двойная, расширенная
- Округление: к ближайшему четному, к нулю, к +/- бесконечности
- Флаг X совместимость с 287
- Маски соответствуют исключениям в слове состояния



## Регистр признаков (тагов)

- Состояние регистров
  - 0 нормализованное число с плавающей точкой
  - 1 число ноль
  - -2 особые числа (NaN,  $\pm \infty$ , денормализованное число)
  - 3 регистр свободен
- Нумерация соответствует физическим регистрам



Инициализация x87 инструкция FINIT

Tag = 0xFFFF

## NASM и числа с плавающей точкой

```
db - 0.2
                            ; «Четверть»
dw - 0.5
                            ; IEEE 754r/SSE5
                             половинная точность
dd 1.2
                             одинарная точность
dd 1.222_222_222
                             допускается использовать
                            ; знак подчеркивания
dd 0x1p+2
                            1.0x2^2 = 4.0
                            1.0x2^32 = 4294967296.0
dq 0x1p+32
dq 1.e10
                            ; 10 000 000 000.0
dq 1.e+10
                            ; синоним для 1.e10
dq 1.e-10
                            ; 0.000 000 000 1
dt 3.141592653589793238462 ; число Пи
do 1.e+4000
                            ; IEEE 754r четверная точность
```

IEEE 754r – опубликован в 2008 году

## NASM и числа с плавающей точкой

- \_\_float8\_\_
- \_\_float16\_\_\_
- \_\_float32\_\_
- \_\_float64\_\_\_
- \_\_float80m\_\_
- \_\_float80e\_\_\_
- \_\_float1281\_\_\_
- \_\_float128h\_\_\_

- \_\_Infinity\_\_\_
- \_\_NaN\_\_\_

```
dq +1.5, -__Infinity__, __NaN__
mov eax, __float32__(3.1415926)
```

## Сложение двух чисел

```
%include 'io.inc'
                                CMATN:
                                   finit
                                   fld dword [x]
section .data
 x dd 11.2
                                   fld dword [y]
 y dd 0.7
                                   faddp
                                   fstp dword [z]
section .bss
                                   PRINT HEX 4, z
 z resd 1
                                   NFWI TNF
                                   xor eax, eax
section .text
                                   ret
global CMAIN
```

```
x

11.2 ~ 1.119999980926513671875E1

0x41333333 = 01000001 00110011 00110011
```

## Сложение двух чисел

```
%include 'io.inc'
                                CMATN:
                                   finit
                                   fld dword [x]
section .data
 x dd 11.2
                                   fld dword [y]
 y dd 0.7
                                   faddp
                                   fstp dword [z]
section .bss
                                   PRINT_HEX 4, z
 z resd 1
                                   NFWI TNF
                                   xor eax, eax
section .text
                                   ret
global CMAIN
```

```
y

0.7 ~ 6.99999988079071044921875E-1

0x3F333333 = 00111111 00110011 00110011
```

## Сложение двух чисел

```
%include 'io.inc'
                                CMATN:
                                   finit
                                   fld dword [x]
section .data
 x dd 11.2
                                   fld dword [y]
 y dd 0.7
                                   faddp
                                   fstp dword [z]
section .bss
                                   PRINT HEX 4, z
 z resd 1
                                   NFWI TNF
                                   xor eax, eax
section .text
                                   ret
global CMAIN
```

```
z

0x413E6666 = 01000001 00111110 01100110 01100110

1.18999996185302734375E1 ~ 11.9
```

#### Печать числа

```
%include 'io.inc'
                              CMATN:
section .data
                                 ; пролог функции
x dd 11.2
                                 sub esp, 20
 y dd 0.7
                                 fld
                                        dword [x]
section .bss
                                 fld
                                         dword [y]
                                 faddp
 z resd 1
                                fst dword [z]
                                 fstp
                                        dword [esp + 4]
section .rodata
                                        dword [esp], lc
 lc db '%f', 10, 0
                                mov
                                         printf
                                 call
section .text
CEXTERN printf
                                 add
                                        esp, 20
global CMAIN
                                 ; эпилог функции
```

#### Печать числа

```
printf печатает мусор
                                        где ошибка?!?!!1
%include 'io.inc'
                               CMAIN:
section .data
                                  ; пролог функции
x dd 11.2
                                  sub
                                          esp, 20
 y dd 0.7
                                  fld
                                          dword [x]
section .bss
                                  fld
                                          dword [y]
                                  faddp
 z resd 1
                                  fst
                                          dword [z]
                                  fstp
                                          dword [esp + 4]
section .rodata
                                          dword [esp], lc
 lc db '%f', 10, 0
                                  mov
                                  call
                                          printf
section .text
CEXTERN printf
                                  add
                                          esp, 20
global CMAIN
                                  ; эпилог функции
                                                          30
```

#### Печать числа

```
ISO/IEC 9899:1999
                                        § 6.5.2.2 абзац №6
%include 'io.inc'
                               CMAIN:
section .data
                                  ; пролог функции
 x dd 11.2
                                  sub
                                          esp, 20
 y dd 0.7
                                  fld
                                          dword [x]
section .bss
                                  fld
                                          dword [y]
                                  faddp
 z resd 1
                                  fst
                                          dword [z]
                                  fstp
section .rodata
                                          qword [esp + 4]
                                          dword [esp], lc
 lc db '%f', 10, 0
                                  mov
                                  call
                                          printf
section .text
CEXTERN printf
                                  add
                                          esp, 20
global CMAIN
                                  ; эпилог функции
                                                          31
```

## Порядок действий имеет значение

```
%include 'io.inc'
                              CMATN:
                                 ; … пролог функции
section .data
                                 sub esp, 20
x dq 3.14
                                 fld
y dq 1e50
                                         qword [x]
 z dq -1e50
                                 fld
                                         qword [y]
                                 f1d
                                         qword [z]
                                 faddp
section .bss
                                 faddp
 r resq 1
                                 fst
                                         qword [r]
                                 fstp
                                         qword [esp + 4]
section .rodata
 lc db '%f', 10, 0
                                         dword [esp], lc
                                 mov
                                 call
                                         printf
section .text
CEXTERN printf
                                 add esp, 20
global CMAIN
                                 ; … эпилог функции
```

## Порядок действий имеет значение

```
CMATN:
                              CMATN:
   ; … пролог функции
                                 ; … пролог функции
   sub
                                 sub
      esp, 20
                                         esp, 20
   fld
                                 fld
           qword [y]
                                         qword [x]
           qword [z]
   fld
                                 fld
                                         qword [y]
           qword [x]
   fld
                                 fld
                                         qword [z]
   faddp
                                 faddp
   faddp
                                 faddp
   fst
                                 fst
           qword [r]
                                         qword [r]
   fstp
          qword [esp + 4]
                                 fstp
                                         qword [esp + 4]
           dword [esp], lc
                                         dword [esp], lc
  mov
                                 mov
           printf
                                         printf
   call
                                 call
   add
        esp, 20
                                 add
                                         esp, 20
   ; … эпилог функции
                                 ; … эпилог функции
```

## Распределение слагаемых

```
section .data
                             CMATN:
x dq 1e200
                                ; … пролог функции
y dq 1e200
                                sub esp, 20
z dq 1e200
                                fld
                                        qword [x]
section .bss
                                fld
                                        qword [y]
                                fsubp
 r resq 1
                                fld
                                        qword [z]
                                fmulp
section .rodata
                                fst
lc db '%lf', 10, 0
                                        qword [r]
                                fstp
                                        qword [esp + 4]
                                        dword [esp], lc
                                mov
                                        printf
                                call
                                add esp, 20
                                 ; … эпилог функции
```

## Распределение слагаемых

```
section .data
                             CMATN:
                                 ; ... пролог функции
x dq 1e200
y dq 1e200
                                sub esp, 20
z dq 1e200
                                fld
                                        qword [x]
section .bss
                                fld
                                        qword [z]
                                fmulp
 r resq 1
                                fld
                                        qword [y]
                                fld
                                        qword [z]
section .rodata
                                fmulp
lc db '%lf', 10, 0
                                fsubp
                                 ; вызов printf
                                add
                                        esp, 20
                                 ; … эпилог функции
```

## Распределение слагаемых

```
section .data
                         fstcw word [cw]
                         and word [cw], 11111111_11000000b
x dq 1e200
y dq 1e200
                         fldcw word [cw]
z dq 1e200
                         fld
                                qword [x]
section .bss
                         fld
                                 qword [z]
                         fmulp
r resq 1
                         fld
cw resw 1
                                 qword [y]
                         fld
                                 qword [z]
section .rodata
                         fmulp
lc db '%lf', 10, 0
                         fsubp
                         ; вызов printf
CMATN:
                         add esp, 20
   ; … пролог функции
  sub esp, 20
                         ; … эпилог функции
```

## Польская обратная запись

```
section .text
• (w + x + y + z) / 4
                                   •
•
•
• w x + y + z + 4 /
                                   fld
                                           qword [w]
                                   fld
                                           qword [x]
                                   faddp
                                   fld
                                           qword [y]
                                   faddp
                                   fld
                                           qword [z]
                                   faddp
                                   fild
                                           dword [d]
                                   fdivp
section .data
w dq 1e10
                                   •
•
•
x dq
      1e10
y dq 1e10
 z dq 1e10
```

d dd

4

## Предопределенные константы

• На «верхушку» стека регистров (ST0) помещается определенная константа

```
\begin{array}{lll} - & \text{FLD1} & +1.0 \\ - & \text{FLDL2T} & \log_2 10 \\ - & \text{L2E} & \log_2 e \\ - & \text{FLDPI} & \pi \\ - & \text{FLDLG2} & \log_{10} 2 \end{array}
```

 $log_e 2$ 

+0.0

- FLDLN2

– FLDZ

## Сравнение чисел

```
_Bool isLe(double x, float y) {
   return x <= y;
}</pre>
```

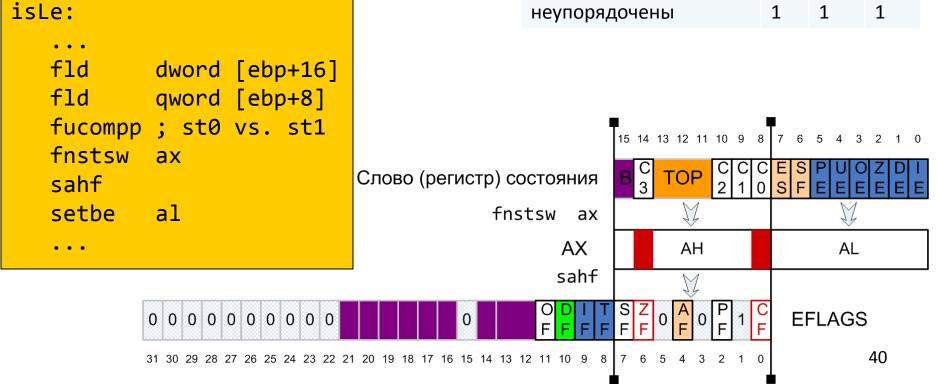
| Результат     | С3 | C2 | C0 |
|---------------|----|----|----|
| сравнения     | 0  | 0  | 0  |
| St0 > St1     | 0  | 0  | 0  |
| St0 < St1     | 0  | 0  | 1  |
| St0 == St1    | 1  | 0  | 0  |
| неопределенно | 1  | 1  | 1  |

```
isLe:
       push ebp
      mov ebp, esp
      fld dword [ebp+16]
      fld qword [ebp+8]
      fucompp ; st0 vs. st1
      fnstsw ax
      sahf
      setbe al
      pop ebp
      ret
```

## Извлечение результатов сравнения

- C3  $\rightarrow$  ZF, C0  $\rightarrow$  CF
- Можно использовать условные коды, применяемые при сравнении беззнаковых чисел

| Результат сравнения | <b>C3</b> | C2 | CO |
|---------------------|-----------|----|----|
| ST(0) > ST(i)       | 0         | 0  | 0  |
| ST(0) < ST(i)       | 0         | 0  | 1  |
| ST(0) == ST(i)      | 1         | 0  | 0  |
| неупорядочены       | 1         | 1  | 1  |



## Функции: возвращение числа с плавающей точкой

```
void caller(double *p) {
                                       float inverse(double x) {
   *p = inverse(*p);
                                          return 1/x;
}
caller:
                                       inverse:
   push
           ebp
                                          push
                                                  ebp
           ebp, esp
                                                ebp, esp
   mov
                                          mov
           esp, 8
                                          fld1
   sub
           eax, dword [ebp+8]
                                          fld
   mov
                                                   qword [ebp+8]
           dword [eax]
   fld
                                          fdivp
   fstp
           qword [esp]
                                          pop ebp
   call
           inverse
                                          ret
           eax, dword [ebp+8]
   mov
   fstp
           qword [eax]
   leave
```

ret