# Floating Point Precision

Excursus

# Representation of Numbers

- Unsigned integers are traditionally represented as a string of zeroes and ones in the 2-adic system

  - E.g. 0100 0111 =
    $$1 \times 2^6 + 1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 = 64 + 4 + 2 + 1 = 71$$

- Integers need to incorporate a sign ±

  - Explicit sign ("signed magnitude representation") leads to inefficient hardware addition / subtraction

  - Use two's complement or one's complement

# Representation of Numbers

- In hardware, integers take up 16, 32, 64, or even 128 bits

- Limits range to

  - 32768 (16 bits)

  - 2147483648 (32 bits)

  - 9223372036854775808 (64 bits)

  - 170141183460469231731687303715884105728 (128b)

- But larger integers are used

# Representation of Numbers

- Overflow:

  - An result of an addition / subtraction / multiplication exceeds the range

  - Depending on platform, can be misinterpreted:

    - E.g.: Adding two large numbers results in a negative number

# Representation of Numbers

- Arbitrary precision integers

  - Overflow happens because results of calculations do not fit into the number of bits assigned for integers

    - This can be a function of the architecture

    - Arbitrary precision integers combine storage for several integers to store a single integer

- Python uses arbitrary precision integers

# Representation of Numbers

- Example:

```
=============================== RESTART: Shell ===============================
>>> 2**2**2**2**2
    Squeezed text (247 lines).

>>>
```

# Representation of Numbers

- Example:

  - Double click on the message

```
================================ RESTART: Shell ================================
>>> 2**2**2**2**2
    2003529930406846464979072351560255750447825475569751419265016973710894059556311
    4895061308809333481010382343429072631818229493821188126688695063647615470291650
    4191635158796634721944293092798208430910485599057015931895963952486337236720300
    2919592156108764948889254090805911457037675208500206671563702366126359747144807
    1117158809141357427209671901518362825606180914588526998261414250301233911082736
    0384387644904320596037912449090570756031403507616256247603186379312648470374378
    2954973770981604614413308692118102485959152380195331030292162800160568670105651
```

- We just calculate $2^{65536}$.

# Representation of Numbers

- Python does this automatically

  - Result: Integer calculation in Python are always exact

- **Nota Bene:**

  - There are Python modules that do not use arbitrary precision integers

# Representation of Floating Point Numbers

- Rational numbers are represented as floating point numbers

  - Stored as sign significant $\times$ base$^{\text{exponent}}$

  - You should know this as the scientific notation for the decimal numbers

    - E.g. Planck's constant $6.62607004 \times 10^{-34} \dfrac{\text{m}^2\text{kg}}{\text{sec}}$

# Representation of Floating Point Numbers

- The significant and the exponent can store limited information

  - This means that some numbers cannot be represented exactly

  - In the decadic system:

    - 1/17 is 0.<u>076923</u>0769230769130769230<u>76923</u>076923

    - with an infinite repetition of the same pattern

    - Or $\pi$ = 3.141592653589793…

# Representation of Floating Point Numbers

- Computers (almost universally) use the binary system

  - But we have the same phenomena

- Consequences:

  - Normal mathematical identities are no longer true

    - E.g. $x(y - z) = xy - xz$

    - E.g. $(\sqrt{x})^2 = x$

```
>>> (3**0.5)**2
2.9999999999999996
```

# Representation of Floating Point Numbers

- Python (Cython):

    - Uses 8 bytes or 64 bits to represent a floating point number

        - Industry standard for high precision floating point numbers

    - There are packages for C++ or Java available for higher precision numbers

    - Python similarly has wrapper modules that make higher precision available

# Representation of Floating Point Numbers

- In theory, it is impossible to use exact precision for floating point calculations