

# Data Structures and Algorithms

Course at D-MATH (CSE) of ETH Zurich

Felix Friedrich

FS 2019

## 1. Introduction

Overview, Algorithms and Data Structures, Correctness, First Example

### Goals of the course

- Understand the design and analysis of fundamental algorithms and data structures.
- An advanced insight into a modern programming model (with C++).
- Knowledge about chances, problems and limits of the parallel and concurrent computing.

### Contents

#### data structures / algorithms

The notion invariant, cost model, Landau notation

algorithms design, induction

searching, selection and sorting

amortized analysis

dynamic programming

Minimum Spanning Trees, Fibonacci Heaps

shortest paths, Max-Flow

Fundamental algorithms on graphs,

dictionaries: hashing and search trees

van-Emde Boas Trees, Splay-Trees

#### programming with C++

RAII, Move Konstruktion, Smart Pointers,

Templates and generic programming

Exceptions

functors and lambdas

promises and futures

threads, mutex and monitors

#### parallel programming

parallelism vs. concurrency, speedup (Amdahl/Gustavson), races, memory reordering, atomic registers, RMW (CAS,TAS), deadlock/starvation

## Algorithm

## 1.2 Algorithms

[Cormen et al, Kap. 1; Ottman/Widmayer, Kap. 1.1]

Algorithm: well defined computing procedure to compute *output* data from *input* data

23

24

### example problem

**Input:** A sequence of  $n$  numbers  $(a_1, a_2, \dots, a_n)$   
**Output:** Permutation  $(a'_1, a'_2, \dots, a'_n)$  of the sequence  $(a_i)_{1 \leq i \leq n}$ , such that  $a'_1 \leq a'_2 \leq \dots \leq a'_n$

#### Possible input

$(1, 7, 3), (15, 13, 12, -0.5), (1) \dots$

Every example represents a *problem instance*

The performance (speed) of an algorithm usually depends on the problem instance. Often there are “good” and “bad” instances.

25

### Examples for algorithmic problems

- Tables and statistics: sorting, selection and searching
- routing: shortest path algorithm, heap data structure
- DNA matching: Dynamic Programming
- evaluation order: Topological Sorting
- autocompletion and spell-checking: Dictionaries / Trees
- Fast Lookup : Hash-Tables
- The travelling Salesman: Dynamic Programming, Minimum Spanning Tree, Simulated Annealing

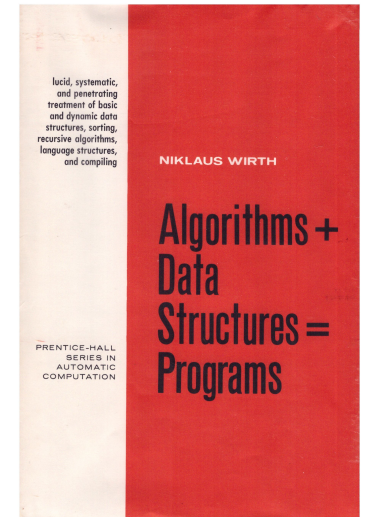
26

## Characteristics

- Extremely large number of potential solutions
- Practical applicability

## Data Structures

- A data structure is a particular way of *organizing data* in a computer so that they can be *used efficiently* (in the algorithms operating on them).
- Programs = algorithms + data structures.



27

28

## Efficiency

Illusion:

- If computers were infinitely fast and had an infinite amount of memory ...
- ... then we would still need the theory of algorithms (only) for statements about correctness (and termination).

Reality: resources are bounded and not free:

- Computing time → Efficiency
- Storage space → Efficiency

**Actually, this course is nearly only about efficiency.**

## Hard problems.

- NP-complete problems: no known efficient solution (the existence of such a solution is very improbable – but it has not yet been proven that there is none!)
- Example: travelling salesman problem

**This course is *mostly* about problems that can be solved efficiently (in polynomial time).**

29

30

## 2. Efficiency of algorithms

Efficiency of Algorithms, Random Access Machine Model, Function Growth, Asymptotics [Cormen et al, Kap. 2.2,3,4.2-4.4 | Ottman/Widmayer, Kap. 1.1]

## Efficiency of Algorithms

### Goals

- Quantify the runtime behavior of an algorithm independent of the machine.
- Compare efficiency of algorithms.
- Understand dependence on the input size.

31

32

## Programs and Algorithms

## Technology Model

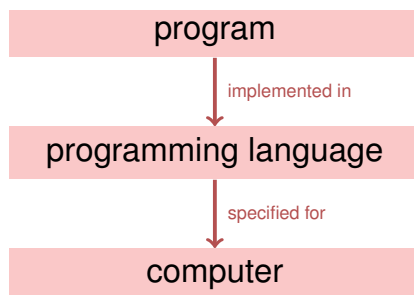
### *Random Access Machine (RAM)*

- Execution model: instructions are executed one after the other (on one processor core).
- Memory model: constant access time (big array)
- Fundamental operations: computations (+, -, ·, ...) comparisons, assignment / copy on machine words (registers), flow control (jumps)
- Unit cost model: fundamental operations provide a cost of 1.
- Data types: fundamental types like size-limited integer or floating point number.

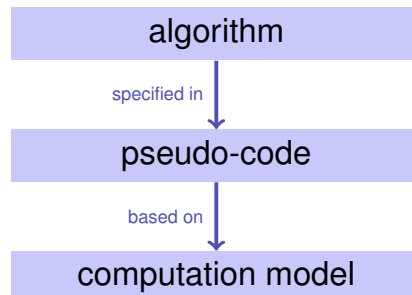
33

34

### Technology



### Abstraction



## Size of the Input Data

Typical: number of input objects (of fundamental type).

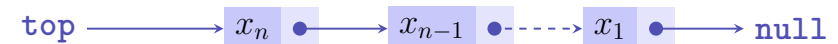
Sometimes: number bits for a *reasonable / cost-effective* representation of the data.

fundamental types fit into word of size :  $w \geq \log(\text{sizeof}(\text{mem}))$  bits.

## Pointer Machine Model

We assume

- Objects bounded in size can be dynamically allocated in constant time
- Fields (with word-size) of the objects can be accessed in constant time 1.



35

36

## Asymptotic behavior

An exact running time of an algorithm can normally not be predicted even for small input data.

- We consider the asymptotic behavior of the algorithm.
- And ignore all constant factors.

### Example

An operation with cost 20 is no worse than one with cost 1  
Linear growth with gradient 5 is as good as linear growth with gradient 1.

## Algorithms, Programs and Execution Time

Program: concrete implementation of an algorithm.

Execution time of the program: measurable value on a concrete machine. Can be bounded from above and below.

### Beispiel

3GHz computer. Maximal number of operations per cycle (e.g. 8).  $\Rightarrow$  lower bound.  
A single operations does never take longer than a day  $\Rightarrow$  upper bound.

From the perspective of the *asymptotic behavior* of the program, the bounds are unimportant.

37

38

## Superficially

## 2.2 Function growth

$\mathcal{O}$ ,  $\Theta$ ,  $\Omega$  [Cormen et al, Kap. 3; Ottman/Widmayer, Kap. 1.1]

Use the asymptotic notation to specify the execution time of algorithms.

We write  $\Theta(n^2)$  and mean that the algorithm behaves for large  $n$  like  $n^2$ : when the problem size is doubled, the execution time multiplies by four.

39

40

## More precise: asymptotic upper bound

provided: a function  $g : \mathbb{N} \rightarrow \mathbb{R}$ .

Definition:<sup>1</sup>

$$\mathcal{O}(g) = \{f : \mathbb{N} \rightarrow \mathbb{R} \mid \exists c > 0, \exists n_0 \in \mathbb{N} : \forall n \geq n_0 : 0 \leq f(n) \leq c \cdot g(n)\}$$

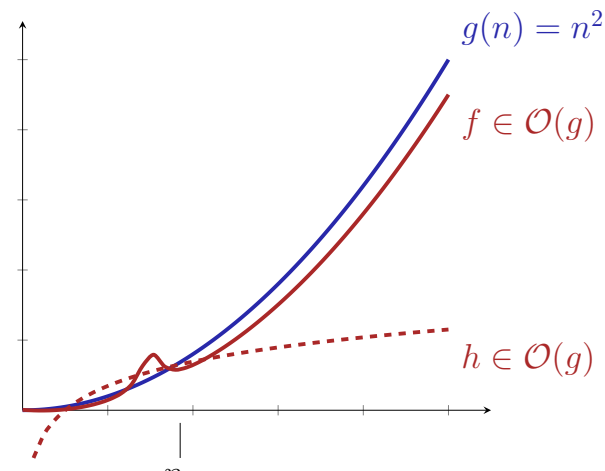
Notation:

$$\mathcal{O}(g(n)) := \mathcal{O}(g(\cdot)) = \mathcal{O}(g).$$

<sup>1</sup>Ausgesprochen: Set of all functions  $f : \mathbb{N} \rightarrow \mathbb{R}$  that satisfy: there is some (real valued)  $c > 0$  and some  $n_0 \in \mathbb{N}$  such that  $0 \leq f(n) \leq c \cdot g(n)$  for all  $n \geq n_0$ .

41

## Graphic



42

## Examples

$$\mathcal{O}(g) = \{f : \mathbb{N} \rightarrow \mathbb{R} \mid \exists c > 0, \exists n_0 \in \mathbb{N} : \forall n \geq n_0 : 0 \leq f(n) \leq c \cdot g(n)\}$$

$f(n)$	$f \in \mathcal{O}(?)$	Example
$3n + 4$	$\mathcal{O}(n)$	$c = 4, n_0 = 4$
$2n$	$\mathcal{O}(n)$	$c = 2, n_0 = 0$
$n^2 + 100n$	$\mathcal{O}(n^2)$	$c = 2, n_0 = 100$
$n + \sqrt{n}$	$\mathcal{O}(n)$	$c = 2, n_0 = 1$

$$f_1 \in \mathcal{O}(g), f_2 \in \mathcal{O}(g) \Rightarrow f_1 + f_2 \in \mathcal{O}(g)$$

43

44

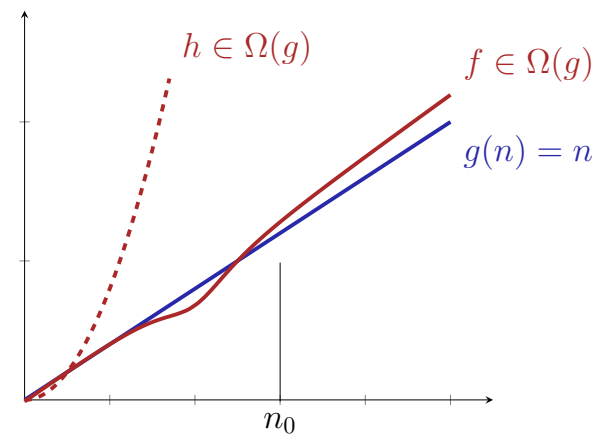
## Converse: asymptotic lower bound

Given: a function  $g : \mathbb{N} \rightarrow \mathbb{R}$ .

Definition:

$$\Omega(g) = \{f : \mathbb{N} \rightarrow \mathbb{R} \mid \exists c > 0, \exists n_0 \in \mathbb{N} : \forall n \geq n_0 : 0 \leq c \cdot g(n) \leq f(n)\}$$

## Example



45

46

## Asymptotic tight bound

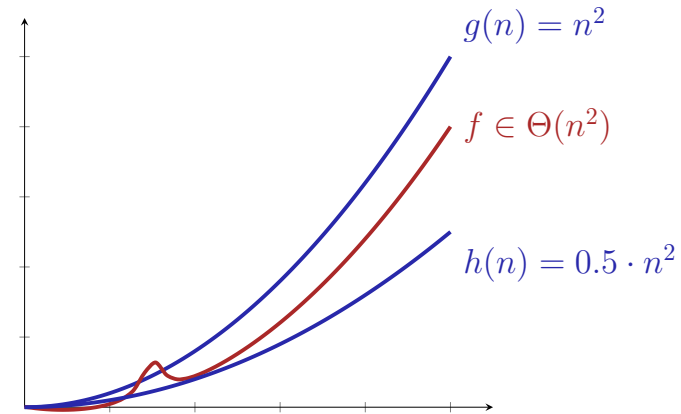
Given: function  $g : \mathbb{N} \rightarrow \mathbb{R}$ .

Definition:

$$\Theta(g) := \Omega(g) \cap \mathcal{O}(g).$$

Simple, closed form: exercise.

## Example



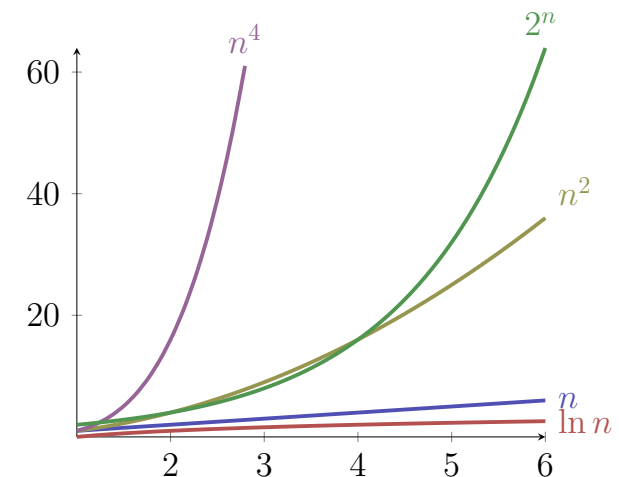
47

48

## Notions of Growth

$\mathcal{O}(1)$	bounded	array access
$\mathcal{O}(\log \log n)$	double logarithmic	interpolated binary sorted sort
$\mathcal{O}(\log n)$	logarithmic	binary sorted search
$\mathcal{O}(\sqrt{n})$	like the square root	naive prime number test
$\mathcal{O}(n)$	linear	unsorted naive search
$\mathcal{O}(n \log n)$	superlinear / loglinear	good sorting algorithms
$\mathcal{O}(n^2)$	quadratic	simple sort algorithms
$\mathcal{O}(n^c)$	polynomial	matrix multiply
$\mathcal{O}(2^n)$	exponential	Travelling Salesman Dynamic Programming
$\mathcal{O}(n!)$	factorial	Travelling Salesman naively

## Small $n$

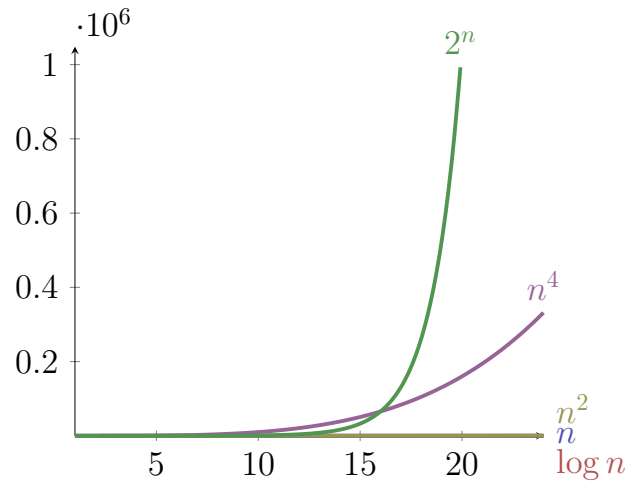


49

50

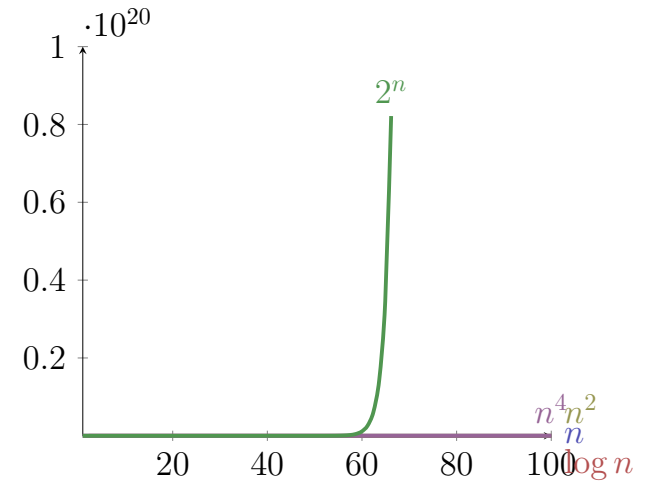


## Larger $n$



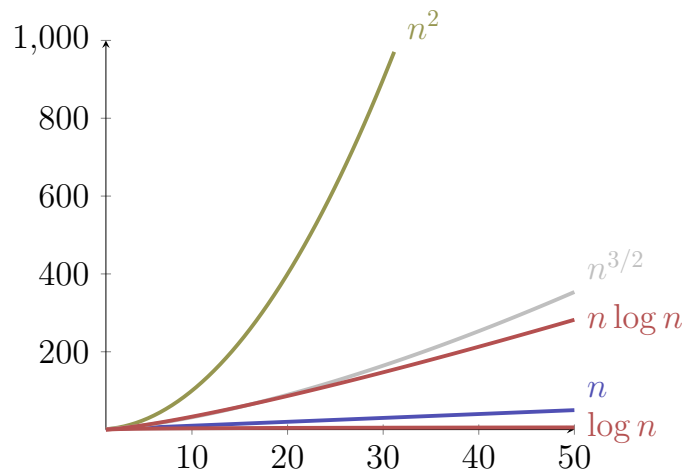
51

## "Large" $n$



52

## Logarithms



53

## Time Consumption

Assumption 1 Operation =  $1\mu s$ .

problem size	1	100	10000	$10^6$	$10^9$
$\log_2 n$	$1\mu s$	$7\mu s$	$13\mu s$	$20\mu s$	$30\mu s$
$n$	$1\mu s$	$100\mu s$	$1/100s$	$1s$	17 minutes
$n \log_2 n$	$1\mu s$	$700\mu s$	$13/100\mu s$	$20s$	8.5 hours
$n^2$	$1\mu s$	$1/100s$	1.7 minutes	11.5 days	317 centuries
$2^n$	$1\mu s$	$10^{14}$ centuries	$\approx \infty$	$\approx \infty$	$\approx \infty$

54

## Useful Tool

### Theorem

Let  $f, g : \mathbb{N} \rightarrow \mathbb{R}^+$  be two functions, then it holds that

- 1  $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0 \Rightarrow f \in \mathcal{O}(g), \mathcal{O}(f) \subsetneq \mathcal{O}(g).$
- 2  $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = C > 0$  ( $C$  constant)  $\Rightarrow f \in \Theta(g).$
- 3  $\frac{f(n)}{g(n)} \xrightarrow{n \rightarrow \infty} \infty \Rightarrow g \in \mathcal{O}(f), \mathcal{O}(g) \subsetneq \mathcal{O}(f).$

## About the Notation

Common casual notation

$$f = \mathcal{O}(g)$$

should be read as  $f \in \mathcal{O}(g).$

Clearly it holds that

$$f_1 = \mathcal{O}(g), f_2 = \mathcal{O}(g) \not\Rightarrow f_1 = f_2!$$

### Beispiel

$n = \mathcal{O}(n^2), n^2 = \mathcal{O}(n^2)$  but naturally  $n \neq n^2.$

*We avoid this notation where it could lead to ambiguities.*

55

56

## Reminder: Efficiency: Arrays vs. Linked Lists

- Memory: our `avec` requires roughly  $n$  ints (vector size  $n$ ), our `llvec` roughly  $3n$  ints (a pointer typically requires 8 byte)
- Runtime (with `avec = std::vector`, `llvec = std::list`):

```

prepending (insert at front) [100,000x]:
  > avec: 675 ms
  > llvec: 10 ms
appending (insert at back) [100,000x]:
  > avec: 2 ms
  > llvec: 9 ms
removing first [100,000x]:
  > avec: 675 ms
  > llvec: 4 ms
removing last [100,000x]:
  > avec: 0 ms
  > llvec: 4 ms

removing randomly [10,000x]:
  > avec: 3 ms
  > llvec: 113 ms
inserting randomly [10,000x]:
  > avec: 16 ms
  > llvec: 117 ms
fully iterate sequentially (5000 elements) [5,000x]:
  > avec: 354 ms
  > llvec: 525 ms
    
```

## Asymptotic Runtimes

With our new language ( $\Omega, \mathcal{O}, \Theta$ ), we can now *state the behavior of the data structures and their algorithms more precisely*

Typical asymptotic running times (Anticipation!)

Data structure	Random Access	Insert	Next	Insert After Element	Search
<code>std::vector</code>	$\Theta(1)$	$\Theta(1) A$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$
<code>std::list</code>	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$
<code>std::set</code>	–	$\Theta(\log n)$	$\Theta(\log n)$	–	$\Theta(\log n)$
<code>std::unordered_set</code>	–	$\Theta(1) P$	–	–	$\Theta(1) P$

$A$  = amortized,  $P$ =expected, otherwise worst case

57

58

# Complexity

*Complexity* of a problem  $P$ : minimal (asymptotic) costs over all algorithms  $A$  that solve  $P$ .

Complexity of the single-digit multiplication of two numbers with  $n$  digits is  $\Omega(n)$  and  $\mathcal{O}(n^{\log_3 2})$  (Karatsuba Ofman).

## Example:

Problem	Complexity	$\mathcal{O}(n)$	$\mathcal{O}(n)$	$\mathcal{O}(n^2)$
		↑	↑	↑
Algorithm	Costs <sup>2</sup>	$3n - 4$	$\mathcal{O}(n)$	$\Theta(n^2)$
		↓	↕	↕
Program	Execution time	$\Theta(n)$	$\mathcal{O}(n)$	$\Theta(n^2)$