



算法设计与分析

第一章 绪论

哈尔滨工业大学
丁小欧

dingxiaou@hit.edu.cn



本讲内容

1.1 什么是算法？

1.2 计算机科学中算法的位置

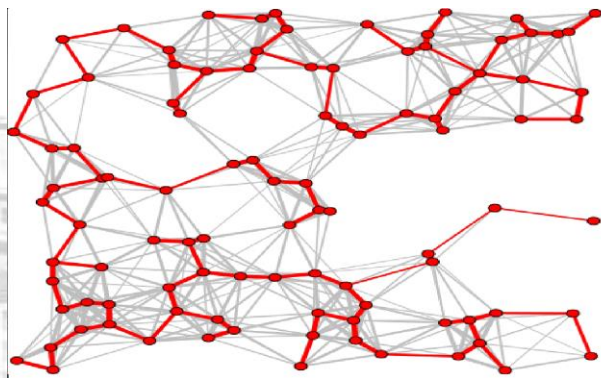
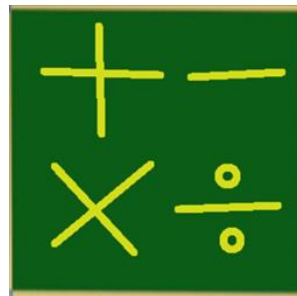
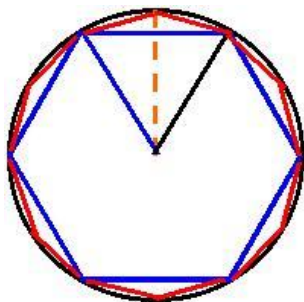
1.3 算法分析引论

1.4 算法设计引论

什么是算法?

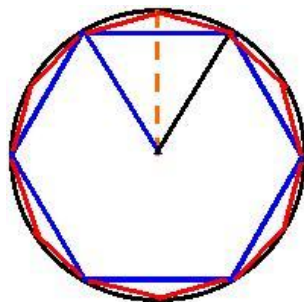
- 在数学和计算机科学之中，算法/算则法 (Algorithm) 为一个计算的具体步骤，常用于计算、数据处理和自动推理。(Wikipedia)
- 算法的例子

- 刘徽割圆术
- 四则运算
- 最小生成树
- 快速排序



什么是算法？

- 在数学和计算机科学之中，算法/算则法（Algorithm）为一个计算的具体步骤，常用于计算、数据处理和自动推理。(Wikipedia)
- 算法的例子



- 刘徽割圆术
- 四则运算
- 最小生成树
- 快速排序

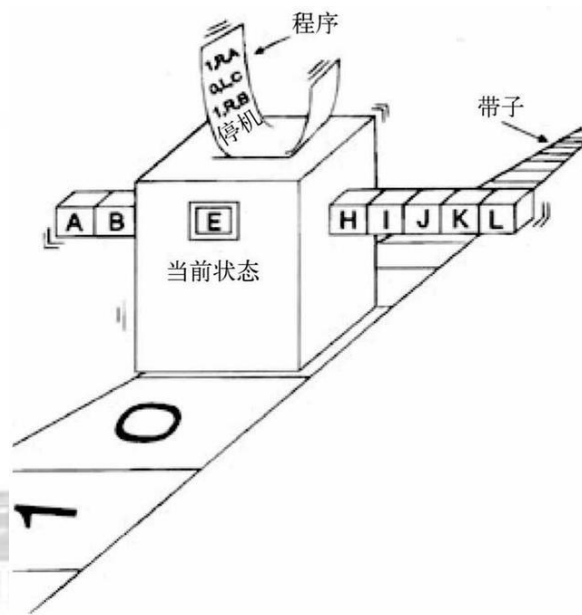
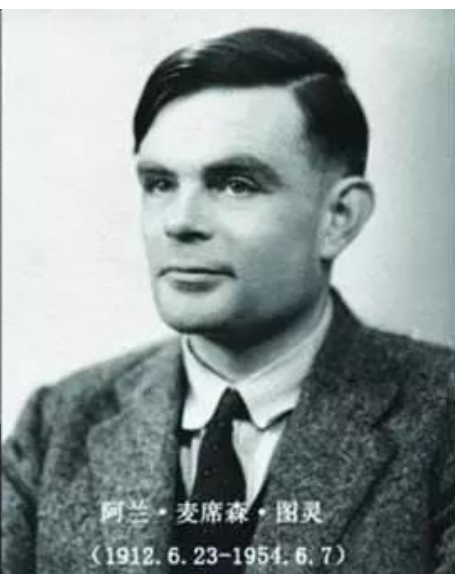
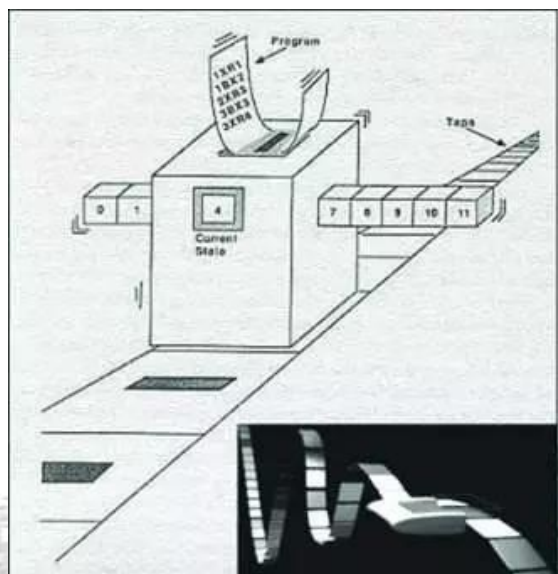
魏晋时期的数学家刘徽首创割圆术，为计算圆周率建立了严密的理论和完善的算法，所谓割圆术，就是不断倍增圆内接正多边形的边数求出圆周率的方法

计算的定义

—可由一个给定计算模型机械地执行的规则或计算步骤序列，称为该计算模型的一个计算

—注意

- 一个计算机程序是一个计算（计算模型是计算机）
- 计算可能永远不停止——不是算法



算法的定义

算法是一个满足下列条件的计算：

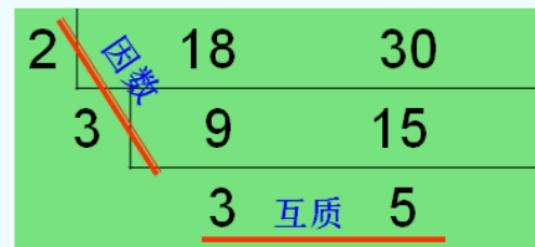
- 有穷性/终止性：有限步内必须停止，
- 确定性：每一步都是严格定义和确定的动作，
- 能行性：每一个动作都能够被精确地机械执行，
- 输入：有一个满足给定约束条件的输入，
- 输出：满足给定约束条件的结果。

关于算法

- “算法”的来源
 - 中文名称：周髀算经
 - 英文名称
 - “Algorithm” 来自于9世纪波斯数学家花拉子米 (al-Khwarizmi)
 - “算法” 原为 “algorism”，即 “al-Khwarizmi” 的音转，意思是 “花拉子米” 的运算法则
 - 在18世纪演变为 “algorithm”
- 最早的算法
 - 欧几里德的 “求最大公因子算法”

关于算法

- “算法”的来源
 - 中文名称：周髀算经
 - 英文名称
 - “Algorithm” 来自于9世纪波斯数学家花拉子米 (al-Khwarizmi)
 - “算法” 原为 “algorism”，即 “al-Khwarizmi” 的音转，意思是 “花拉子米” 的运算法则
 - 在18世纪演变为 “algorithm”
- 最早的算法
 - 欧几里德的 “求最大公因子算法”



2	18	30
3	9	15
	3	5

因數

互質

关于算法

- 算法，不一定只是给计算机才能运行的。只要满足了上述五个条件，由人工自己的智慧也可以执行算法
- 在计算机出现之前，就有了算法，现在有了计算机，就需要更多的算法。
- 算法是计算的核心

问题的定义

- 算法的目的是求解问题。什么是问题？
- 问题
 - 设 Input 和 Output 是两个集合。一个问题是一个关系 $R \subseteq \text{Input} \times \text{Output}$ ，Input 称为问题 R 的输入集合，Input 的每个元素称为 R 的一个输入，Output 称为问题 R 的输出或结果集合，Output 的每个元素称为 R 的一个结果。
 - 注意
 - 问题定义了输入和输出的关系

问题的例子

SORT问题定义如下:

- 输入集合 $\text{Input} = \{ \langle a_1, \dots, a_n \rangle \mid a_i \text{ 是整数} \}$
- 输出集合 $\text{Output} = \{ \langle b_1, \dots, b_n \rangle \mid b_i \text{ 是整数, } b_1 \leq \dots \leq b_n \}$
- 问题 $\text{SORT} = \{ (\langle a_1, \dots, a_n \rangle, \langle b_1, \dots, b_n \rangle) \mid \langle a_1, \dots, a_n \rangle \in \text{Input}, \langle b_1, \dots, b_n \rangle \in \text{Output}, \{a_1, \dots, a_n\} = \{b_1, \dots, b_n\} \}$

• 问题实例

– 问题P的一个实例是P的一个二元组。

– 注意

- 一个算法面向一个问题，而不是仅求解一个问题的某个或某几个实例
- 要有用严谨的算法语言表述问题的思维意识、直击问题的内涵

算法示例

- 问题定义

- $\text{Input} = \{ \langle a_1, \dots, a_n \rangle \mid a_i \text{ 是整数} \}$
- $\text{output} = \{ \langle b_1, \dots, b_n \rangle \mid b_i \text{ 是整数, 且 } b_1 \leq \dots \leq b_n \}$
- $R = \{ (\langle a_1, \dots, a_n \rangle, \langle b_1, \dots, b_n \rangle) \mid \langle a_1, \dots, a_n \rangle \in \text{Input}, \langle b_1, \dots, b_n \rangle \in \text{output}, \{a_1, \dots, a_n\} = \{b_1, \dots, b_n\} \}$

- 算法的思想—扑克牌游戏

算法演示

$A[1, \dots, n] = 5, 2, 4, 6, 1, 3$

$A[1, \dots, n] = 5, 2, 4, 6, 1, 3$

$A[1, \dots, n] = 2, 5, 4, 6, 1, 3$

$A[1, \dots, n] = 2, 4, 5, 6, 1, 3$

$A[1, \dots, n] = 2, 4, 5, 6, 1, 3$

$A[1, \dots, n] = 1, 2, 4, 5, 6, 3$

$A[1, \dots, n] = 1, 2, 3, 4, 5, 6$

算法描述（伪代码：体现算法逻辑和流程）

Insertion-sort (A)

Input: $A[1, \dots, n]$ = n 个数

Output: $A[1, \dots, n]$ = n 个sorted数

FOR $j=2$ To n Do

$key \leftarrow A[j]$;

$i \leftarrow j-1$

 WHILE $i > 0$ AND $A[i] > key$ Do

$A[i+1] \leftarrow A[i]$;

$i \leftarrow i-1$;

$A[i+1] \leftarrow key$;

- 实例: $A[1, \dots, n] = 5, 2, 4, 6, 1, 3$

伪代码中的约定

- 缩进
- 循环
- 赋值
- 非全局变量

阅读《算法导论》第2章“伪代码中的一些约定”进行了解

本讲内容

1.1 什么是算法?

1.2 计算机科学中算法的位置

1.3 算法分析引论

1.4 算法设计引论

算法是计算机科学基础的重要主题

- 70年代前

- 计算机科学基础的主题没有被清楚地认清。

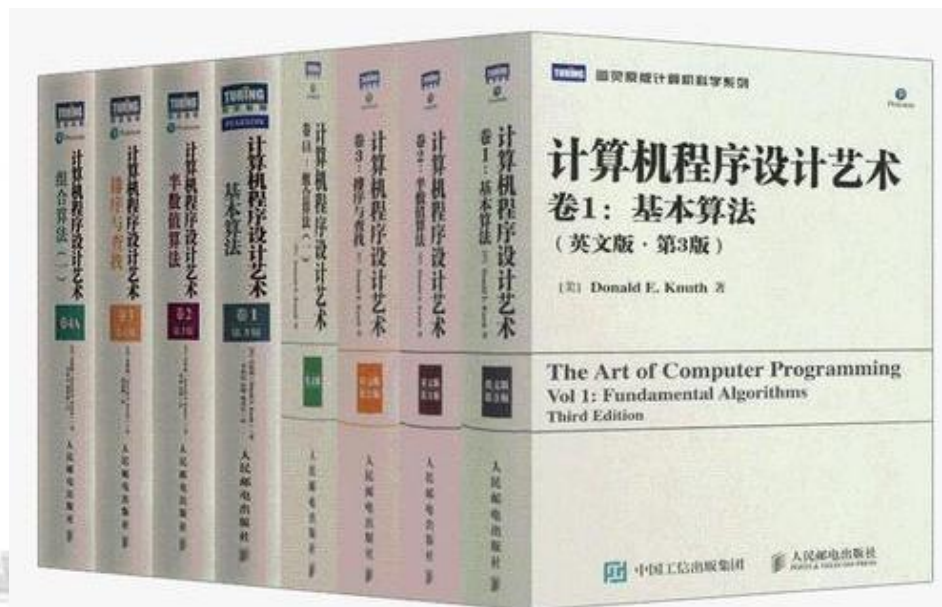
- 70年代

- Knuth出版了《The Art of Computer Programming》

- 以算法研究为主线

- 确立了算法为计算机科学基础的重要主题

- 1974年获得图灵奖

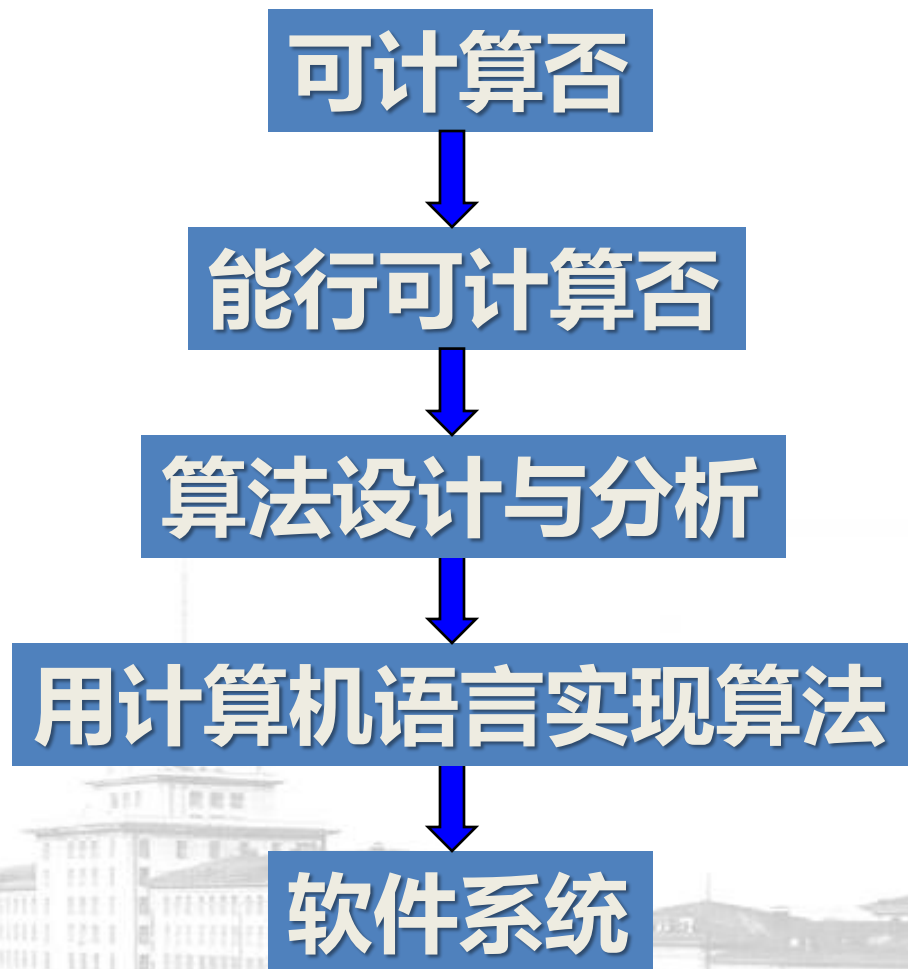


算法是计算机科学基础的重要主题

- 70年代前
 - 计算机科学基础的主题没有被清楚地认清。
- 70年代
 - Knuth出版了《The Art of Computer Programming》
 - 以算法研究为主线
 - 确立了算法为计算机科学基础的重要主题
 - 1974年获得图灵奖
- 70年代后
 - 算法作为计算机科学核心推动了计算机科学技术飞速发展
 - 算法上的贡献是一种重要的学术贡献

计算机科学的体系

- 解决一个计算问题的过程



• 可计算理论

- 计算模型

- 可计算问题/不可计算问题

- 计算模型的等价性--图灵/Church命题

• 计算复杂性理论

- 在给定的计算模型下研究问题的复杂性

 - 固有复杂性

 - 上界

 - 下界

 - 平均

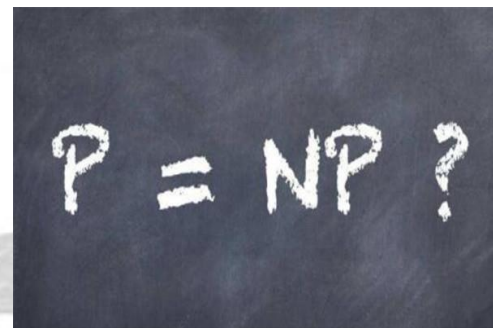
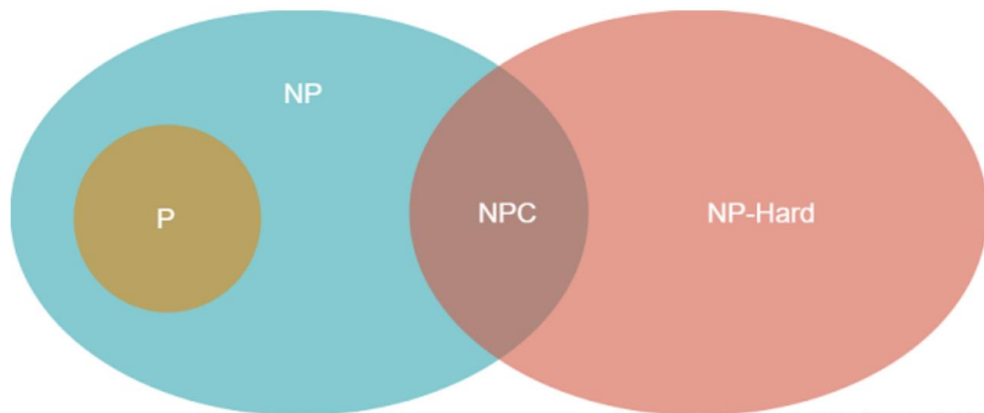
 - 复杂性问题的分类: $P=NP?$

 - 抽象复杂性研究

• 计算复杂性理论： $P=NP?$

引申知识——NP完全问题的魅力：

- 迄今为止没有对一个NP完全问题的有效算法，但也没有证否；
- NP完全问题集的非凡性质：如果任何一个NP完全问题存在有效算法，那么所有NP完全问题都存在有效算法；
- 有一些NP完全问题类似于（又不完全等同于）一些有着已知有效算法的问题



- **算法设计和分析**
 - 可计算问题的算法的设计与分析
 - 设计算法的理论、方法和技术
 - 分析算法的理论、方法和技术
- **计算机软件**
 - 系统软件
 - 工具软件
 - 应用软件



本讲内容

- 1.1 什么是算法?
- 1.2 计算机科学中算法的位置
- 1.3 算法分析引论**
- 1.4 算法设计引论

算法的分析

- 通过分析（若干个方面）手段而非实验手段对算法有整体认识
- 主要包括算法正确性分析和复杂性分析



算法的正确性分析

- **算法正确性**

- 一个算法是正确的，如果它对于每一个输入都最终停止,而且产生正确的输出。

- **不正确算法：**

- ①不停止(在某个输入上)

- ②对所有输入都停止，但对某输入产生不正确结果

- **近似算法**

- ①对所有输入都停止

- ②产生近似正确的解或产生不多的不正确解

• 算法正确性证明

- 证明算法对所有输入都停止
- 证明对每个输入都产生正确结果

– 调试程序 ≠ 程序正确性证明:

**程序调试只能说明程序是否有错,
不能证明程序有无错误!**



插入排序的正确性

- 循环不变量

- 在每次循环的开始, 子数组 $A[1..j-1]$ 包含原来数组中 $A[1..j-1]$, 但是已经有序

- 证明

- 初始化: $j=2$, $A[1..j-1]=A[1..1]=A[1]$, 已经有序.
- 维护: 每一层循环维护循环不变量.
- 终止: $j=n+1$, so $A[1..j-1]=A[1..n]$ 有序.

算法的复杂性分析

- **目的：**
 - 预测算法对不同输入所需资源量
- **复杂性测度：**
 - 时间，空间， I/O， 通信， 能源等，
 - 是输入大小的函数
- **用途：**
 - 为求解一个问题选择最佳算法、最佳设备
- **需要的数学基础**
 - 离散数学， 组合数学， 概率论， 代数等
- **需要的数学能力**
 - 建立算法复杂性的数学模型
 - 数学模型化简

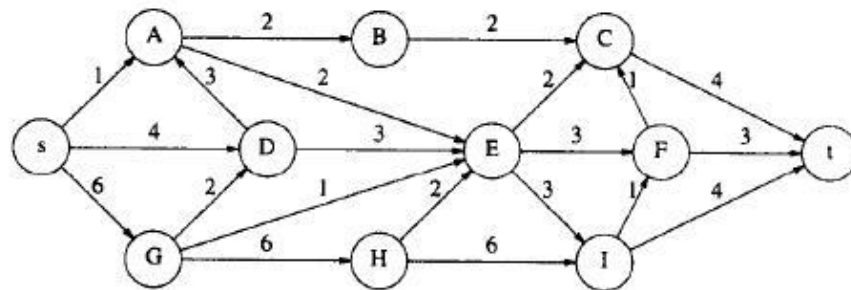
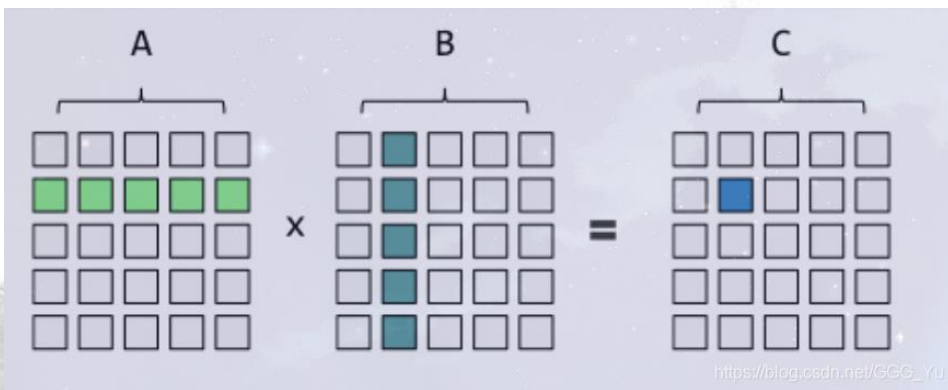
算法复杂性分析的度量

- 输入的大小

- 设Input是问题R的输入集合，R的输入大小是一个函数 $F: \text{Input} \rightarrow N$, N 是正整数集合。

示例:

- 矩阵问题的输入大小=矩阵的维数
 - 图论问题的输入大小=图的边数/结点数



算法复杂性分析的度量

- **时间复杂性**

- 一个算法对特定输入的时间复杂性是该算法对该输入产生结果需要的原子操作或“步”数

- **注意**

- 时间复杂性是输入大小的函数
 - 我们假设每一步的执行需要常数时间，实际上每步需要的时间量可能不同。



算法复杂性分析的度量

- **空间复杂性**

- 一个算法对特定输入的空间复杂性是该算法对该输入产生结果所需要的存储空间大小。
- （在算法运行的整个过程中，所需要的最大的内存的大小）



算法复杂性分析的度量

- 最坏复杂性

- 设Input是问题R的输入集合, $Complexity(X)$ 是求解R的算法A的复杂性函数, $Size(y)$ 是确定R中输入大小的函数, A的最坏复杂性是

$$\text{Max}\{Complexity(size(y)) \mid y \in \text{Input}\}$$

- 最小复杂性

$$\text{Min}\{Complexity(size(y)) \mid y \in \text{Input}\}$$

- 平均复杂性

- 设 $y \in \text{Input}$, y作为算法A的输入出现的概率是 p_y , A的平均复杂性为

$$\sum_{y \in \text{Input}} p_y \times complexity(size(y))$$

算法分析的模型

- 随机访问模型 (*Random-Access-Model, RAM*)
 - 单处理机，串行执行，无并发
 - 基本数据类型
 - 基本操作(每个操作常数时间)
- 并行多处理机模型(*PRAM*)



插入排序的分析

INSERTION-SORT(A)

```
1.  for  $j = 2$  to length[A]
2.      do  $key \leftarrow A[j]$ 
3.      //insert  $A[j]$  to sorted sequence  $A[1..j-1]$ 
4.       $i \leftarrow j-1$ 
5.      while  $i > 0$  and  $A[i] > key$ 
6.          do  $A[i+1] \leftarrow A[i]$ 
7.           $i \leftarrow i-1$ 
8.       $A[i+1] \leftarrow key$ 
```

代价 次数

c_1	n
c_2	$n-1$
0	$n-1$
c_4	$n-1$
c_5	$\sum_{j=2}^n t_j$
c_6	$\sum_{j=2}^n (t_j-1)$
c_7	$\sum_{j=2}^n (t_j-1)$
c_8	$n-1$

(t_j 是对 j 来说循环中执行的次数)

总时间代价 $T(n) =$ 代价的和 \times 每行执行次数

$$= c_1 n + c_2(n-1) + c_4(n-1) + c_5 \sum_{j=2}^n t_j + c_6 \sum_{j=2}^n (t_j-1) + c_7 \sum_{j=2}^n (t_j-1) + c_8(n-1)$$

插入排序的分析(续)

- 最好代价: 有序的数组

- $t_j=1$, 且6和7行执行0次

- $T(n) = c_1n + c_2(n-1) + c_4(n-1) + c_5(n-1) + c_8(n-1)$
 $= (c_1 + c_2 + c_4 + c_5 + c_8)n - (c_2 + c_4 + c_5 + c_8) = cn + c'$

- 最坏代价: 逆序数组

- $t_j=j$,

- $\sum_{j=2}^n t_j = \sum_{j=2}^n j = n(n+1)/2 - 1$, 且 $\sum_{j=2}^n (t_j - 1) = \sum_{j=2}^n (j - 1) = n(n-1)/2$

- $T(n) = c_1n + c_2(n-1) + c_4(n-1) + c_5(n(n+1)/2 - 1) + c_6(n(n-1)/2) + c_7(n(n-1)/2) + c_8(n-1)$
 $= ((c_5 + c_6 + c_7)/2)n^2 + (c_1 + c_2 + c_4 + c_5/2 - c_6/2 - c_7/2 + c_8)n - (c_2 + c_4 + c_5 + c_8)$
 $= an^2 + bn + c$

- 平均代价: 随机数

- 平均来看, $t_j = j/2$. $T(n)$ 与 n^2 同阶, 和最坏情况相同.

本讲内容

- 1.1 什么是算法?
- 1.2 计算机科学中算法的位置
- 1.3 算法分析引论
- 1.4 算法设计引论**

算法设计模式

- 如何把算法写出？（算法设计和生活中的思维类比）
- 暴力搜索
- 分治法
- 图搜索与枚举
 - 分支界限
 - 回溯
- 随机化方法

算法实现方法

- 递归与迭代
- 顺序、并行与分布式
- 确定性与非确定性
- 近似求解与精确求解
- 量子算法

算法实现方法：

- 顺序、并行与分布式

- 处理器功率密度随时钟速度超线性增加，时钟速度过快会导致芯片有熔化危险
- 芯片被设计成多个处理“核”，为每秒执行更多计算

- 并行vs分布式

- 并行：同任务、紧耦合
- 分布式：不同任务、松耦合

最优化算法设计方法

- 线性规划
- 动态规划
- 贪心法
- 启发式方法



本章小结

- 算法满足的五个性质
 - 终止性、确定性、能行性、输入、输出
- 问题P的一个实例是P的一个二元组，一个算法面向一个问题，而不是仅求解一个问题的某个或某几个实例
- 算法语言——伪代码

本章小结

- 承上启下的算法：
 - 可计算理论、算法复杂性理论
 - 程序设计、软件设计与开发



本章小结

- 算法分析的两个重要方面：
 - 正确性分析：能停止、正确的解
 - 复杂性分析：对所消耗资源的分析

