

DATA STRUCTURES AND ALGORITHMS

Textbook:

*Fundamentals of Data Structure in C++,
second edition, Silicon Press*

Instructor:

吕建华, COSE

Email: lujianhua@seu.edu.cn

地 址: 计算机楼230

Teaching assistants:

刘明瑶 : 1298639127@qq.com

吴景贤 : 79039140@qq.com

Total class hours: 64

week 1-16

Total lab. Hours: 16

**Every Wednesday Evening of Week 4, 7, 11,
14, 6:00 – 9:30, Room 268(60), 235(?),262(?)**

Assignments and Projects:

- **Should be handed to teaching assistants.**
- **Deadline: In TWO weeks after assignments.**

Evaluation:

Course Attendance: 10%,

Exercises and Projects: 20%,

**Final Examination (Textbook and Course
Notes allowed): 70%**

References:

- 1 金远平, 数据结构(C++描述), 清华大学出版社, 2005
- 2 T. A. Standish, Data Structures, Algorithms & Software Principles in C, Addison-Wesley Publishing Company, 1994

Tips

- **Make good use of your time in class**
 - **Listening**
 - **Thinking**
 - **Taking notes**
- **Utilize your free time**
 - **Go over**
 - **Programing**
- **Take a pen and some paper with you**
 - **Notes**
 - **Exercises**

Prerequisites:

Programming Language: C, C++

Pointer in C & C++

In Computer science, a **data structure** is a particular way of storing and organizing data in a computer so that it can be used **efficiently**.

物有本末，事有终始。
知所先后，则近道矣。

Sorting

- Rearrange $a[0], a[1], \dots, a[n-1]$ into ascending order. When done, $a[0] \leq a[1] \leq \dots \leq a[n-1]$
- $8, 6, 9, 4, 3 \Rightarrow 3, 4, 6, 8, 9$

Sort Methods

- **Insertion Sort**
- Bubble Sort
- Selection Sort
- Counting Sort
- Shell Sort
- Heap Sort
- Merge Sort
- Quick Sort
-

Insert An Element

- Given a sorted list/sequence, insert a new element
- Given 3, 6, 9, 14
- Insert 5
- Result 3, 5, 6, 9, 14

Insert an Element

- 3, 6, 9, 14 insert 5
- Compare new element (5) and last one (14)
- Shift 14 right to get 3, 6, 9, , 14
- Shift 9 right to get 3, 6, , 9, 14
- Shift 6 right to get 3, , 6, 9, 14
- Insert 5 to get 3, 5, 6, 9, 14

Insert An Element

```
// insert t into a[0:i-1]
int j;
for (j = i - 1; j >= 0 && t < a[j]; j--)
    a[j + 1] = a[j];
a[j + 1] = t;
```

Insertion Sort

- Start with a sequence of size 1
- Repeatedly insert remaining elements

Insertion Sort

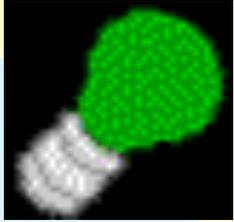
- Sort 7, 3, 5, 6, 1
- Start with 7 and insert 3 \Rightarrow 3, 7
- Insert 5 \Rightarrow 3, 5, 7
- Insert 6 \Rightarrow 3, 5, 6, 7
- Insert 1 \Rightarrow 1, 3, 5, 6, 7

Insertion Sort

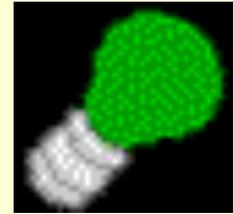
```
for (int i = 1; i < a.length; i++)  
{  
  // insert a[i] into a[0:i-1]  
  // code to insert comes here  
}
```

Insertion Sort

```
for (int i = 1; i < a.length; i++)  
  {  
    // insert a[i] into a[0:i-1]  
    int t = a[i];  
    int j;  
    for (j = i - 1; j >= 0 && t < a[j]; j--)  
      a[j + 1] = a[j];  
    a[j + 1] = t;  
  }
```



Insertion Sort



```
for (int i = 1; i < a.length; i++)  
  {  
    // insert a[i] into a[0:i-1]  
    int t = a[i];  
    int j;  
    for (j = i - 1; j >= 0 && t < a[j]; j--)  
      a[j + 1] = a[j];  
    a[j + 1] = t;  
  }
```

Basic Concepts

Purpose:

Provide the tools and techniques necessary to design and implement **large-scale software systems**, including:

- Data abstraction and encapsulation
- Algorithm specification and design
- Performance analysis and measurement
- Recursive programming

Overview: System Life Cycle

(1) Requirements

specifications of purpose

input

output

(2) Analysis

break the problem into manageable pieces

bottom-up

top-down

Overview: System Life Cycle

(3) Design

a **SYSTEM?** (from the designer's angle)

data objects

operations on them

TO DO

abstract data type

algorithm specification and design

Example: scheduling system of university

??

??

(4) Refinement and coding
representations for data object
algorithms for operations
components reuse

(5) Verification and maintenance
correctness proofs
testing
error removal
update

Data Abstraction and Encapsulation

Data Encapsulation or information Hiding is the concealing of the implementation details of a data object from the outside world.

Data Abstraction is the separation between the *specification* of a data object and its *implementation*.

DVD example.

A Data Type is a collection of *objects* and a set of *operations* that act on those objects.

predefined and user-defined:

char, int, arrays, structs, classes.

An Abstract Data Type (ADT) is a data type with the specification of the objects and the specification of the operations on the objects being **separated** from the representation of the objects and the implementation of the operations.

Benefits of data abstraction and data encapsulation:

(1) Simplification of software development

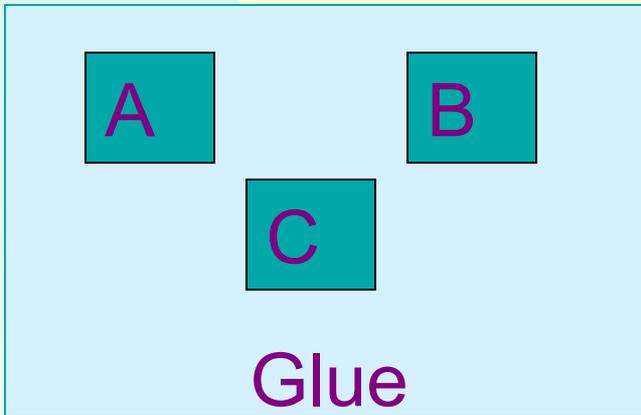
Applicaton : **data types A, B, C & Code Glue**

(a) a team of 4 programmers

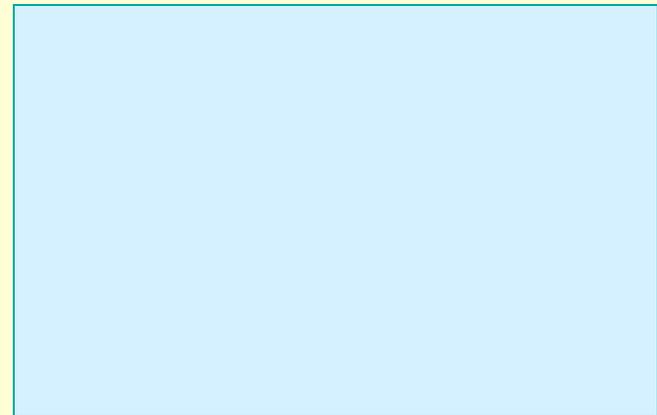
(b) a single programmer

Testing and debugging

Code with data abstraction



Code without data abstraction



Unshaded areas represent code to be searched for bugs.

(3) Reusability

data structures implemented as distinct entities of a software system

(4) Modifications to the representation of a data type

a change in the internal implementation of a data type will not affect the rest of the program as long as its interface does not change.

Algorithm Specification

An **algorithm** is finite set of instructions that, if followed, accomplishes a particular task.

Must satisfy the following criteria:

- (1) Input** Zero or more quantities externally supplied.
- (2) Output** At least one quantity is produced.
- (3) Definiteness** Clear and unambiguous.
- (4) Finiteness** Terminates after a finite number of steps.
- (5) Effectiveness** Basic enough, feasible

Compare: algorithms and programs

Finiteness

Recursion

Exercises: P32-2, P33-14

Performance Analysis and Measurement

Definition:

The **Space complexity** of a program is the amount of memory it needs to run to completion.

The **Time complexity** of a program is the amount of computer time it needs to run to completion.

(1) Priori estimates --- **Performance analysis**

(2) Posteriori testing--- **Performance measurement**

Performance Analysis

Space complexity

The space requirement of program P:

$$S(P) = c + S_p(\text{instance characteristics})$$

We concentrate solely on S_p .

Performance Analysis

Example 1.10

```
float Rsum (float *a, const int n) //compute  $\sum_{i=0}^{n-1} a[i]$ 
recursively
{
    if (n <=0) return 0;
    else return (Rsum(a,n-1)+a[n-1]);
}
```

The instances are characterized by

n

each call requires 4 words (n, a, return value, return address)

the depth of recursion is

n+1

$S_{\text{rsum}}(n) =$

$4(n+1)$

Time complexity

Run time of a program P:

$$T(P) = c + t_p(\text{instance characteristics})$$

A **program step** is loosely defined as a syntactically or semantically meaningful segment of a program that has an execution time that is **independent** of instance characteristics.

In P41-43 of the textbook, there is an detailed assignment of step counts to statements in C++.

Step Count

A step is an amount of computing that does not depend on the instance characteristic n

10 adds, 100 subtracts, 1000 multiplies
can all be counted as a single step

n adds cannot be counted as 1 step

Our main concern:

how many steps are needed by a program to solve a particular problem instance?

2 ways:

(1) count

(2) table

Example 1.12

```
count=0;
float Rsum (float *a, const int n)
{
    count++; // for if
    if (n <=0) {
        count++; // for return
        return 0;
    }
    else {
        count++; // for return
        return (Rsum(a,n-1)+a[n-1]);
    }
}
```

$$\begin{aligned}t_{\text{Rsum}}(0) &= 2, \\t_{\text{Rsum}}(n) &= 2 + t_{\text{Rsum}}(n-1) \\&= 2 + 2 + t_{\text{Rsum}}(n-2) \\&\cdot \\&\cdot \\&\cdot \\&= 2n + t_{\text{Rsum}}(0) \\&= 2n + 2\end{aligned}$$

Examp Let us use a table to count its total steps.

```
1 void F
```

	Line	s/e	frequency	total steps
2 { // cc	1	0	1	0
3 if (n	2	0	1	0
4 else	3	1 (n > 1)	1	1
5 in	4	0	1	0
6 fo	5	2	1	2
7 {	6	1	n	n
8	7	0	n-1	0
9	8	1	n-1	n-1
10	9	1	n-1	n-1
11 } //				
12 c				
13 } //				
14 }				

10	1	n-1	n-1
11	0	n-1	0
12	1	1	1
13	0	1	0
14	0	1	0

So

for $n > 1$, $t_{\text{Fibonacci}}(n) = 4n + 1$,

for $n = 0$ or 1 , $t_{\text{fibonacci}}(n) = 2$

Sometime, the instance characteristics is related with the content of the input data set.

e.g., BinarySearch.

Hence:

- **best-case**
- **worst-case,**
- **average-case.**

Asymptotic Notation

Because of the inexactness of what a step stands for, we are mainly concerned with the magnitude of the number of steps.

Definition [O]: $f(n)=O(g(n))$ iff there exist positive constants c and n_0 such that $f(n) \leq c g(n)$ for all $n, n>n_0$.

Example 1.13: $3n+2=O(n), 6*2^n+n^2=O(2^n),\dots$

Note $g(n)$ is an **upper bound**.

$n = O(n^2)$, $n = O(2^n)$, ...,

for $f(n) = O(g(n))$ to be informative, $g(n)$ should be

as small as possible.

In practice, the coefficient of $g(n)$ should be 1. We never say $O(3n)$.

Theory 1.2: if $f(n) = a_m n^m + \dots + a_1 n + a_0$, then $f(n) = O(n^m)$.

When the complexity of an algorithm is actually, say, $O(\log n)$, but we can only show that it is $O(n)$ due to the limitation of our knowledge, it is OK to say so. This is one benefit of O notation as upper bound.

Self-study:

Ω --- low bound

Θ --- equal bound

A Few Comparisons

Function #1

Function #2

$$n^3 + 2n^2$$



$$100n^2 + 1000$$

$$n^{0.1}$$



$$\log n$$

$$n + 100n^{0.1}$$



$$2n + 10 \log n$$

$$5n^5$$



$$n!$$

$$n^{-15} 2^n / 100$$



$$1000n^{15}$$

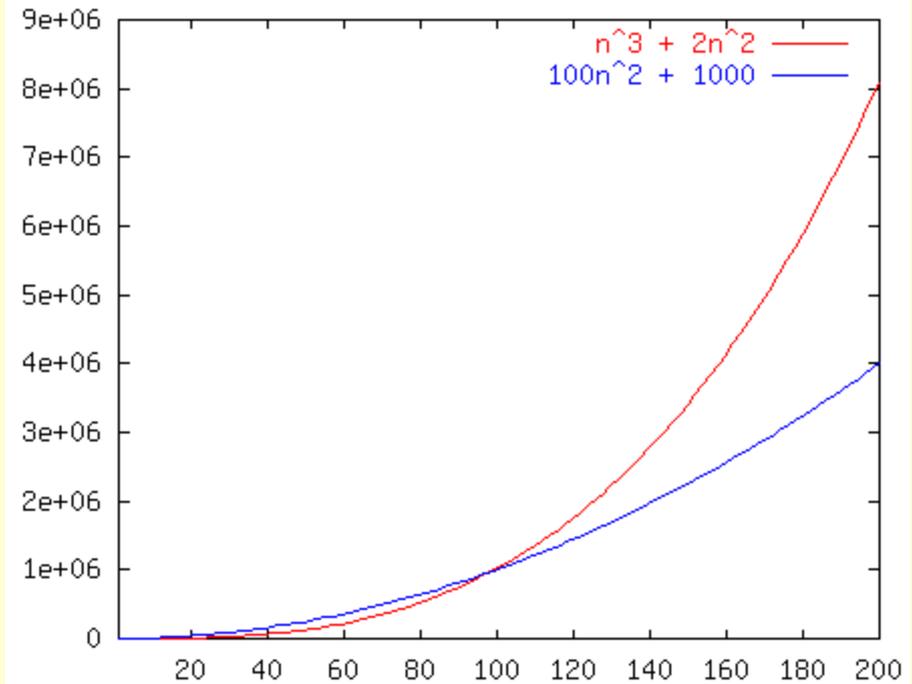
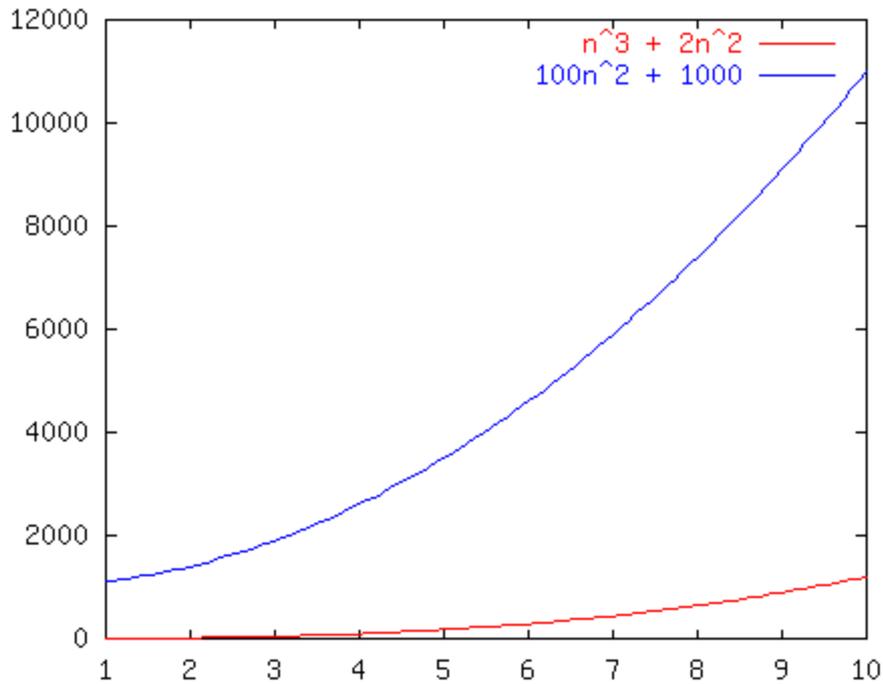
$$8^{2 \log n}$$



$$3n^7 + 7n$$

Race I

$n^3 + 2n^2$ vs. $100n^2 + 1000$

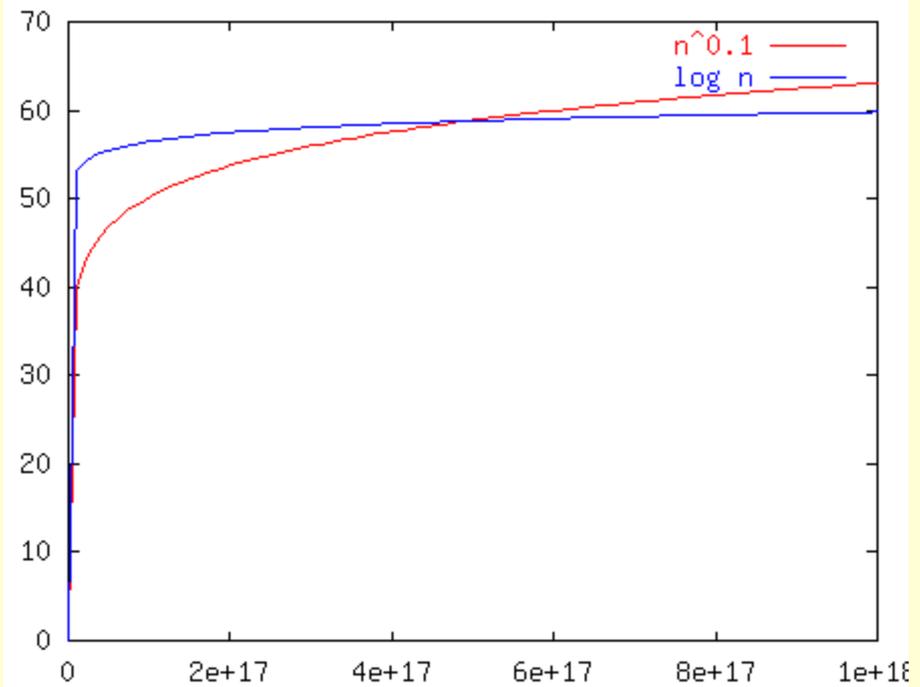
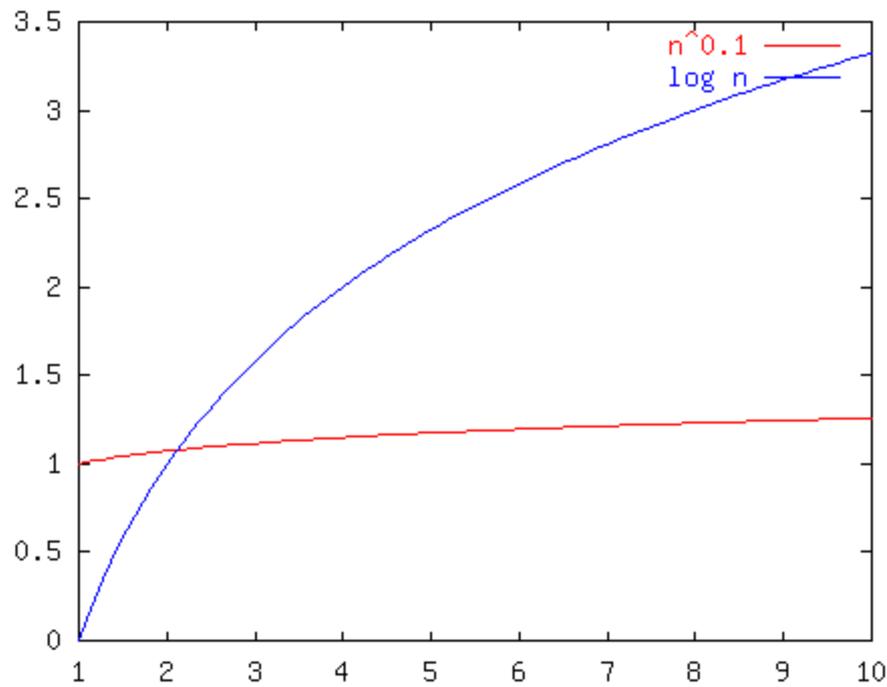


Race II

$n^{0.1}$

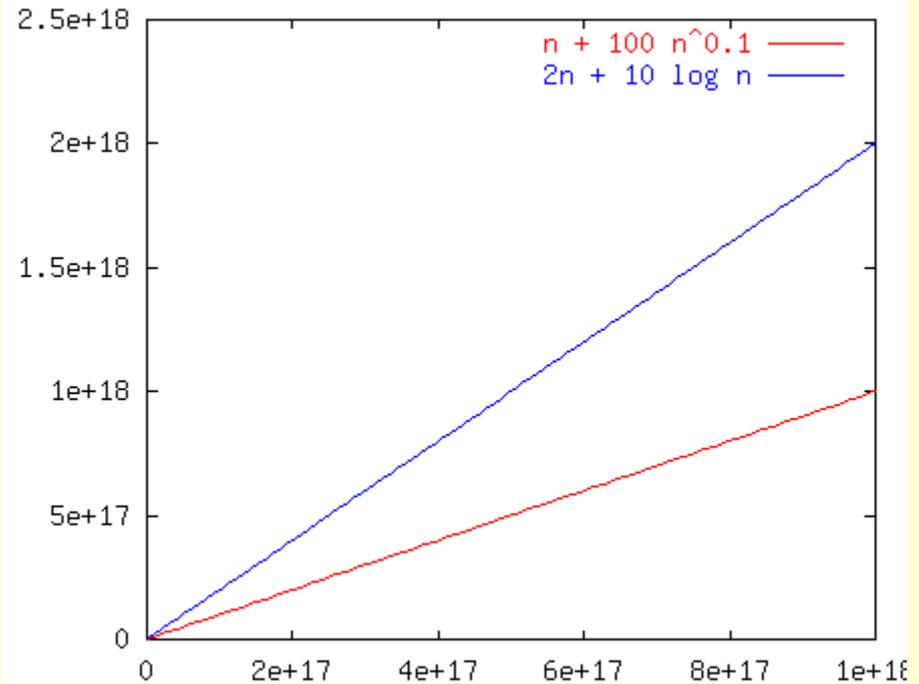
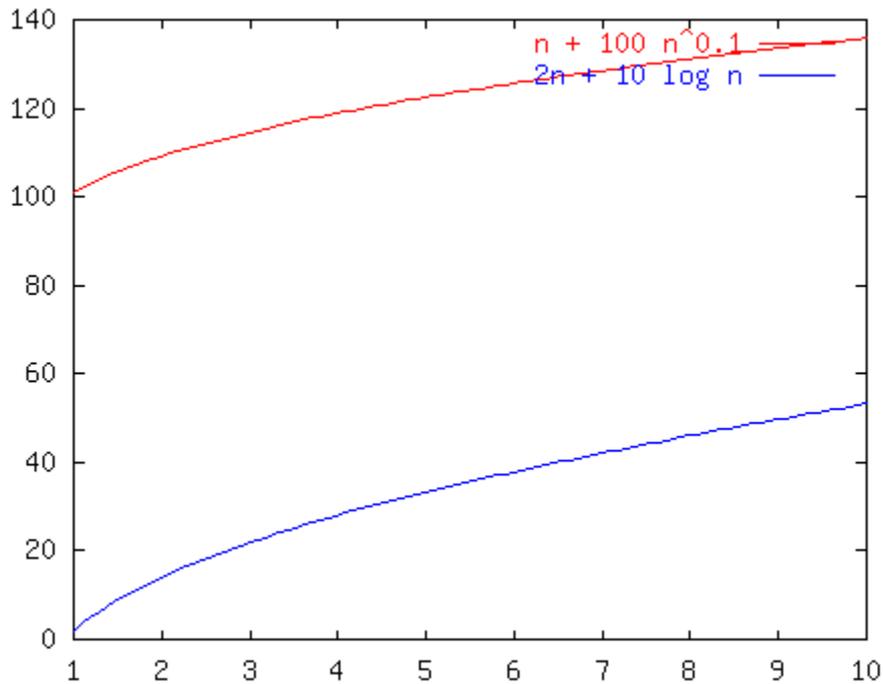
vs.

$\log n$



Race III

$n + 100n^{0.1}$ vs. $2n + 10 \log n$

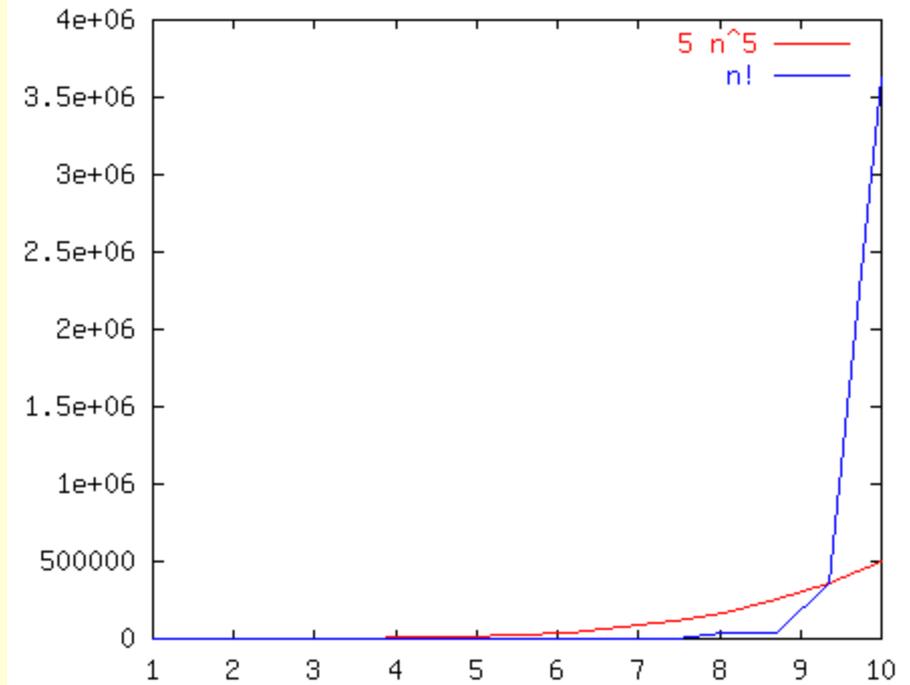
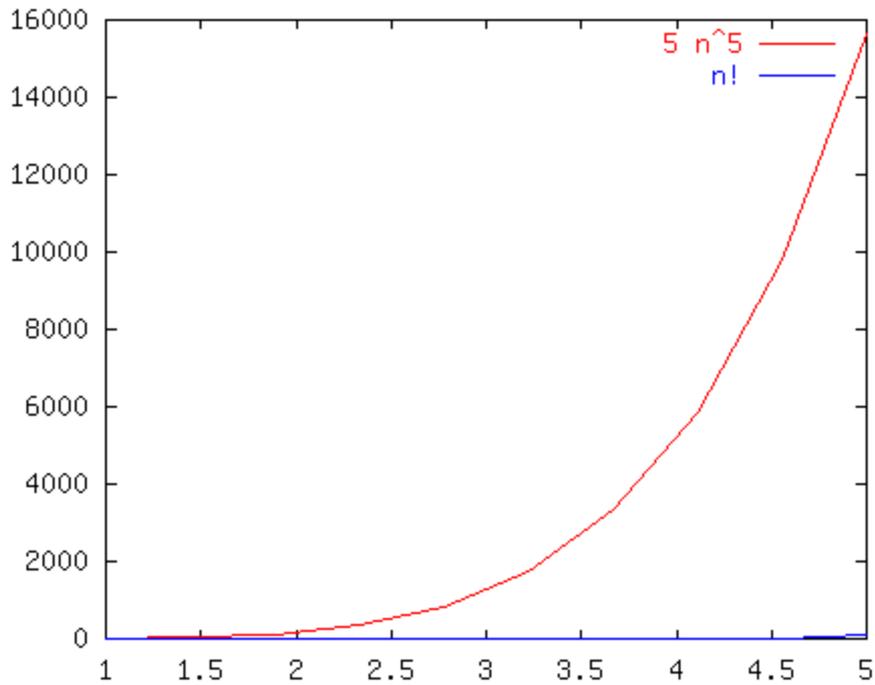


Race IV

$5n^5$

vs.

$n!$

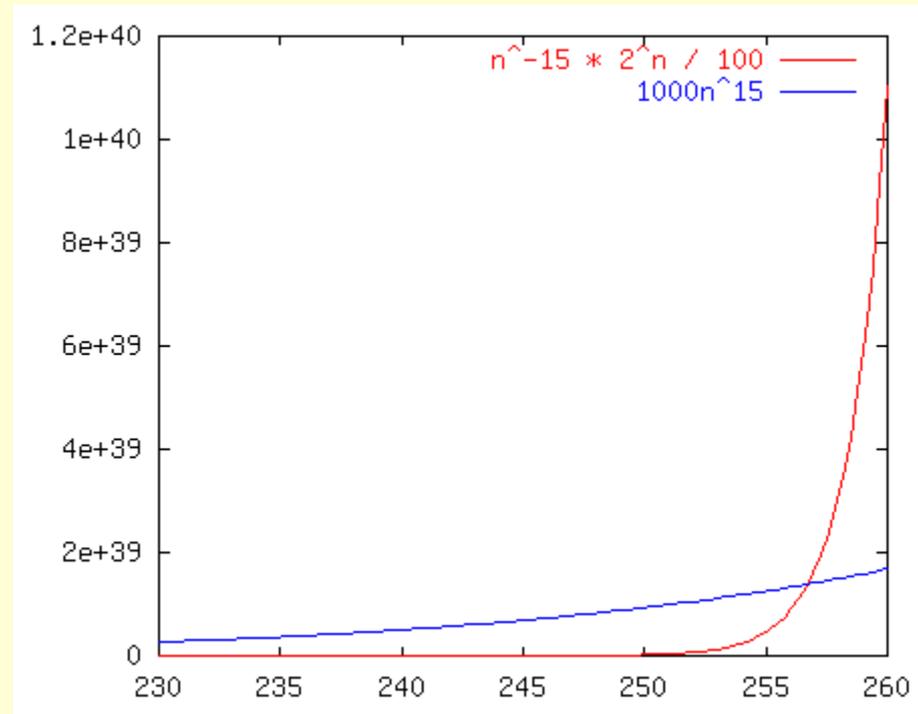
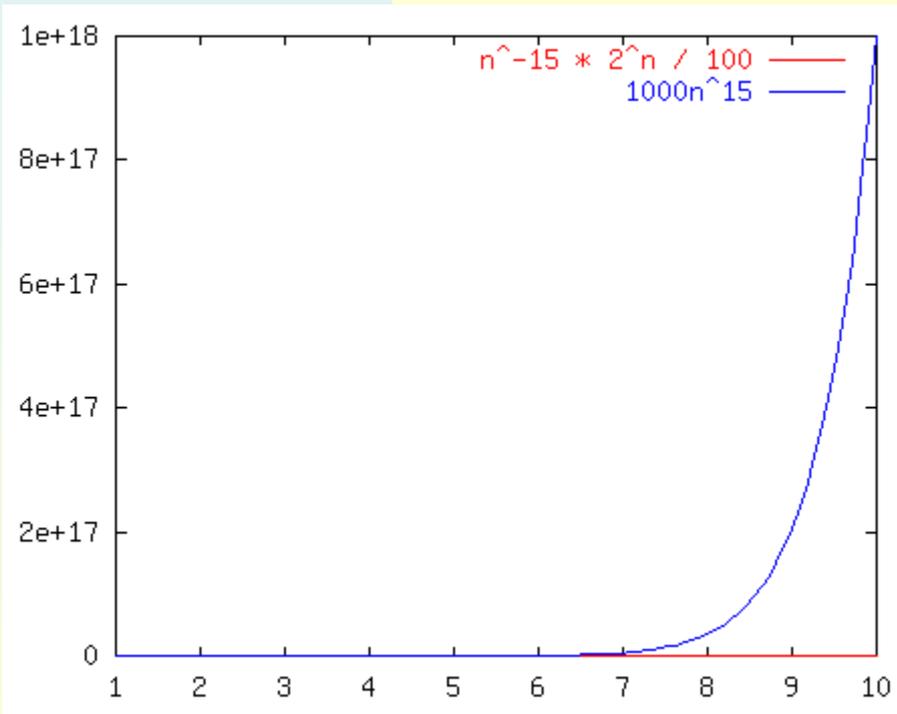


Race V

$$n^{-15} 2^n / 100$$

vs.

$$1000n^{15}$$

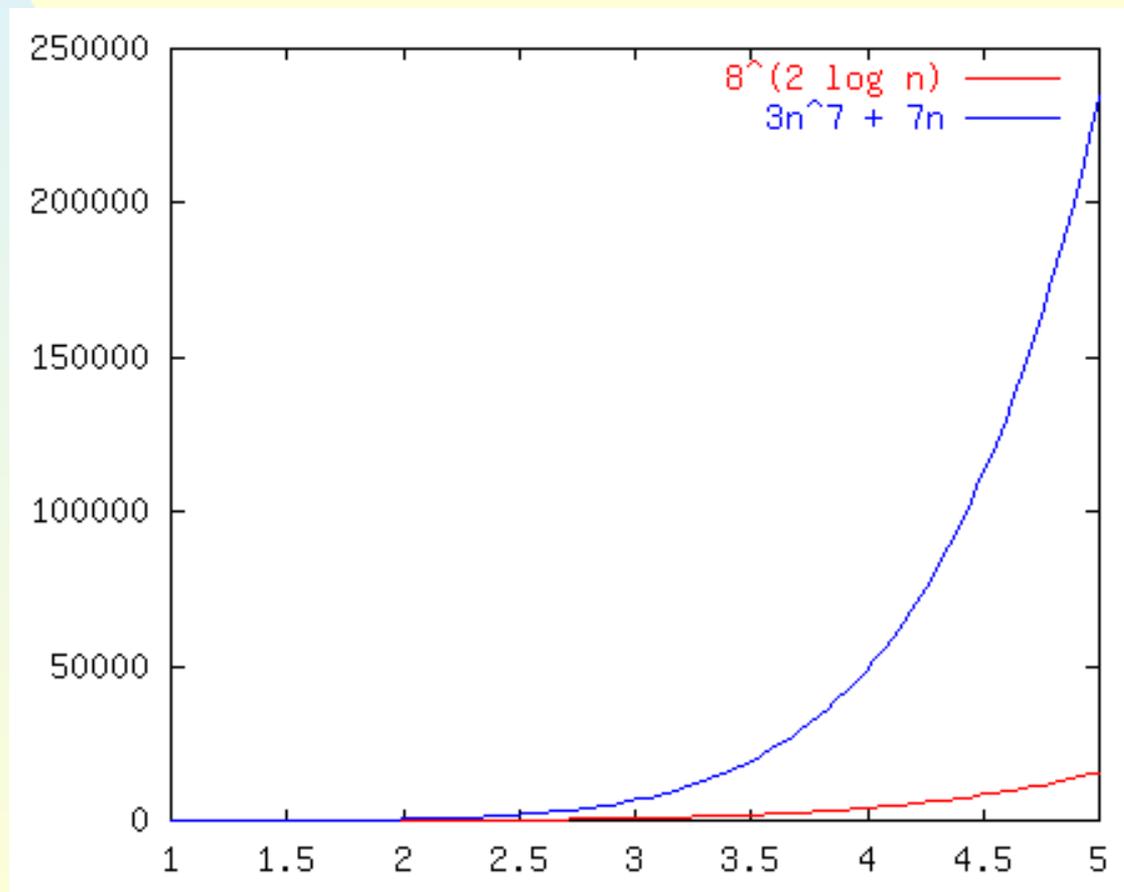


Race VI

$$8^{2 \log(n)}$$

vs.

$$3n^7 + 7n$$



The Losers Win

Function #1

$$n^3 + 2n^2$$

$$n^{0.1}$$

$$n + 100n^{0.1}$$

$$5n^5$$

$$n^{-15}2^n/100$$

$$8^{2\log n}$$

Function #2

$$100n^2 + 1000$$

$$\log n$$

$$2n + 10 \log n$$

$$n!$$

$$1000n^{15}$$

$$3n^7 + 7n$$

Better algorithm!

$$O(n^2)$$

$$O(\log n)$$

$$\mathbf{TIE} O(n)$$

$$O(n^5)$$

$$O(n^{15})$$

$$O(n^6)$$

Common Names

constant: $O(1)$

logarithmic: $O(\log n)$

linear: $O(n)$

log-linear: $O(n \log n)$

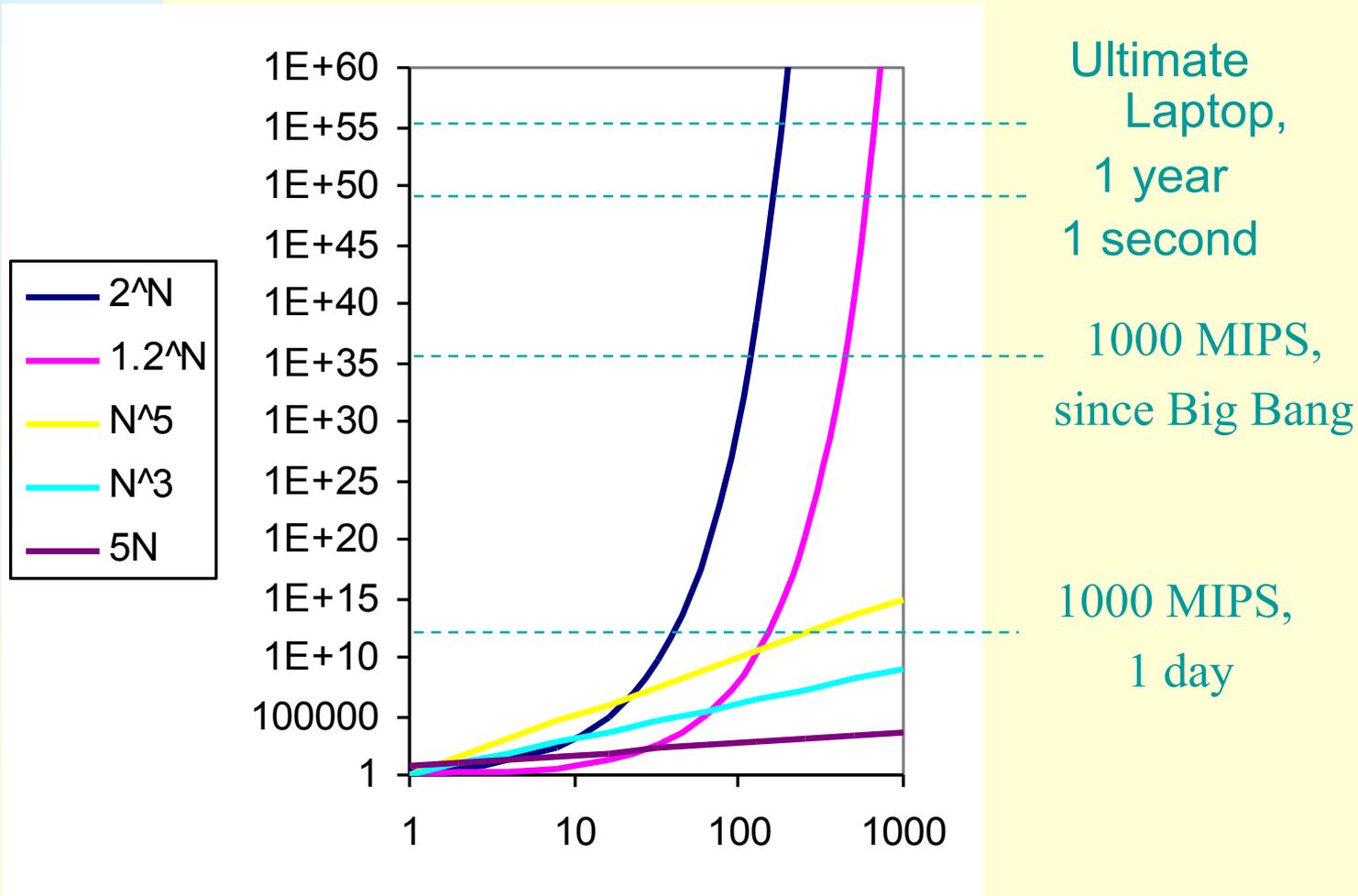
quadratic: $O(n^2)$

polynomial: $O(n^k)$ (k is a constant)

exponential: $O(c^n)$ (c is a constant > 1)

Practical Complexity

How the various functions grow with n?



n	$f(n)=n$	$f(n)=n\log_2 n$	$f(n)=n^2$	$f(n)=n^4$	$f(n)=n^{10}$	$f(n)=2^n$
10	.01 μ s	.03 μ s	.1 μ s	10 μ s	10s	1 μ s
20	.02 μ s	.09 μ s	.4 μ s	160 μ s	2.84h	1 ms
30	.03 μ s	.15 μ s	.9 μ s	810 μ s	6.83d	1 s
40	.04 μ s	.21 μ s	1.6 μ s	2.56ms	121d	18m
50	.05 μ s	.28 μ s	2.5 μ s	6.25ms	3.1y	13 d
100	.1 μ s	.66 μ s	10 μ s	100 ms	3171y	$4 \cdot 10^{13}$ y
10^3	1 μ s	9.66 μ s	1ms	16.67m		
10^4	10 μ s	130 μ s	100ms	115.7d		
10^5	100 μ s	1.66ms	10s	3171y		

Table 1.8: Times on a 1-billion-steps-per-second computer

Performance Measurement

Performance measurement is concerned with obtaining the **actual** space and time requirements of a program.

To time a short event it is necessary to **repeat** it several times and divide the total time for the event by the number of repetitions.

Let us look at the following program:

```
int SequentialSearch (int *a, const int n, const int x )  
{ // Search a[0:n-1].  
    int i;  
    for (i=0; i < n && a[i] != x; i++;)  
        if (i == n) return -1;  
    else return i;  
}
```

```

void TimeSearch ( )
{
    int a[1000], n[20];
    const long r[20] = {300000, 300000, 200000, 200000,
    100000, 100000, 100000, 80000, 80000, 50000, 50000,
    25000, 15000, 15000, 10000, 7500, 7000, 6000, 5000,
    5000 };

    for ( int j=0; j<1000; j++ ) a[j] = j+1; //initialize a
    for ( j=0; j<10; j++ ) { //values of n
        n[j] = 10*j; n[j+10] = 100*(j+1 );
    }

    cout << “ n  total  runTime” << endl;
}

```

```

for ( j=0; j<20; j++ ) {
    long start, stop;
    time (&start);           // start timer
    for ( long b=1; b<=r[j]; b++ )
        int k = seqsearch(a, n[j], 0 ); //unsuccessful search
    time (&stop);           // stop timer
    long totalTime = stop - start;
    float runTime = (float) (totalTime) / (float)(r[j]);
    cout << " " << n[j] << " " << totalTime << " " << runTime
        << endl;
    }
}

```

The results of running *TimeSearch* are as in the next slide.

n	total	runTime	n	total	runTime
0	241	0.0008	100	527	0.0105
10	533	0.0018	200	505	0.0202
20	582	0.0029	300	451	0.0301
30	736	0.0037	400	593	0.0395
40	467	0.0047	500	494	0.0494
50	565	0.0056	600	439	0.0585
60	659	0.0066	700	484	0.0691
70	604	0.0075	800	467	0.0778
80	681	0.0085	900	434	0.0868
90	472	0.0094	1000	484	0.0968

Times in hundredths of a second, the plot of the data can be found in Fig. 1.7.

Issues to be addressed:

- (1) Accuracy of the clock**
- (2) Repetition factor**
- (3) Suitable test data for worst-case or average performance**
- (4) Purpose: comparing or predicting?**
- (5) Fit a curve through points**

Exercises:

P72-10