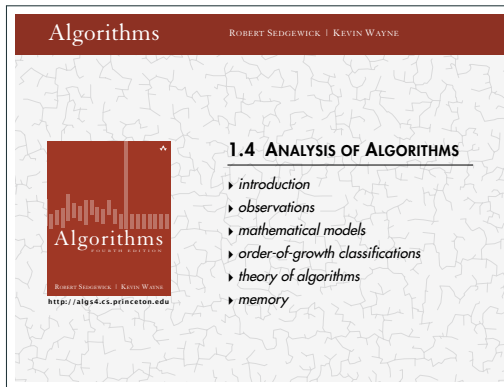# CS2010: ALGORITHMS AND DATA STRUCTURES

## Lecture 2: Analysis of Algorithms

Vasileios Koutavas
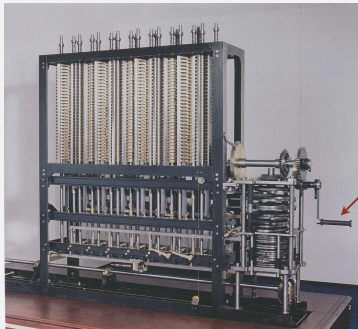
School of Computer Science and Statistics
Trinity College Dublin

→ Parts from S&W 1.4
→ Estimate the performance of algorithms by
  → Experiments & Observations
  → Precise Mathematical Calculations

" *As soon as an Analytic Engine exists, it will necessarily guide the future course of the science. Whenever any result is sought by its aid, the question will arise—By what course of calculation can these results be arrived at by the machine in the shortest time?* " — *Charles Babbage (1864)*



how many times do you have to turn the crank?

**Analytic Engine**

→ **Good programmer:** to predict the performance of our programs.

→ **Good client:** to choose between alternative algorithms/implementations.

→ **Good manager:** to provide guarantees to clients / avoid client complaints.

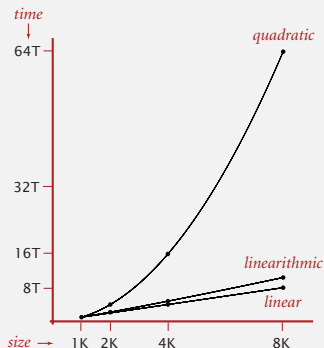→ **Good theoritician:** to understand the nature of computing.

## Some algorithmic successes

Discrete Fourier transform.
- Break down waveform of $N$ samples into periodic components.
- Applications: DVD, JPEG, MRI, astrophysics, ....
- Brute force: $N^2$ steps.
- FFT algorithm: $N \log N$ steps, enables new technology.
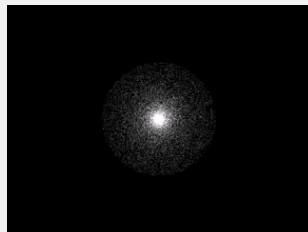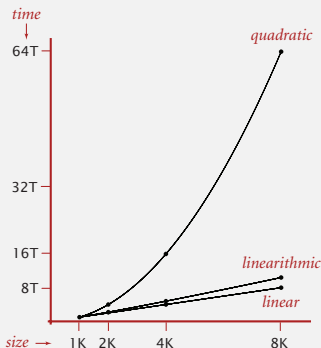
**Friedrich Gauss**
**1805**

## Some algorithmic successes

N-body simulation.

- Simulate gravitational interactions among $N$ bodies.
- Brute force: $N^2$ steps.
- Barnes-Hut algorithm: $N \log N$ steps, enables new research.

**Andrew Appel**
**PU '81**

Q. Will my program be able to solve a large practical input?



Insight. [Knuth 1970s] Use scientific method to understand performance.

## Scientific method applied to analysis of algorithms

A framework for predicting performance and comparing algorithms.

Scientific method.
- Observe some feature of the natural world.
- Hypothesize a model that is consistent with the observations.
- Predict events using the hypothesis.
- Verify the predictions by making further observations.
- Validate by repeating until the hypothesis and observations agree.

Principles.
- Experiments must be reproducible.
- Hypotheses must be falsifiable.

Feature of the natural world. Computer itself.

# Experimental Approach:

# Measuring Precise Running Time

3-Sum. Given $N$ distinct integers, how many triples sum to exactly zero?

```
% more 8ints.txt
8
30 -40 -20 -10 40 0 10 5

% java ThreeSum 8ints.txt
4
```

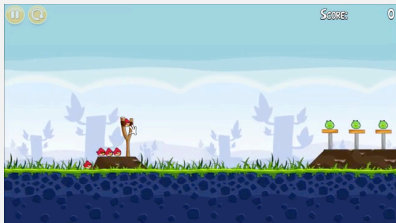| | a[i] | a[j] | a[k] | sum |
|---|---|---|---|---|
| 1 | 30 | -40 | 10 | 0 |
| 2 | 30 | -20 | -10 | 0 |
| 3 | -40 | 40 | 0 | 0 |
| 4 | -10 | 0 | 10 | 0 |



Context. Deeply related to problems in computational geometry.

## 3-Sum: brute-force algorithm

```
public class ThreeSum
{
   public static int count(int[] a)
   {
      int N = a.length;
      int count = 0;
      for (int i = 0; i < N; i++)
         for (int j = i+1; j < N; j++)
            for (int k = j+1; k < N; k++)         ←——— check each triple
               if (a[i] + a[j] + a[k] == 0)       ←——— for simplicity, ignore
                  count++;                              integer overflow
      return count;
   }

   public static void main(String[] args)
   {
      In in = new In(args[0]);
      int[] a = in.readAllInts();
      StdOut.println(count(a));
   }
}
```

The input of `ThreeSum` is an array of size *N*.

Suppose we care only about 100-element arrays.

There are many different 100-element arrays.

```java
public class ThreeSum
{
   public static int count(int[] a)
   {
      int N = a.length;
      int count = 0;
      for (int i = 0; i < N; i++)
         for (int j = i+1; j < N; j++)
            for (int k = j+1; k < N; k++)
               if (a[i] + a[j] + a[k] == 0)
                  count++;
      return count;
   }
}
```

The input of `ThreeSum` is an array of size *N*.

Suppose we care only about 100-element arrays.

There are many different 100-element arrays.

Q. Is the running time of `ThreeSum` dependent on which 100-element array we provide as input?

```java
public class ThreeSum
{
   public static int count(int[] a)
   {
      int N = a.length;
      int count = 0;
      for (int i = 0; i < N; i++)
         for (int j = i+1; j < N; j++)
            for (int k = j+1; k < N; k++)
               if (a[i] + a[j] + a[k] == 0)
                  count++;
      return count;
   }
}
```

## Measuring the running time

Q. How to time a program?

A. Manual.



```
% java ThreeSum 1Kints.txt
```



*tick tick tick*

70

```
% java ThreeSum 2Kints.txt
```



*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*

528

```
% java ThreeSum 4Kints.txt
```



*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*
*tick tick tick tick tick tick tick tick*

4039

## Measuring the running time

Q. How to time a program?

A. Automatic.

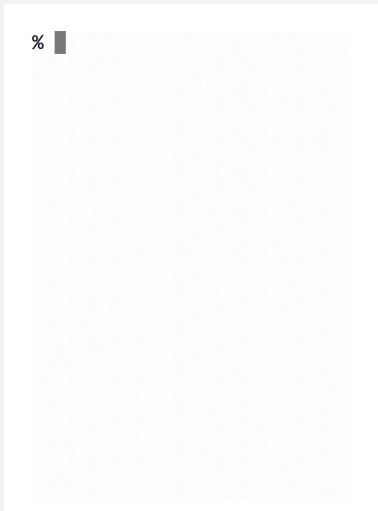| public class Stopwatch | (part of stdlib.jar) |
| --- | --- |
| Stopwatch() | *create a new stopwatch* |
| double elapsedTime() | *time since creation (in seconds)* |

```
public static void main(String[] args)
{
   In in = new In(args[0]);
   int[] a = in.readAllInts();
   Stopwatch stopwatch = new Stopwatch();
   StdOut.println(ThreeSum.count(a));
   double time = stopwatch.elapsedTime();
   StdOut.println("elapsed time " + time);
}
```

## Empirical analysis

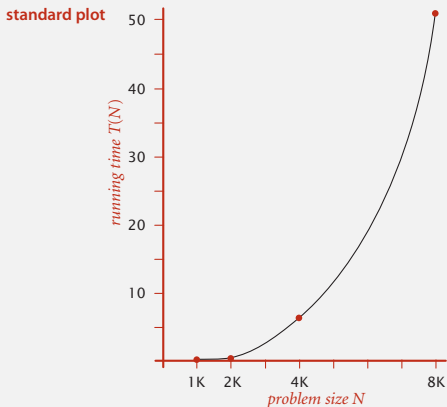Run the program for various input sizes and measure running time.

## Empirical analysis

Run the program for various input sizes and measure running time.

| N | time (seconds) [†] |
|---|---|
| 250 | 0.0 |
| 500 | 0.0 |
| 1,000 | 0.1 |
| 2,000 | 0.8 |
| 4,000 | 6.4 |
| 8,000 | 51.1 |
| 16,000 | ? |

## Data analysis

Standard plot. Plot running time $T(N)$ vs. input size $N$.



standard plot

*running time T(N)*

*problem size N*

## Data analysis

Log-log plot. Plot running time $T(N)$ vs. input size $N$ using log-log scale.



$\lg(T(N)) = b \lg N + c$
$b = 2.999$
$c = -33.2103$

$T(N) = a N^b$, where $a = 2^c$

Regression. Fit straight line through data points: $a N^b$.

Hypothesis. The running time is about $1.006 \times 10^{-10} \times N^{2.999}$ seconds.

Try out the experimental analysis:

https://docs.google.com/spreadsheets/d/
1WnihyK6g1pYdcT2ndZOqNNRkTitXkWKnOrTgCnM-bw8/edit?usp=sharing

## Prediction and validation

Hypothesis. The running time is about $1.006 \times 10^{-10} \times N^{2.999}$ seconds.

"order of growth" of running time is about $N^3$ [stay tuned]

Predictions.

- 51.0 seconds for $N = 8,000$.
- 408.1 seconds for $N = 16,000$.

Observations.

| N | time (seconds) [†] |
|---|---|
| 8,000 | 51.1 |
| 8,000 | 51.0 |
| 8,000 | 51.1 |
| 16,000 | 410.8 |

**validates hypothesis!**

## Doubling hypothesis

Doubling hypothesis.  Quick way to estimate $b$ in a power-law relationship.

Run program, doubling the size of the input.

| N | time (seconds) [†] | ratio | lg ratio |
|---|---|---|---|
| 250 | 0.0 | | – |
| 500 | 0.0 | 4.8 | 2.3 |
| 1,000 | 0.1 | 6.9 | 2.8 |
| 2,000 | 0.8 | 7.7 | 2.9 |
| 4,000 | 6.4 | 8.0 | 3.0 ← lg (6.4 / 0.8) = 3.0 |
| 8,000 | 51.1 | 8.0 | 3.0 |

$$\frac{T(2N)}{T(N)} = \frac{a(2N)^b}{aN^b}$$
$$= 2^b$$

seems to converge to a constant b ≈ 3

Hypothesis.  Running time is about $a\,N^b$ with $b =$ lg ratio.
Caveat.  Cannot identify logarithmic factors with doubling hypothesis.

## Doubling hypothesis

Doubling hypothesis. Quick way to estimate $b$ in a power-law relationship.

Q. How to estimate $a$ (assuming we know $b$) ?
A. Run the program (for a sufficient large value of $N$) and solve for $a$.

| N | time (seconds) [†] |
|-------|-------------------|
| 8,000 | 51.1 |
| 8,000 | 51.0 |
| 8,000 | 51.1 |

$51.1 = a \times 8000^3$
$\Rightarrow a = 0.998 \times 10^{-10}$

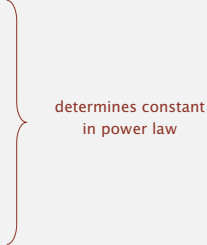Hypothesis. Running time is about $0.998 \times 10^{-10} \times N^3$ seconds.

almost identical hypothesis
to one obtained via linear regression

## Experimental algorithmics

System independent effects.

- Algorithm.
- Input data.

determines exponent
in power law

System dependent effects.

- Hardware:  CPU, memory, cache, …
- Software:  compiler, interpreter, garbage collector, …
- System:  operating system, network, other apps, …

determines constant
in power law

Bad news.  Difficult to get precise measurements.

Good news.  Much easier and cheaper than other sciences.

e.g., can run huge number of experiments

This was the experimental approach to algorithm analysis.

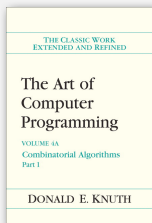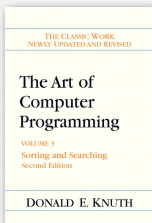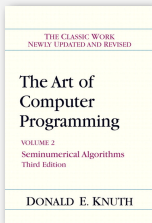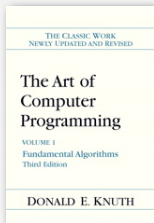Is there a mathematical approach where we can do **calculations** instead of experiments?

# Mathematical Approach 1:

# Calculating Precise Running Time

## Mathematical models for running time

Total running time:  sum of cost × frequency for all operations.

- Need to analyze program to determine set of operations.
- Cost depends on machine, compiler.
- Frequency depends on algorithm, input data.



**Donald Knuth
1974 Turing Award**

In principle, accurate mathematical models are available.

## Cost of basic operations

**Challenge.** How to estimate constants.

| operation | example | nanoseconds † |
|---|---|---|
| integer add | `a + b` | 2.1 |
| integer multiply | `a * b` | 2.4 |
| integer divide | `a / b` | 5.4 |
| floating-point add | `a + b` | 4.6 |
| floating-point multiply | `a * b` | 4.2 |
| floating-point divide | `a / b` | 13.5 |
| sine | `Math.sin(theta)` | 91.3 |
| arctangent | `Math.atan2(y, x)` | 129.0 |
| ... | ... | ... |

† Running OS X on Macbook Pro 2.2GHz with 2GB RAM

## Cost of basic operations

Observation. Most primitive operations take constant time.

| operation | example | nanoseconds † |
|-----------|---------|---------------|
| variable declaration | `int a` | $c_1$ |
| assignment statement | `a = b` | $c_2$ |
| integer compare | `a < b` | $c_3$ |
| array element access | `a[i]` | $c_4$ |
| array length | `a.length` | $c_5$ |
| 1D array allocation | `new int[N]` | $c_6 N$ |
| 2D array allocation | `new int[N][N]` | $c_7 N^2$ |

Caveat. Non-primitive operations often take more than constant time.

novice mistake: abusive string concatenation

## Example: 1-Sum

Q. How many instructions as a function of input size $N$?

```
int count = 0;
for (int i = 0; i < N; i++)
    if (a[i] == 0)
        count++;
```
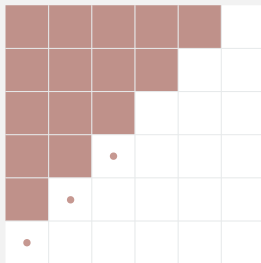
N array accesses

| operation | frequency |
|---|---|
| variable declaration | 2 |
| assignment statement | 2 |
| less than compare | $N + 1$ |
| equal to compare | $N$ |
| array access | $N$ |
| increment | $N$ to $2N$ |

### Example: 2-SUM

Q. How many instructions as a function of input size $N$?

```
int count = 0;
for (int i = 0; i < N; i++)
    for (int j = i+1; j < N; j++)
        if (a[i] + a[j] == 0)
            count++;
```
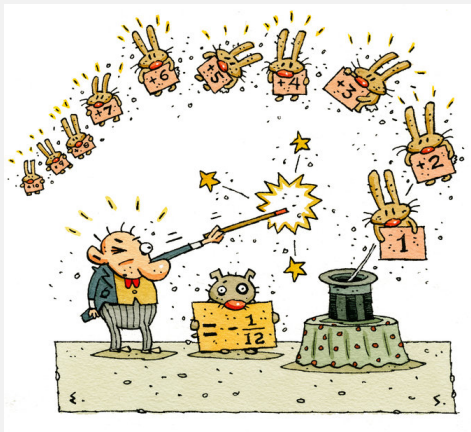
$$0 + 1 + 2 + \ldots + (N-1) = \frac{1}{2}N(N-1)$$
$$= \binom{N}{2}$$

Pf. [ n even]



$$0 + 1 + 2 + \ldots + (N-1) = \frac{1}{2}N^2 - \frac{1}{2}N$$

half of square     half of diagonal

$$1 + 2 + 3 + 4 + \ldots = -\frac{1}{12}$$



http://www.nytimes.com/2014/02/04/science/in-the-end-it-all-adds-up-to.html

## Example: 2-Sum

Q. How many instructions as a function of input size $N$?

```
int count = 0;
for (int i = 0; i < N; i++)
    for (int j = i+1; j < N; j++)
        if (a[i] + a[j] == 0)
            count++;
```

$$0 + 1 + 2 + \ldots + (N-1) \;=\; \frac{1}{2} N (N-1)$$
$$=\; \binom{N}{2}$$

| operation | frequency |
|---|---|
| variable declaration | $N + 2$ |
| assignment statement | $N + 2$ |
| less than compare | $\frac{1}{2} (N + 1)(N + 2)$ |
| equal to compare | $\frac{1}{2} N (N-1)$ |
| array access | $N(N-1)$ |
| increment | $\frac{1}{2} N(N-1)$ to $N(N-1)$ |

tedious to count exactly

30

$$T_N = c_1 A + c_2 B + c_3 C + c_4 D + c_5 E$$

Where

$c_1$ :cost of array access      $A$ :number of array accesses

$c_2$ :cost of integer addition      $B$ :number of integer additions

$c_3$ :cost of integer comparison      $C$ :number of integer comparisons

$c_4$ :cost of increment      $D$ :number of increments

$c_5$ :cost of assignment      $E$ :number of assignments

$$T_N = c_1A + c_2B + c_3C + c_4D + c_5E$$

Where

| | |
|---|---|
| $c_1$ :cost of array access | $A$ :number of array accesses |
| $c_2$ :cost of integer addition | $B$ :number of integer additions |
| $c_3$ :cost of integer comparison | $C$ :number of integer comparisons |
| $c_4$ :cost of increment | $D$ :number of increments |
| $c_5$ :cost of assignment | $E$ :number of assignments |

Q. Advantages / Disadvantages?