

## Analysis of Algorithms —

Between Mathematics and Computer Science Philippe Flajolet, INRIA Rocquencourt, F.





## FIRST,

A glimpse of history ...

Mathematics and Computing, i.e., *algorithms*, = a joint enterprise since the dawn of history.

Thesis: (i) Conceptual advances lead to more complex and efficient algorithms. (ii) Computer age obeys similar principle?

#### Rhind papyrus (ca 1650BC → 1900BC)





lihaxove:
MAXIMAXI
Souther
201201
11 8 AII 811100
leellmp475a

## Egyptians *knew* binary representations and technique of "binary powering"!

1 + 8 + 32					
$41 \times x =$	$(1+2^3)$	$(3+2^5)$	$\times x = 1x$	$x + 2^3$	$x+2^5x$ .
	41	×	59		
	1		59		
	2		118	v	
	4		236		
	8		472		
	16		944	v	
	32		1888		
			= 2419	<b>v</b>	
RSA, PGP	):				
$x^{41} = x^1 \cdot [x^2 \cdot x^4] \cdot x^8 \cdot [x^{16}] \cdot x^{32}.$					

rein. 01.2

The Rhind papyrus contains eighty-seven problems. The papyrus, a scroll about 6 metres long and 1/3 of a metre wide, was written around 1650 BC by the scribe Ahmes who states that he is copying a document which is 200 years older. © History of Mathematics archive @ St Andrews, UK.

## Computing <u>without</u> Computers!



Αγιατόν βΑ διχωτόγμας Ο μαροστόρ θωσιωτόν τόγμας ο ο ο όλο αρο Ν Ο του ζεξι Ο όλα χολι ο όμο το το το τον τον Γιαραι αι Δτ βΑ τόγμα ο στο τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι αι Δτ βΑ τόγμα ο το δι τον τον Γιαραι ο σι τον Γιαραι ο σι τον τον Γιαραι ο σι τον Γιαραι ο σι τον τον Γιαραι ο σι τον τον Γιαραι ο σι 

• Euclid (325BC–265BC) discovers *Euclid's algorithm* and formalizes *geometry*. Archimedes (287BC–212BC) discovers that  $\pi$  *is computable*; cf Viète (1540–1603):

$$\frac{\pi}{2} = \frac{2}{\sqrt{2}} \frac{2}{\sqrt{2+\sqrt{2}}} \frac{2}{\sqrt{2+\sqrt{2}+\sqrt{2}}} \cdots$$

Arithmetics & Algorithms



• <u>Al Kwarizmi</u> (780–850) gives complete set of rules, an "*algorithm"* for the four operations on "hindi" numerals.



 <u>Newton</u> (1643–1727) "De Methodis Serierum et Fluxionum" = Newton's algorithm; "computer algebra" Let P(x, y) = 0; determine  $\frac{d}{dx}y(x)$ ? Cf: Newton 1671; here, Buffon's translation.

#### METHODE

12

#### PROBLEME I.

Etant donnée la Relation des Quantités Fluentes, trouver la Relation de leurs Fluxions.

#### SOLUTION.

I. DISPOSEZ l'Equation par laquelle la Relation donnée est exprimée fuivant les Dimensions de l'une de ses Quantités Flueutes x par exemple, & multipliez ses Termes par une Progression Arithmetique quelconque, & ensuite par  $\frac{x}{x}$  faites cette Opération séparément pour chacune des Quantités Fluentes; après quoi égalez à zero la somme de tous les produits, & vous aurez l'Equation cherchée.

II. EXEMPLE I. Si la Relation des Quantités Fluentes x & yeff  $k^3 - ax^3 + axy - y^3 = 0$ , difposez d'abord les Termes suivant x, & ensuite suivant y, & multipliez-les comme vous voyez.

Multipliez  $x^3 - ax^2 + axy - y^3 - y^3 + axy - ax^2$ par  $\frac{3x}{x} + \frac{2x}{x} + \frac{2x}{x} + \frac{2x}{x} - \frac{3y}{y} + \frac{3y$ 

 Euler, Gauß and others apeal to computation a lot!

- Mathematics and computing largely progress together till the XIX-th century.

## Computing <u>without</u> Computers!— $\pi$

1 Rhind papyrus 2 Archimedes 3 Vitruvius 4 Chang Hong 5 Ptolemy 6 Wang Fan 7 Liu Hui 8 Tsu Ch'ung Chi 9 Aryabhata 10 Brahmagupta 11 Al-Khwarizmi 12 Fibonacci 13 Madhava 14 Al-Kashi 15 Otho 16 Viète 17 Romanus 19 Van Ceulen 20 Newton 21 Sharp 22 Seki Kowa 24 Machin 25 De Lagny 26 Takebe 27 Matsunaga 28 von Vega 29 Rutherford	2000 BC 250 BC 20 BC 130 150 250 263 480 499 640 800 1220 1400 1430 1573 1593 1593 1593 1593 1593 1593 1593 159	3.16045 (= 4(8/9)2) 3.1418 (average of the bounds) 3.125 (= 25/8) 3.1622 (= sqrt10) 3.14166 3.155555 (=142/45) 3.14159 3.141592920 (= 355/113) 3.1416 (=62832/2000) 3.1622 (= sqrt10) 3.1416 3.141818 3.14159265359 3.14159265358979 3.141592653589793 35 D 16 D 71 D 10 D 100 D 127 D, 112 correct 41 D 50 D 140 D, 136 correct 208 D, 152 correct
27 Matsunaga	1739	50 <b>D</b>
20 VULL VEGU	1794	208  D, 150 correct
30 Strassnitzky Dase	1024	200 <b>D</b> , 132 Contect
31 Clausen	1844	200 <b>D</b> 248 <b>D</b>
321 ehmann	1047	240 <b>D</b> 261 <b>D</b>
33 Dutharford	1853	
31 Shanke	1000	$707 \mathbf{D}$ 527 correct
	10/4	

Source: http://www-gap.dcs.st-and.ac.uk/~history/HistTopics/Pi\_chronology.html



#### PROGRESS = Geometry + Arithmetics + Analysis.



## Computing <u>with</u> Computers! — $\pi$

ENIAC, 1949: 1120**D**; 1000 IPS; Supercomputer 2002:  $2 \cdot 10^{12}$  **D**;  $10^{12}$  IPS (Instruction Per Second)





Moore's law

ENIAC 1949 :  $\frac{1120\mathbf{D}}{1000IPS} = 1.12;$  Kanada 2002 :  $\frac{2 \cdot 10^{12} \mathbf{D}}{10^{12} IPS} = 2.00$ 

··· is only **half** of the story.

Computation cost is **superlinear**  $\implies$  **better algorithms** are needed!!

Initially: 
$$\approx \mathcal{O}(n^2)$$
,  $\frac{\pi}{4} = 4 \arctan \frac{1}{5} - \arctan \frac{1}{239}$ .

- Subquadratic multiplication (Karatsuba)
- Fast Fourier transform
- Arithmetic-geometric mean; elliptic functions ...
- Superguadratically convergent algorithms

```
Finally: \approx \mathcal{O}\left(n(\log n)^2\right)
```



An aside: 'miraculous' Bailey-Borwein-Plouffe alg.

$$\pi = \sum_{n=0}^{\infty} \left( \frac{4}{8n+1} - \frac{2}{8n+4} - \frac{1}{8n+5} - \frac{1}{8n+6} \right) \left( \frac{1}{16} \right)^n$$

The forty-trillionth bit of Pi is '0' 10100000111110011111111001101110001 = A0F9FF371D17593E

$$\pi = \int_0^{1/\sqrt{2}} \frac{4\sqrt{2} - 8x^3 - 4\sqrt{2}x^4 - 8x^5}{1 - x^8} \, dx.$$

Experimental maths in the computer age: Found originally by PSLQ algorithm, finding dependencies between high precision evaluations applied to an inspired guess.

Cf. CECM site on *Experimental Mathematics* at Vancouver & Borwein's pages.

A curiosity

$$B := 4 \sum_{k=1}^{500\ 000} \frac{(-1)^{k-1}}{2k-1},$$

B = 3.1415906535897932404626433832695028841972913993'

 $\pi = 3.1415926535897932384626433832795028841971693993'$ 

Yet another case

#### Integer factorization challenge

The problem of decomposing  $15 = 3 \times 5$  is *not* known to be in class *P*olynomial time.

Triggered by Public Key Cryptosystems based on arithmetic structures, a la RSA.



#### (©Richard Brent.)

Probabilistic algorithms start largely with Rabin in 1976: here (almost) all the algorithms are randomized —they make bets...

## Analysis of Algorithms = an indispensable companion!

Some algorithms are more efficient than others. — By how much ? Why?  $\sim$  Optimizations

"Subliminal" in classical math.

— Trial division for factoring and Erastothenes' sieve are costly.

— Newton's algorithm for root finding doubles the number of digits at each stage  $\neq$  fixed-point iteration only adds a fixed amount.



With respect to the time employed  $\cdots$  in turns of the handle:

and  $\frac{\frac{20(12+n)+15(p+4)}{382}}{\frac{20(12+n)+15(p+d+4)}{382}}$ 

for Mult without stepping for Mult with stepping



#### Burks, Goldstine, von Neumann, 1946 (US Army)

"The logical design of an electronic computing instrument"

"We shall show that for a sum of binary words, each of length n, the length of the largest carry sequence is <u>on the average</u> not in excess of  $2 \log n$ ."



 $\sim$  Feller, Knuth: Runs of good luck in coin tossings...

### Next ...

# The Saga of Digital Trees

1. Pioneers

**1950's**: Scientific computing meets information processing  $\rightarrow$  *non-numerical data*, esp. Sorting & Searching.

First algorithms deal with sorting and searching.

Radix-exchange sort (H&I)



Compare-exchange based on successive *bits* of data. place 0's on left, 1's on right; recurse.

### The trie splitting process (Fredkin)



Separate recursively based on successive bits of data.

Journal of the ACM Vol. 6 (April 1959)

#### Radix Exchange—An Internal Sorting Method for Digital Computers\*

PAUL HILDEBRANDY AND HAROLD ISBITZ

System Development Corporation, Santa Monica, California

This note describes a new technique—Radix Exchange. The technique is faster than Inserting by the ratio  $(\log_2 n)/n$ Its speed compares favorably with internal merging and it has the significant advantage of requiring essentially no working area...

Communications of the ACM Vol. 3 (August 1960)

### **Techniques**

### **Trie Memory**<sup>\*</sup>

EDWARD FREDKIN, Bolt Beranek and Newman, Inc., Cambridge, Mass.

#### Nexuses and Nexus Chains

In order to permit more general description of trie memory, it is helpful to substitute, for the system of successive addresses used in connection with the illustrations in Fig. 1, a system of directed connections that we may call nexuses. A brief explication of nexuses and nexus chains will facilitate further discussion of trie memory.

#### Don Knuth (b. 1938)



### What is the number of turns of the handle?

???

At CalTech around 1965, cooperation of Knuth & de Bruijn

In The Art of Computer Programming 1973



#### Page 131 of Knuth's TAOCP, Vol. 3 (1973) — The original derivation

 $\blacksquare$  Decompose  $\Longrightarrow$  Divide & Conquer recurrence :

$$C_n = n + \sum_{k=0}^n \frac{1}{2^n} \binom{n}{k} (C_k + C_{n-k}).$$

♦ Solve binomial recurrence & reorganize.

♡ Asymptotics : cleverly use Gamma function

$$e^{-x} = \frac{1}{2i\pi} \int_{1-i\infty}^{1+i\infty} \Gamma(s) \, x^{-s} ds$$

Miraculous factorizations occur, residues fly all around, and ...

The Big Theorem of P. 131 $\pm$  of Knuth's Vol. 3

 Tries and Radix-exchange sort have expected cost

(path length, bit comparisons)

$$\sim \boxed{n \log_2 n} + n \left( rac{\gamma - 1}{\log 2} - rac{1}{2} + f(n) 
ight)$$

"where f(n) is the rather strange function  $\cdots$ Furthermore

f(n) < 0.000001725

thus we may safely ignore f(n) for practical purposes."

• Size has expectation (with fluctuations!)

$$\sim \boxed{\frac{n}{\log 2}} + n \cdot \widehat{f}(n)$$

$$f(n) = \frac{1}{\log 2} \sum_{k \neq 0} \Gamma\left(-1 - \frac{2ik\pi}{\log 2}\right) \exp\left(2ik\pi \log_2 n\right)$$



A complicated math exercise. An isolated problem.

An expected outcome ( $\pm$ ):  $O(n \log n)$  by easy probabilistic argument.

 $\diamond$  A useless answer with  $10^{-6}$  fluctuations!

With Moore's law, anyhow, etc.

# The Saga of Digital Trees

2. Analysis

Some "modern" views: Trabb Pardo 1978, Greene 1980, F.–Régnier–Sedgewick–Sotteau 1985, F.–Gourdon-Dumas 1995.

## Methodological advances

Symbolic methods: Combinatorics is reflected by algebra of generating functions

Mainstream methods of enumerative combinatorics ( $\geq 1980$ ) replace recurrences.

$$\left\{f_n\right\} \longrightarrow f(z) := \sum_n f_n z^n.$$

 $\rightarrow$  Difference equations for expected trie costs:

$$\phi(z) = 2e^{z/2}\phi\left(rac{z}{2}
ight) + \operatorname{toll}(z).$$

Semiclassical: Iteration, coefficient extraction,...

## **Methodological advances**

 $\bigotimes \underline{\mathsf{Mellin transforms}}$ 

$$f \xrightarrow{\mathcal{M}} f^{\star} := \int_0^\infty f(x) x^s \frac{dx}{x}$$



Real asymptotics from *complex* singularities. Factorizes linear superposition of models





#### Work from 1965++ yields a systematic approach

### Algorithms

 $\Downarrow$ 

#### Algebra of Costs Gen. Functions

#### $\downarrow$

#### Asymptotic estimates from singularities

applicable to a major combinatorial process of computer science.

#### Knuth's and others' results inform us on *shape* of certain trees: Binary trie (uniform bits)



#### Continued fraction trie



Weyl tree by Devroye,



versus 'beta tree'



# The Saga of Digital Trees

## 3. Data Bases

### Adaptive hashing schemes

Tries are very versatile.

— They can be paginated (bucketted): stop splitting at b.

— They can be combined with *hashing* to cope with non-uniformity of data.

Near 1977-78, several groups discover the virtues of dynamic hashing. Idea: *Split buckets instead of chaining them*. (Larson;

Fagin-Nievergelt-Pippenger-Strong; Litwin)

Expected size of *b*-tree is  $\frac{n}{b \log 2}$  + fluctu, corresponding to 69% filling ratio.

Compare with similar ratio for B-trees (Yao)

2 accesses suffice for very large DB.

Extendible Hashing transforms the index into a perfect tree  $\equiv$  array that can be paginated.

*H*=height Index size  $\equiv 2^H$ 

(Yao, Régnier, F., ca 1980)

 $\mathbb{E}(H) \sim \left(1 + \frac{1}{b}\right) \log_2 n; \qquad \mathbb{E}(2^H) \approx 4^{\text{fluctu}} n^{1+1/b}.$ 



Height: One of the very first intrusions of saddle point method in Analysis of Algorithms.

$$[z^n]f(z) = \frac{1}{2i\pi} \oint f(z) \frac{dz}{z^{n+1}}.$$



→ Jacquet & Szpankowski's "analytic de-Poissonization": analyse under probabilistic model with "imaginary probabilities"!





## Skip lists

From VSAM's to skip lists

Idea 1 (old): build indexes of indexes of indexes ... Idea 2: balance  $\rightarrow$  B-trees Idea 2': randomize!  $\equiv$  Pugh's skip lists



Much easier to maintain than balanced structures!

Analysis by Papadakis + Munro, Poblete Kirschenhofer, Martínez, Prodinger entirely based on trie technology.

## Probabilistic counting algorithms

Can you estimate to 5% the number of different words in Shakespeare given a pencil and one sheet of paper?

Yes. F.+Martin (1985) for data base query optimization.

Ideas: hash to get uniformity; observe bit patterns.  $0 \dots = 50\%$  of times;  $10 \dots = 25\%$ ;  $110 \dots = 12.5\%$ Try  $2^{K}$  where  $\overbrace{11 \dots 1}^{K} 0 \dots$  is longest initial run of 1's. The best observable known is trie-like and has accuracy



for m words of memory

(+"stochastic averaging"); 0.78 is a Mellin constant.

Works in distributed environment: Yellow pages of New York U San Francisco by phone line!

Data mining applications. Quick running counts in routers (Durand 2003) based on other trie observables.

# The Saga of Digital Trees

3. Protocols

♡ 1970: the shared communication channel



**Ethernet:** Try; wait  $\times 1$ ,  $\times 2$ ,  $\times 4$ , etc

 $\sim$  Aldous 1987: Ethernet is unstable!

V 1977: The **Tree/Stack protocol**CTM = Capetanakis, Tsybakov, Mikhailov



= A digital trie but with a flow of arrivals!

Erroneous analyses missed the wobbles.
 Variance by Kirschenhofer, Prodinger et al.= Mellin
 + modular forms.

**Tree protocol**  $\implies$  Poisson GenFun solves (p+q=1)

 $\psi(z) - \psi(\lambda + pz) - \psi(\lambda + qz) = \text{toll}(z).$ 

A non-commutative iteration semigroup with a globally invariant measure.

**Theorem.** Stable till  $\lambda_{max} = 0.36017$  root of:

$$-\frac{1}{2} = \frac{e^{-2x}}{1-2x} \sum_{i\geq 0} 2^{i}g\left(\frac{x}{2^{i}}\right),$$
$$g(y) := e^{-2y} \left( \left(e^{-y}(1-y) - 1 + 2y(1+y)\right) \right).$$

Analyses by Fayolle, F., Hofri, Jacquet, Mathys  $\implies$ 

Ternary tree algorithms gives 10% better throughput
 Protocol is hyperstable at all arrival rates.

The IEEE 802.14 norm...a failed success story!

Also analyses by Greenberg+F+Ladner: tree protocol modified to attain 93% of optimal:  $\lambda_{max} = 0.4672$ .

#### Leader Election:

#### (i) The leftmost branch of a trie



### (ii) The leftmost border of a trie



Analyses by Fill, Mahmoud, Szpankowski, Prodinger, F+Sedgewick; includes distributions.

 $(2 \log n)$  rounds;  $\log_2 n$  rounds.

# The Saga of Digital Trees

4. Text and compression

## Tries meet texts again!

Szpankowski's *Analysis of Algorithms on Sequences*.





Random text: kwnbpr hwnqqcpq yt nxgfhsd agghos fhskla zmmxnz kasiweyzkcn ejhjsal ehrdjn...

*≠* "Natural" language text: Cale Pismo przez Boga
jest natchnione i pozyteczne do nauki, do wykrywania
bledow...

Can be compressed!

 Lempel & Ziv invent LZ compression (1977+) based on building adaptive dictionaries.
 a|b|r|ac|ad|ab|ra|abr|aca|d|abra|abrac|ada|br|aa|br|acad|abraa|... Turns out to be related to digital search trees.

- Régnier-Jacquet (1987) do distributional analysis of tries under Bernoulli models.
- Szpankowski-Jacquet (1990) do average-case analysis of tries under *Markovian dependencies*.
- Jacquet–Szpankowski–Louchard (1995+) extend distributional analysis to DST's:

 $\frac{\partial}{\partial z}F(z,u) = F(z,pu)F(z,qu) + \text{fudge}$ 

#### $\sim$ Combines everything:

algebra of trie costs, Mellin, analytic dePoissonization...

©© Complete characterizations of *Lempel-Ziv* algorithms, notably: redundancy.

# The Saga of Digital Trees

## 5. Geometry & Dynamical Systems





- "Thermodynamic formalism" by Ruelle (1970)
- Operators & Euclid's alg. by Babenko, D. Mayer (1977.)
- Related to information theory & tries by Vallée (1995+)
- T is a transformation. Iterates?

Transfer operator:  $\mathcal{G}_{s}[f](x) := \sum_{h \in T^{-1}} (h'(x))^{s} f \circ h(x).$ 

Vallée: Spectra & functional analysis serve to generate probabilities of prefixes  $\rightsquigarrow$  tries.

Tries under dynamic source models;  $\rightarrow$  Bentley-Sedgewick's Ternary Search Tries Entropy for size, depth path-length; Eigenvalue  $\lambda(2)$  for height, etc.



Applies to <u>continued fraction</u> representations & algs: → Накмем Algorithm (Gosper, 1972); 2D orientation = Avnaim, Boissonnat, Devillers, Preparata, Yvinec 1997.



$$K_0 n \log n + K_1 n + Q(n) + K_2 + o(1),$$

$$K_{0} = \frac{6 \log 2}{\pi^{2}}, \quad K_{1} = 18 \frac{\gamma \log 2}{\pi^{2}} + 9 \frac{(\log 2)^{2}}{\pi^{2}} - 72 \frac{\log 2\zeta'(2)}{\pi^{4}} - \frac{1}{2}.$$
  
$$\heartsuit \heartsuit \heartsuit \quad Q \text{ depends on Riemann hypothesis!!!}$$

# The Saga of Digital Trees

6. Everywhere...

## Random Trie Encounters

Vol 2., F+Gourdon+Panario

© Quadtries and geometry, multiD search Rivest–Bentley–Samet

© Other probabilistic counting algorithms Morris–Freivalds, Wegner's, etc

Sinary Decision Diagrams (BDD's) by Bryant (!?)
 = Fully developed tries + common subtree
 factoring...

♡ Hierarchical data compression by J. Kieffer

 $\heartsuit$  Level compressed tries  $\implies$  fast lookup in routers! Nilsson et al.

### Finally ...

## Where are we?

Analysis of algorithms as of now:

#### **Complex Models**

... become more and more tractable.

♡ A large number of basic algorithms have been analysed. Cf Sedgewick's book.

 Symbolic Methods help translate complex probabilistic models into gen. functions.
 Analytic Combinatorics = an extensive calculus of asymptotic properties based on singularities.
 A unified theory of basic random combinatorial structures and algorithms.

Fruitful connections with computer algebra.
 Automatic counting, automatic asymptotics, automatic random generation.

#### Two basic principles $\mapsto$ "dictionaries"

#### **SYMBOLIC METHODS**

Generating functions



 $z + z^{2} + z^{3} + 2 z^{4} + 2 z^{5} + 4 z^{6} + 5 z^{7} + 9 z^{8} + \cdots$ 

 $f(z) = z + f(z^2 + z^3 + z^4)$ 

#### **ANALYTIC FUNCTIONS AND SINGULARITIES**



#### The example of TRAINS

#### • Cope with complex structural "specifications"



0.1008557594 (0.5180547070)



### Analytic Combinatorics

= organize *random discrete structures* (cf. stochastic proc.)

= tightly coupled with Analysis of algs.

- Permutations: order stat., search & sort.
- Words: patterns, comput. biology, coding
- DIGITAL TREES
- Allocations: hashing, comb. opt.,...
- Graphs: combinat opt., networks (?)
- Trees: symbolic manipulation, etc.

THERE IS A story about two friends, who were classmates in high school, talking about their jobs. One of them became a statistician ... "And what is this symbol here?" "Oh," said the statistician, "this is pi." "What is that?" "The ratio of the circumference of the circle to its diameter." "Well, now you are pushing your joke too far," said the classmate, "surely the population has nothing to do with the circumference of the circle."

The miracle of the appropriateness of the language of mathematics for the formulation of the laws of physics is a wonderful gift which we neither understand nor deserve. We should be grateful for it and hope that it will remain valid in future research and that it will extend, for better or for worse, to our pleasure, even though perhaps also to our bafflement, to wide branches of learning.

- Eugene Wigner

<sup>&</sup>quot;The Unreasonable Effectiveness of Mathematics in the Natural Sciences," in Communications in Pure and Applied Mathematics, vol. 13, No. I (February 1960).