

Avventure Algoritmiche

Bioinformatica e Sequenziamento di Genomi - I

Alberto Policriti



Udine, 30.10.2019

L'inizio: Crick and Watson, 1953

MOLECULAR STRUCTURE OF NUCLEIC ACIDS

A Structure for deoxyribose Nucleic Acid

WE wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has several features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey¹. They briefly sketch a structural possibility for an advance of ribonucleases. Their model consists of three close spaced chains, with the phosphate near the three gaps, and the bases on the outside. In our opinion, this structure is unsatisfactory for two reasons. (1) We believe that the material which gives the X-ray diagram is the salt, not the free acid. With this the outer hydrogen atoms in a nucleic acid chain would hold the structure together, especially as the negatively charged phosphate near the sea will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another three-chain structure has also been suggested by Bower in the papers. It is a model the phosphate are on the outside and the base on the inside, linked together by hydrogen bonds. This structure is described as a rather 2-chained, and for this reason we shall not comment on it.

We wish to put forward a relatively objective suggestion for the salt of deoxyribose nucleic acid. This structure has two linked chains each coiled round the same axis (see diagram). We have made the usual chemical assumption, namely, that each chain consists of phosphate diester groups joining deoxyribose molecules with 3',5' linkages. The two chains then coil round the same axis, but in a dextral sense to the other.

Both chains follow right-handed helices. The coiling of the two chains is in opposite directions. Each chain closely resembles Pauling's model No. 1; that is, the bases are on the inside of the helix and the phosphates on the outside. The non-adjacency of the sugar and the sugar phosphate is close to Pauling's standard configuration, the sugar being roughly perpendicular to the phosphate. There



Fig. 1. The structure of the salt of deoxyribose nucleic acid. The two chains are coiled round the same axis, but in opposite directions. Each chain closely resembles Pauling's model No. 1; that is, the bases are on the inside of the helix and the phosphates on the outside.

GENETICAL IMPLICATIONS OF THE STRUCTURE OF DEOXYRIBONUCLEIC ACID

By J. D. WATSON and F. H. C. CRICK

Medical Research Council Unit for the Study of the Molecular Structure of Biological Systems, Cavendish Laboratory, Cambridge

Francis Crick
Leslie Watson

THE importance of deoxyribonucleic acid (DNA) within living cells is undisputed. It is found in all dividing cells, largely if not entirely in the nucleus, where it is an essential constituent of the chromosomes. Many lines of evidence indicate that it is the carrier of a part (if not all) the genetic specificity of the chromosomes and thus of the gene itself.

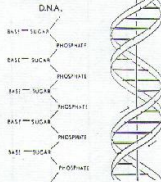
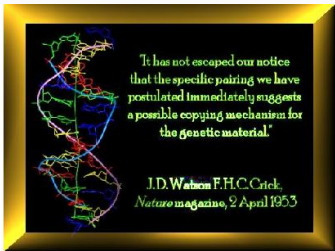


Fig. 1. Chemical formula of a right-handed deoxyribose nucleic acid.



Fig. 2. The space is nearly symmetrical. The two chains are coiled round the same axis, but in opposite directions. Each chain closely resembles Pauling's model No. 1; that is, the bases are on the inside of the helix and the phosphates on the outside.

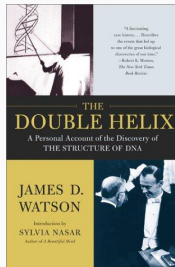
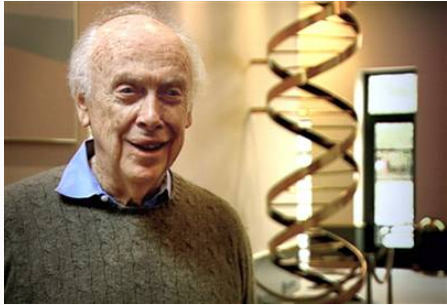
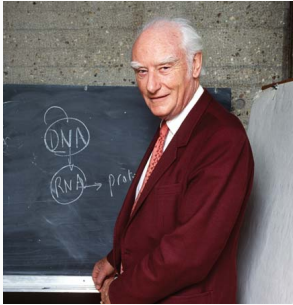


"It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material."

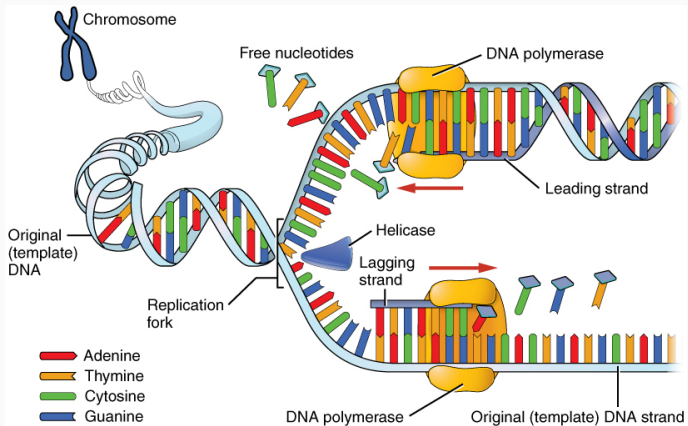
J.D. Watson F.H.C. Crick
Nature magazine, 2 April 1953



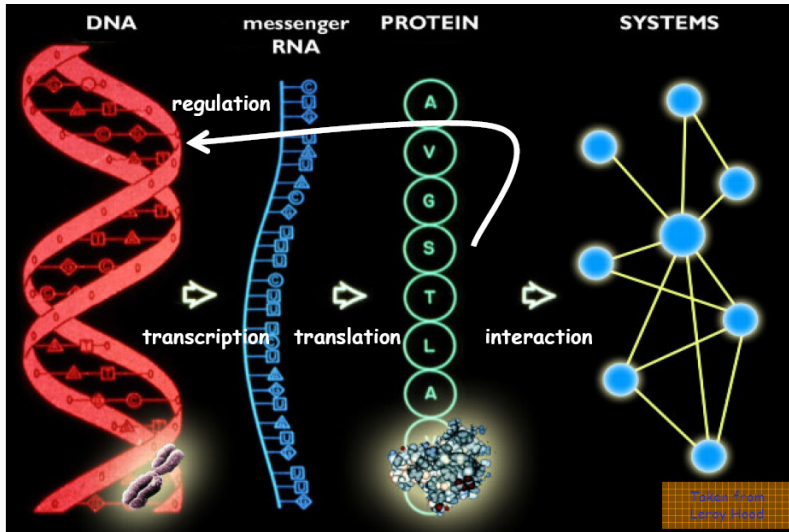
L'inizio: Crick and Watson, 1953



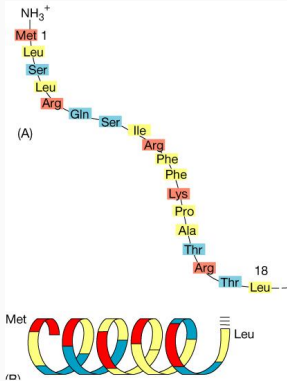
Informazione: codifica, memorizzazione, duplicazione, ...



Informazione: codifica, memorizzazione, duplicazione, ...

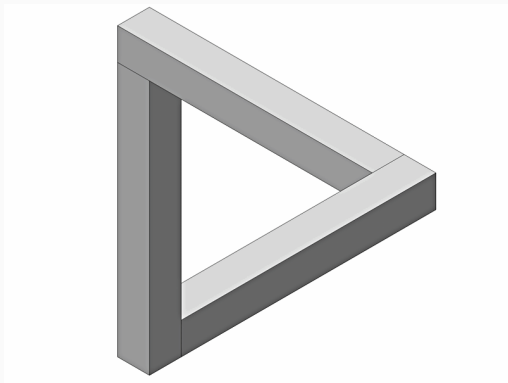


Informazione: codifica, memorizzazione, duplicazione, ...

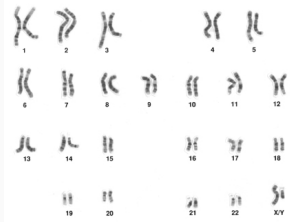
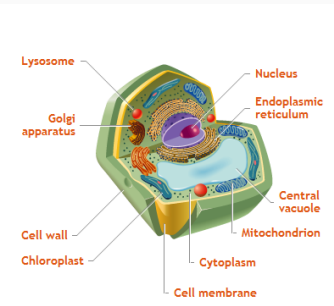
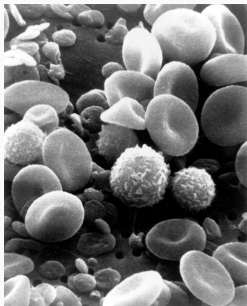


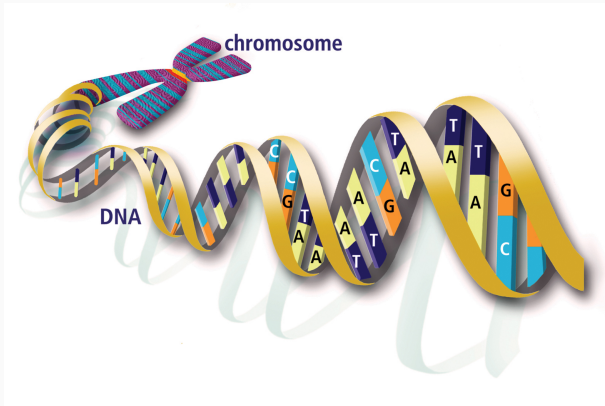
La filosofia è scritta in questo grandissimo libro che continuamente ci sta aperto innanzi a gli occhi (io dico l'universo), ma non si può intendere se prima non s'impara a intender la lingua, e conoscer i caratteri, ne' quali è scritto. Egli è scritto in lingua matematica, e i caratteri son triangoli [A,T], cerchi [C,G], ed altre figure geometriche [Met, Lys, Leu, ...], senza i quali mezi è impossibile a intenderne umanamente parola; senza questi è un aggirarsi vanamente per un oscuro laberinto.

GALILEO GALILEI

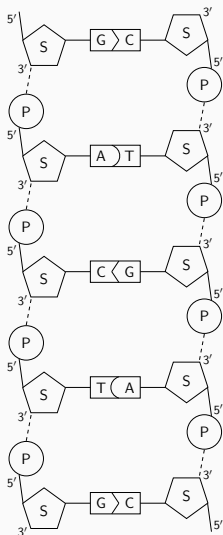


Biologia. Sequenziamento: Dati (per un Biologo)



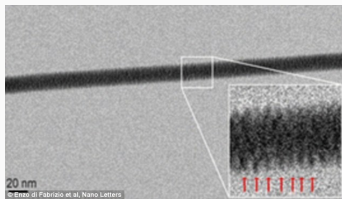


The raw material



- una voluta di DNA (3.4nm) 10.5 basi
- una cellula 2m of DNA
- un organismo umano 10^{13} cellule
- enorme quantità di DNA

Estremamente difficile da manipolare!



II puzzle



1 Genomic DNA



DNA is extracted from the cells of the organism of interest.

2 Library Creation



3 Assembly

```
...ACCGTAAATGGGCTGATCATGCTTAA  
          TGATCATGCTTAAACCCGTGTCATCCTACTG...  
...ACCGTAAATGGGCTGATCATGCTTAAACCCGTGTCATCCTACTG...  
...ACCACCGTAAATGGGC...          ...GCATCCTACTGTACGTTAA...
```

4 Annotation and Analysis

Sequenziamento

- estrazione e frammentazione del DNA
- *sequenziamento* dei frammenti e generazione delle *reads* (letture)

attività del **wet-lab**

Assemblaggio

- *memorizzazione* delle reads (+ altre informazioni)
- *assemblaggio* di reads (corte) e produzione di *contigs* (lunghi)

attività del **dry-lab**

Sequenziamento

- estrazione e frammentazione del DNA
- *sequenziamento* dei frammenti e generazione delle **reads** (letture)

attività del **wet-lab**

Assemblaggio

- *memorizzazione* delle reads (+ altre informazioni)
- *assemblaggio* di reads (corte) e produzione di **contigs** (lunghi)

attività del **dry-lab**

Vincoli

- *Sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads

Input (per l'assemblaggio di un genoma)

GAC TTTGTAACATACAACCTTTAATCACGCTCAATATGACATTATCTGACGCTGATCTGTCTC...

GAC TTTGTAACATACAACCTTTAATCACGCTCAATATGACATTATCTGACGCTGATCTGTCTC...

Input (per l'assemblaggio di un genoma)

GAC TTTGTA ACATACAACCTTTAATCACGCTCAATATGACATTATCTGACGCTGATCTGTCTC...

CTTTAATCAG

ACATACAAC

CTCAATATGACAT

TATCTGACGCTGA

TCTGTCTC

GAC TTTGTA

...

Input (per l'assemblaggio di un genoma)

GAC TTTGTAACATACAACCTTTAATCACGCTCAATATGACATTATCTGACGCTGATCTGTCTC...

CTTTAATCAGG
ACATACAAC
CTCAATATGACAT
TATCTGACGCTGA
TCTGTCTC
GAC TTTGTA

...

?

Input (per l'assemblaggio di un genoma)

GACTTTGTAACATACAACCTTTAATCACGCTCAATATGACATTATCTGACGCTGATCTGTCTC...

Soluzione: più copie (indipendenti/ben-distribuite)

Input (per l'assemblaggio di un genoma)

GACTTTGTAACATACAACCTTTAATCACGCTCAATATGACATTATCTGACGCTGATCTGTCTC...

Soluzione: più copie (indipendenti/ben-distribuite)

GACTTTGTA

Soluzione: più copie (indipendenti/ben-distribuite)

GACTTTGTA

TTTGTAACATAC

Soluzione: più copie (indipendenti/ben-distribuite)

GACTTTGTAACATAACAAC

TTTGTAACATAC

Soluzione: più copie (indipendenti/ben-distribuite)

GACTTTGTAACATAACAAC

TTTGTAACATAACAACCTTAA

Soluzione: più copie (indipendenti/ben-distribuite)

GACTTTGTAACATACAACCTTTAATCACG

TTTGTAACATACAACCTTTAA

Soluzione: più copie (indipendenti/ben-distribuite)

GACTTTGTAACATAACAACCTTTAATCACG . . .

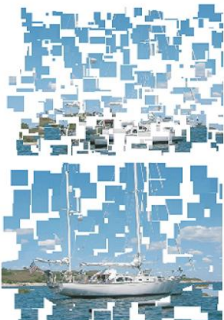
TTTGTAACATAACAACCTTTAATCACGCTCAATA . . .

L'obiettivo finale

TCAGGGCCAGCAGCAAAGTTTCCAGCTGATCATCCTGATGGTGCGCGGTGAGTAAGACACCACCTTCCGG
GAGGTTTTGGCGAAAAACGTCGTAGCGGGCCTGACGGGCAAGATCCTCAAGACTTTGCCGGGCCTTTTTC
TCAATTTCAACACGTTTTGCTATGAATGGTACTGTCAGTACATCGCAAACCTGCGCGCAATGGGTGAGCC'
ATTGGTCAGCCTTAGGATTTAAACCATGATGCACATGCACTGCCTGCAGGGAAAAGCCGTCTCGCTGAGA
TATATAACGGGCCAGCAATTCAAGTAACACGCGGGAATCGAGGCCGCCGCTGAATCCAACAGTAAGTGC
CCGGGCCGGATCAGCTGATCCAGCACAGAAAAGACACTTGCTTCCAATTCATGATGACTCACAACAGATA
CTTCATCGATAAAGTACAATATGCGGATATAGCAGATCCTGTTTGTATTACGTGAACCGGGATAATAA
GCCATATTGCTTGCGGCTCGTTTACACGATAACGATGATCGGATGGAAGAAAATAAAAAAAGCCCTGAGC
GGGGCTTAAAGAAAAGGAAAGCATTAACTTAACGGGTTTCCAGCGGAGCGAAGCTCTTCACAAGATCA
TCAATGGCTTTCATCTGTTGACGGAATGGCTCCAGCTTATCCAGTGGAATGCACTAGGACCATCACAAC
GTGCATGTTCCGGATCAGGATGCGACTCAATAAACAGCCAGCCAGACCAACCGCATTACAGCGGTGC
CAGTTCAGTCACTTGCTCAGCAGCACCCTTGAAGCCGACCCAGCGGTACGACACTGCAGGGAATGC
GTTACGTCAAAAATCACCGGGCTGCCCTGGCTGGCATCTTTCATCACGCCAAGCCAGCATGTCGACCA
CCAGATTGTCATAACCGAAATTCACACCGGCTTACACAGAATAATACGGTCGTTACCGCATTGCGCGAA
TTTCTCAACGATGTTCTTAACTGTCCCGGGCTCAGGAACTGCGGTTTTTTCACGTTAATCACGTTACCT
GTTTTAGCCAGAGCTTCCACCAGATCAGTCTGACGCGCCAAGAAGGCCGGTAATTGCAGCACATCCACG
CTTACCAGATCGGTTTGGCTTGCCATGTTTCATGCACATCTGTGATCAGGCTGACACCGAAGTTTTTCTT
CAGTCTTCAAAAATACGCAGACCTTCTTCCATACCCGGACCACGGTAGGAATGAACTGAGGAGCGGTTG
GCTTTATCCCAGCTGGCTTTAAATACCAAAGGGATGCCAGTTTTTTCAGTTACGGTGACGTAGGTTTCAC
AGATCGACATCGCCAGATCCCGTGACTCCAGTACATTCATGCCACCAAACAGTACGAATGGCTGATCATT
AGCAACATTAATCCGTTAAACTGAACGACTTCTGTTTCATTATCACTCCATCAGTGAAAAATATGG
CGTCATGATCCAGCATCGCCAGCTGTACTTTCAGTACGGCTGCAACCGGATCCTGAGGACATTGTTCAAT
AAAATAGGTGAATCATGAGCTGCCAGCTGCGGACATTCAGCTGTTCATATACCAGCCCGGATCACGG
ATTTTCATAGGGGTCATCGGGGTTTCATCGCCAGCAAACCTCGCTGCAGCGTAACGCCTCCGGCAAACAAC
GACTTTGTAACATAACAACCTTAAATCACGCTCAATAAACGTGACATATCTGACGCTGATCTGTCTCTTT
TGTATATTTACTATCCATCCGGGTTAAGTCACCGAGTGTGGCTCGCAGCATTAATGACTGCTGCTCATAA
TCCCCTCTTACCAGTAAACCGGATCGATTAATAACAGGCTTACTTCCGGGAAAAGCAGCAAGAAGTTAC
CAGGGAAGCAAACACCACCGGCAATTCATCGAATGAGCCAGATGCATTAACACAGTACCCAGCTG
GAGCGGCAAACCCAGTGCCTGCATCAGCACCTGATCCAGCAGGCAGTTTTCCGGTGCATAGTAAATCTGC
CAGTCACCGGAAAAATTGCAATTCATGATAGAACGCATGCAACAGTTGTTTACGACGCCCTTATCATCAG

High “quality” assemblies

- alte “coperture” (molti *genomi equivalenti*)
- molte reads
- reads lunghe
- buon sw (*assembler*)
- buon hw (parallelo?)
- ...



Il consorzio Italo-Francese per il sequenziamento del Genoma Nucleare della Vite



doi:10.1038/nature06148

nature

LETTERS

The grapevine genome sequence suggests ancestral hexaploidization in major angiosperm phyla

The French–Italian Public Consortium for Grapevine Genome Characterization*

The analysis of the first plant genomes provided unexpected evidence for genome duplication events in species that had previously been considered as true diploids on the basis of their genetics^{1–3}. These polyploidization events may have had important consequences in plant evolution, in particular for species radiation and adaptation and for the modulation of functional capacities^{4–6}. Here we report a high-quality draft of the genome sequence of grapevine (*Vitis vinifera*) obtained from a highly homozygous genotype. The

All grapevine varieties are highly heterozygous; preliminary data showed that there was as much as 13% sequence divergence between alleles, which would hinder reliable contig assembly when a whole-genome shotgun strategy was used for sequencing. Our consortium therefore selected the grapevine PN40024 genotype for sequencing. This line, originally derived from Pinot Noir, has been bred close to full homozygosity (estimated at about 93%) by successive selfings, permitting a high-quality whole-genome shotgun assembly.

Tecnologia: la prima “sequencing revolution”: 2005-2013



23.648 ABI3730

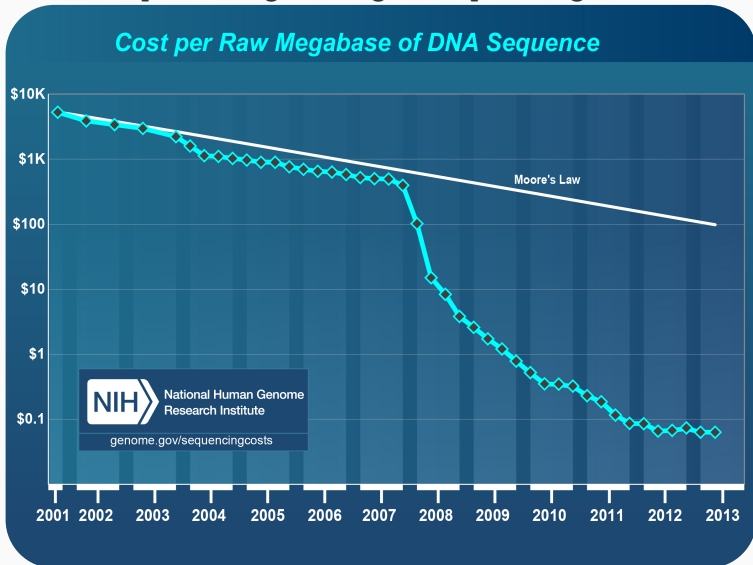
vs



1 HiSeq2000

Tecnologia: la prima "sequencing revolution": 2005-2013

<http://www.genome.gov/sequencingcosts/>



Tecnologia: la seconda rivoluzione: ora!



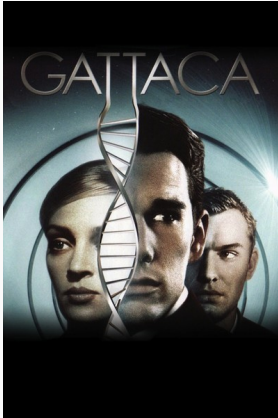
Single Molecule, Real Time
(SMRT) technology



Nano sensing technology



Science fiction?

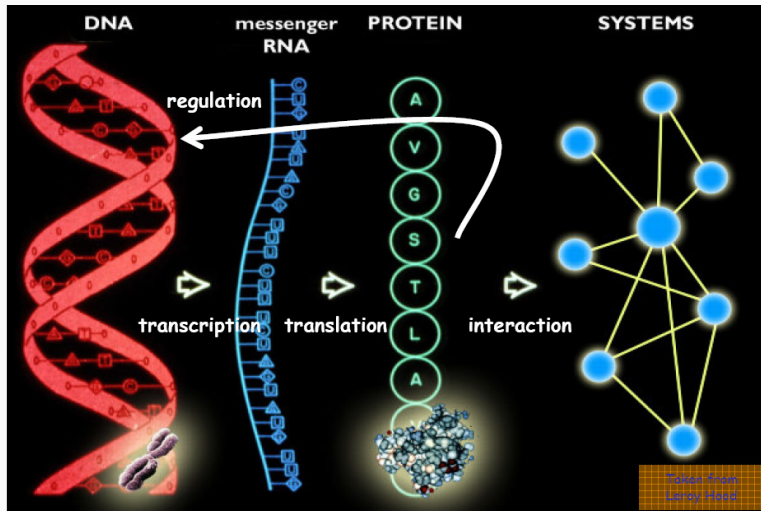


Ethan Hawke as Vincent/Jerome

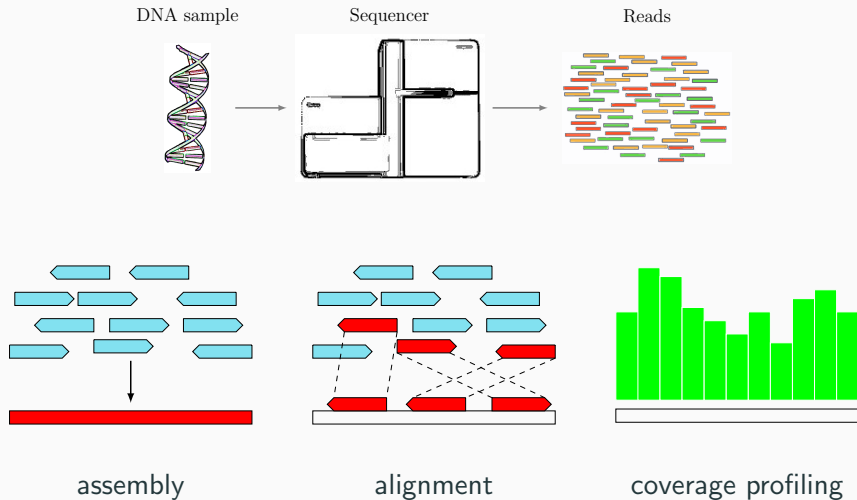
Uma Thurman as Irene

Jude Law as Jerome/Eugene

Written and Directed by *Andrew Niccol*, **1997**



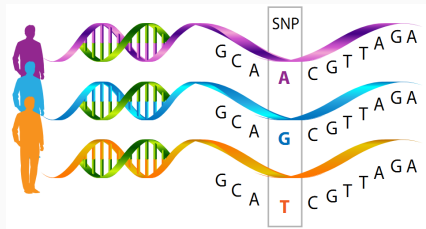
Sequenziamento ed Allineamento



Parliamo di un genoma ma ... siamo diversi!

Due organismi differiscono per un (alto) numero di

Single Nucleotide Polymorphisms (SNP's)



Ri-sequenziamento

Allineamento delle reads del genoma di un individuo, contro una
sequenza di riferimento *R*

per determinare gli SNP's (l' "impronta")

Approssimativamente 12M SNP's nel genoma umano (3.2G basi): ≈ 1 ogni 300 basi.

Quanti genomi?

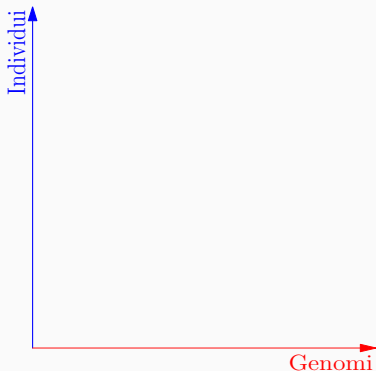
Ogni specie ha il suo genoma.



Quanti genomi?

Ogni specie ha il suo genoma.

Ogni individuo di una data specie ha il suo genoma.



Controllare la complessità: Algoritmi Paralleli e Sequenziali



Sequenziale vs. Parallelo

- **Input** $\langle a_1, \dots, a_n \rangle \implies$ **Output** $\sum_{i=1}^n \sqrt{a_i}$
- **Input** un genoma $R \implies$ **Output** Tutte le α (di lunghezza k) tali che $\alpha\alpha \in R$
- **Input** a set of reads $\mathcal{R} \implies$ **Output** the reference sequence R

Osservazione

Possiamo (dobbiamo) *strutturare* R mentre la carichiamo

Vantaggi? Esempio

Ordinare un insieme di 1G ($\approx 2^{30}$) di numeri, ci consente di cercare *qualunque* elemento in non più di 30 passi

Osservazione

Possiamo (dobbiamo) *strutturare* R mentre la carichiamo

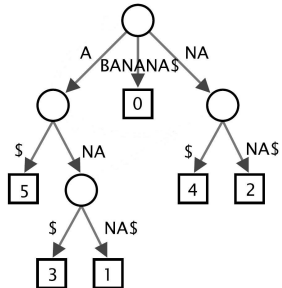
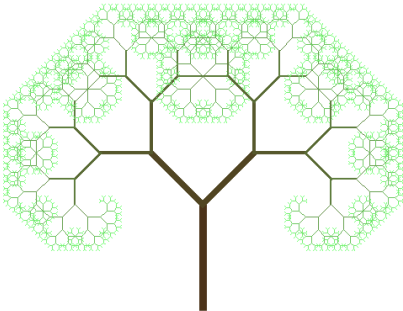
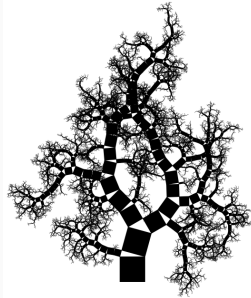
Vantaggi? Esempio

Ordinare un insieme di 1G ($\approx 2^{30}$) di numeri, ci consente di cercare *qualunque* elemento in non più di 30 passi

Domanda

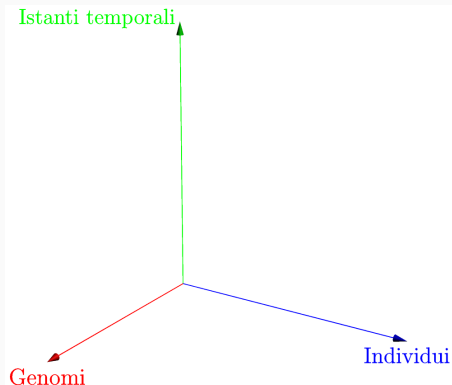
Possiamo fare qualcos di simile (i.e. *ordinare*) un *testo* (come R)?

Controllare la complessità: Strutture Dati



Il tempo

Dobbiamo aggiungere una dimensione: **il tempo**



Monitoring chronic lymphocytic leukemia progression by whole genome sequencing reveals heterogeneous clonal evolution patterns

Anna Schuh,¹ Jennifer Becq,² Sean Humphray,² Adrian Alexa,² Adam Burns,¹ Ruth Clifford,¹ Stephan M. Feller,³ Russell Grocock,² Shirley Henderson,¹ Irina Khrebtukova,⁴ Zoya Kingsbury,² Shujun Luo,⁴ David McBride,² Lisa Murray,² Toshi Menju,^{3,5} Adele Timbs,¹ Mark Ross,² Jenny Taylor,¹ and David Bentley²

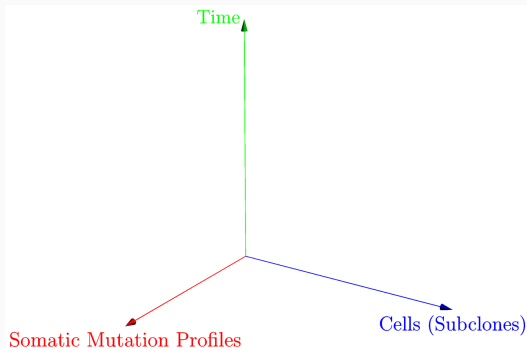
¹Oxford National Institute of Health Research (NIHR) Biomedical Research Centre, University of Oxford, Oxford, United Kingdom; ²Illumina Cambridge Ltd, Saffron Walden, United Kingdom; ³Biologic Systems Architecture Group, Department of Oncology, Weatherall Institute of Molecular Medicine, University of Oxford, Oxford, United Kingdom; ⁴Illumina Inc, Hayward, CA; and ⁵Department of Thoracic Surgery, Graduate School of Medicine, Kyoto University, Kyoto, Japan

Chronic lymphocytic leukemia is characterized by relapse after treatment and chemotherapy resistance. Similarly, in other malignancies leukemia cells accumulate mutations during growth, forming heterogeneous cell populations that are subject to Darwinian selection and may respond differentially to treatment. There is therefore a clinical need to monitor changes in the subclonal composition of cancers during disease progression. Here,

we use whole-genome sequencing to track subclonal heterogeneity in 3 chronic lymphocytic leukemia patients subjected to repeated cycles of therapy. We reveal different somatic mutation profiles in each patient and use these to establish probable hierarchical patterns of subclonal evolution, to identify subclones that decline or expand over time, and to detect founder mutations. We show that clonal evolution patterns are heterogeneous in

individual patients. We conclude that genome sequencing is a powerful and sensitive approach to monitor disease progression repeatedly at the molecular level. If applied to future clinical trials, this approach might eventually influence treatment strategies as a tool to individualize and direct cancer treatment. (*Blood*. 2012; 120(20):4191-4196)

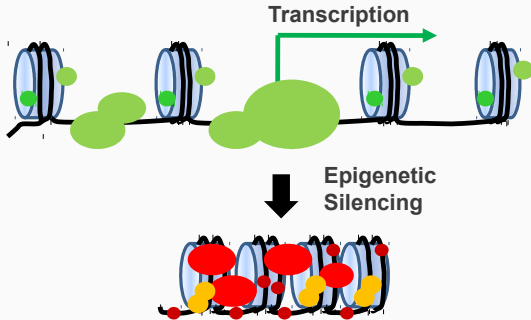
... We defined cellular subpopulations on the basis of somatic mutation profiles ...



Epigenetica (e codifica): la “truffa”

EPIGENETICA (letteralmente, “sopra la genetica”)

Ogni, potenzialmente stabile ed ereditabile, cambiamento nella espressione genica che occorre senza un cambiamento nella sequenza di DNA



Definizioni

- Meccanismi molecolari che convertono l'informazione genetica in fenotipi e caratteri osservabili (C. Waddington 1940)

Molecular mechanisms converting genetic information in phenotypes and observable characters

- Modificazioni della funzione genica **reversibili** ed **ereditabili**, che hanno luogo senza cambiare la sequenza di DNA

- **Reversible** and **inheritable** modifications of the genic function, taking place without changing the DNA sequence

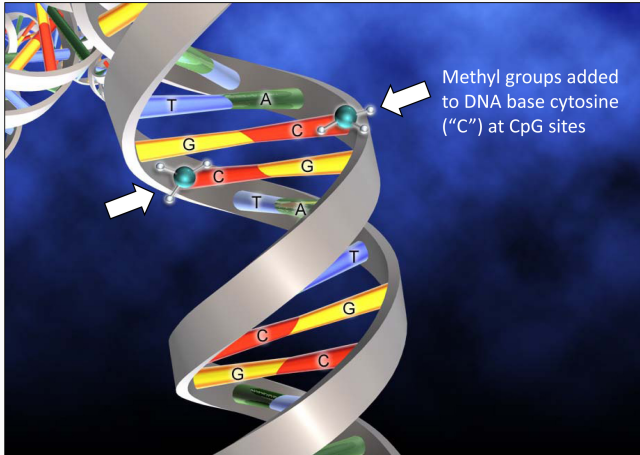
Definizioni

- Meccanismi molecolari che convertono l'informazione genetica in fenotipi e caratteri osservabili (C. Waddington 1940)
- Modificazioni della funzione genica **reversibili** ed **ereditabili**, che hanno luogo senza cambiare la sequenza di DNA

Definizioni

- **Epigenetica** modifiche che alterano l'espressione di geni senza cambiare la loro sequenza
- Gli obiettivi di tali modifiche sono essenzialmente DNA ed **istoni**
- La modifica epigenetica più studiata è la **metitazione** del DNA

Epigenetics



Due (principali) problemi computazionali

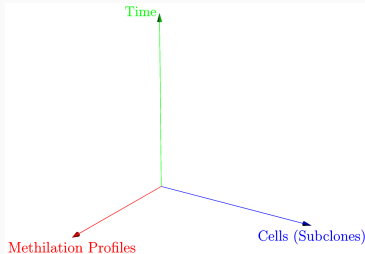
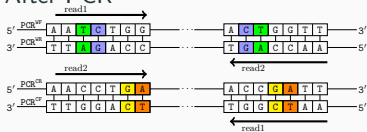
Original DNA fragment



Bisulfite Treated DNA fragment



After PCR



Campo molto stimolante

- Alto impatto sulle Scienze della Vita
- Tecnicamente sofisticato
- Bei problemi computazionali

Sfide

- Algoritmi: manipolazione parallela di sequenze nucleotidiche, compressione, codifica, ...
 - Clustering di *popolazioni* di sequenze
 - Strumenti di analisi per la determinazione di impronte genetiche in serie *temporali*
 - Strumenti di simulazione
-
- Conoscenza *inter-disciplinare*
 - Dati *standard* (abbastanza)
 - Abbiamo (ancora) *tempo*

È probabilmente vero in linea di massima che nella storia del pensiero umano gli sviluppi piú fruttuosi avvengono frequentemente in quei punti di interferenza fra due diverse linee di pensiero.

“Fisica e Filosofia” W. HEISENBERG