

Bioinformatica: DNA e Algoritmi

Alberto Policriti

Dpt. of Mathematics and Informatics, University of Udine.



DI COSA PARLEREMO

IN GENERALE

- Definiamo i termini: DNA & Algoritmi
- Tecnologie (\Rightarrow Problemi) \Rightarrow *Sequenziamento e Assemblaggio*
- **Bioinformatica**: Passato e Futuro

DI COSA PARLEREMO

IN GENERALE

- Definiamo i termini: DNA & Algoritmi
- Tecnologie (\Rightarrow Problemi) \Rightarrow *Sequenziamento e Assemblaggio*
- **Bioinformatica**: Passato e Futuro

DI COSA PARLEREMO

IN GENERALE

- Definiamo i termini: DNA & Algoritmi
- Tecnologie (\Rightarrow Problemi) \Rightarrow *Sequenziamento e Assemblaggio*
- **Bioinformatica**: Passato e Futuro

DI COSA PARLEREMO

IN GENERALE

- Definiamo i termini: DNA & Algoritmi
- Tecnologie (\Rightarrow Problemi) \Rightarrow *Sequenziamento e Assemblaggio*
- **Bioinformatica**: Passato e Futuro

LE AREE

- Algoritmica su stringhe
 - Matematica del discreto
 - String matching esatto e approssimato
- Systems biology
 - Matematica del continuo
 - Automi e biologia
- Computare usando il DNA

CRICK E 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 (DNA). This structure has most features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey.¹ They have used a structural formula available to us in advance of publication. Their model consists of three colored rods along with the phosphorus near the three axis and the bases at the outside. In our opinion, this structure is unsatisfactory for two reasons: (1) We believe that the proposed model gives the X-ray distances in the salt, not the free acid. Without the ionic hydrogen atoms it is not clear what forces would hold the structure together, especially as the negatively charged phosphates near the axis will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another structural structure has also been suggested by Phares in the press. In its model the phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure is described in a rather detailed, and for this reason we could not comment on it.

We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual chemical assumptions, namely that each chain consists of phosphate groups joining deoxyribose ribonucleosides with 3'-5' linkages. The two chains that run their courses are related by a direct perpendicular to the other axis. Both chains follow right-handed helices, but owing to the double sequence of the bases in the two chains run in opposite directions. Each chain, namely according to Berg's model No. 1, has its bases on the inside of the helix and the phosphate on the outside. The correspondence of the sugar and the phosphate is in phase in the two chains. The sugar being roughly perpendicular to the phosphate base. There

¹ This paper is being published in the Proceedings of the National Academy of Sciences, U.S.A. (1953).

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
James Watson

THE importance of deoxyribonucleic acid (DNA) within living cells is undoubted. 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 of if not all the genetic specificity of the chromosomes and thus of the gene itself.

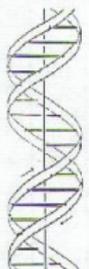
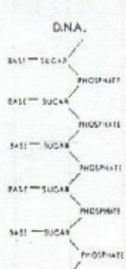


Fig. 1. Chemical formula of a single chain of deoxyribonucleic acid.

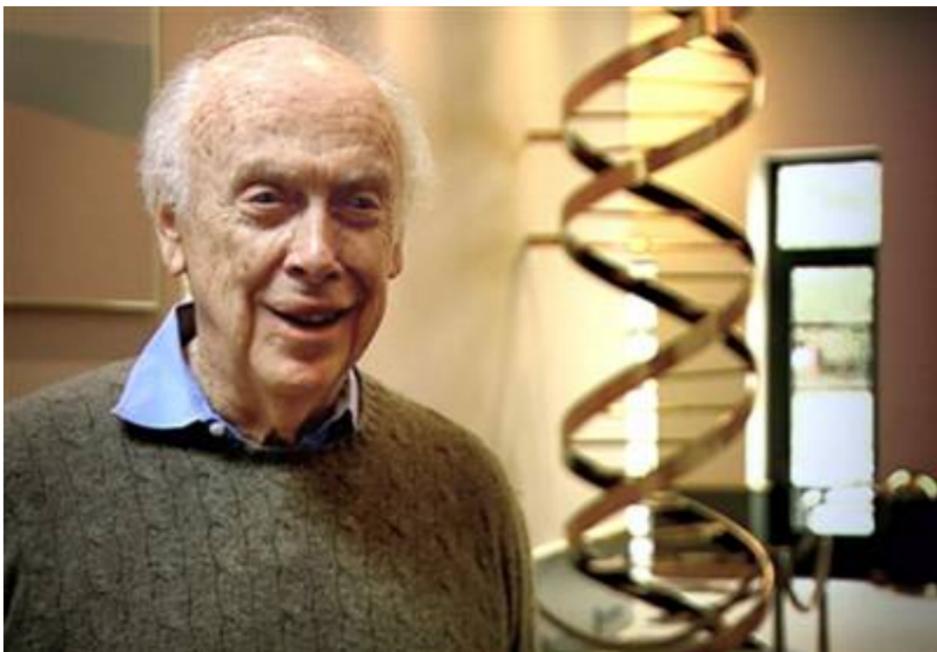
Fig. 2. This space is purely schematic. The two helices revolved at the two phosphate groups shown and the horizontal lines mark the pairs of bases linking the chains together. The vertical line marks the line axis.



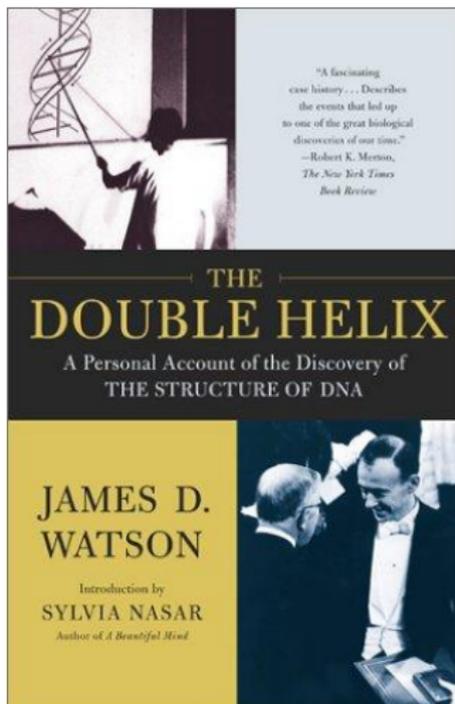
It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material.

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

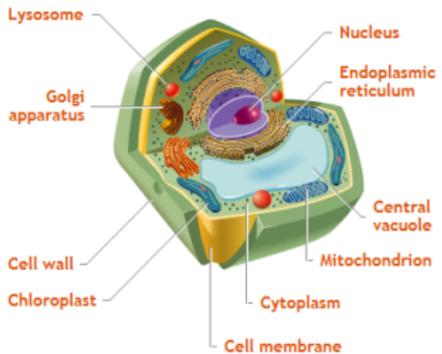
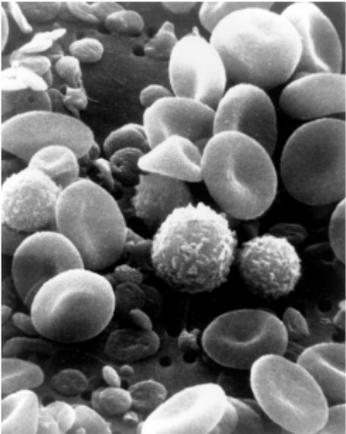




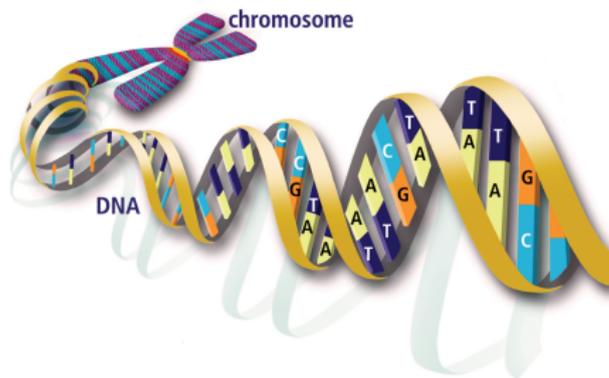
LETTURA INTERESSANTE



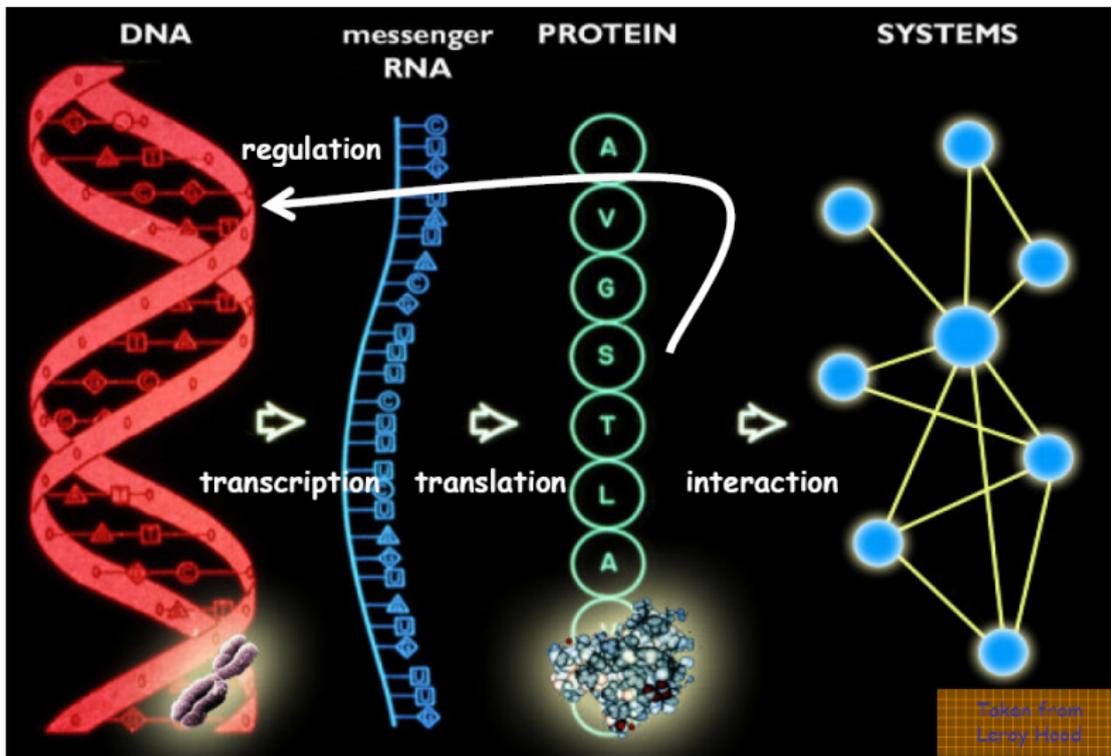
DA COSA PARTIAMO



DA COSA PARTIAMO

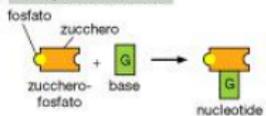


DA COSA PARTIAMO



MEMORIZZAZIONE DELL'INFORMAZIONE

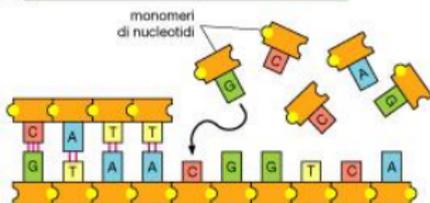
(A) componenti del DNA



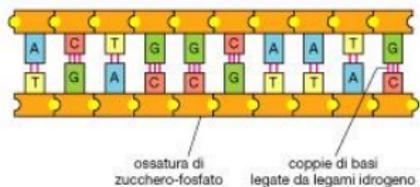
(B) filamento di DNA



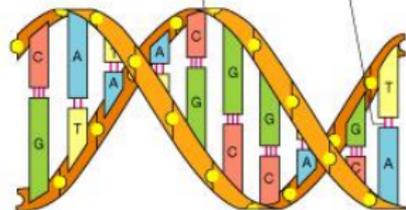
(C) polimerizzazione su stampo del nuovo filamento



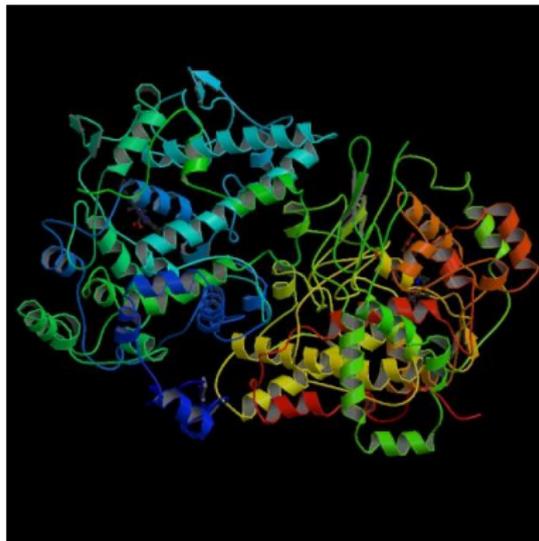
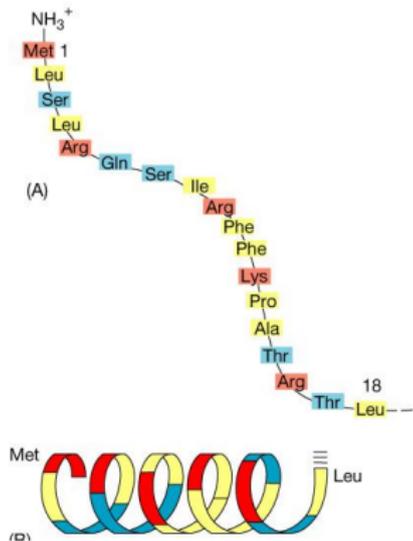
(D) DNA a doppio filamento



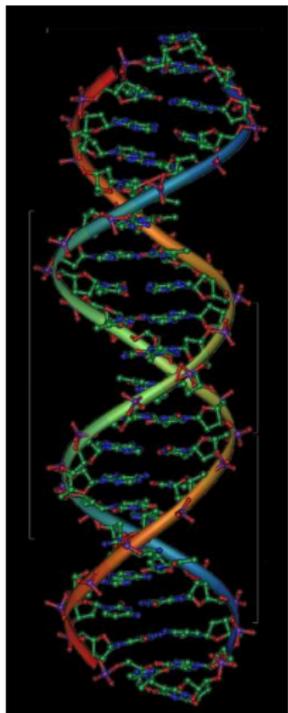
(E) doppia elica di DNA



MEMORIZZAZIONE DELL'INFORMAZIONE



DNA



LE CARATTERISTICHE

- un “turn” di DNA: 3.4nm 10.5 bases
- una cellula: 2m di DNA
- un uomo: 10^{13} cellule
- un mucchio di DNA!

Estremamente delicato da manipolare

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 mezzi è impossibile a intenderne umanamente parola; senza questi è un aggirarsi vanamente per un oscuro laberinto.

Galileo Galilei

IL NOSTRO OBIETTIVO

TCAGGGCCAGCAGCAAAGTTTCCAGCTGATCATCCTGATGGTGC GCGGGTGAGTAAGACACCACCTTCCGG
 GAGGTTTTGGCGAAAAACGTCTGATAGCGGGCTGACGGGCAAGATCCTCAAGACTTTGCCGGCCTTTTTTC
 TCAATTTCAACACGTTTTTGCTATGAATGGTACTGTCTGATACATCGCAAACCTGCGCGCAATGGGTCAGCC'
 ATTTGGTCAGCCTTAGGATTTAAACCATGATGCACATGCAGCTGCTGCCAGGAAAAGCCGTCTCGTGTGAGA
 TATATAACGGGCCAGCAATTC AAGTAACACGCGGGAATCGAGGCCGCCGTGAATCCAACCAGTAAGTGC
 CCGGGCCGGATCAGCTGATCCAGCACAGAAAAGACACTTGGCTTCCAATTCATGATGACTCACAACAGATA
 CTTTCATCGATAAAGTACAATATGCGGATATTAGCAGATCCTGTTTGT CATTACGTGAACCGGGATAATAA
 GCCATATTGCTTGGCGCTCGTTTACACGATAACGATGATCGGATGGAAGAAATAAAAAAAGCCCCTGAGC
 GGGGCTTAAAGAAAAGGAAAGCATTAACTTACGGGTTTTCCAGCGGAGCGAAGCTCTTCAACAAGATCA
 TCAATGGCTTTTCATCTGTTG CAGGAATGGCTCCAGCTTATCCAGTGGTAATGCACTAGGACCATCAACAAC
 GTGCATGTTCCGGATCAGGATGCGACTCAATAAACAAGCCAGCCAGACCAACCGCCATACCAGCGCGTGC
 CAGTT CAGTCACTTGTCTACGACGACCACTTGAAGCCGCACCCAGCGGGT CACGACACTGCAGGGAATGC
 GTTACGTCAAAAATCACCGGGCTGCCCTGGCTGGCATCTTTCATCACGCCGAAGCCAGCATGTCGACCA
 CCAGATTGT CATAACCGAAATTCACACCGCGTTTACACAGAATAATACGGTCGTTACCGCATTCGGCGAA
 TTTCTCAACGATGTTCTTAACTGTCCGGGCTCAGGAAGCTGCGGTTTTTTCAGTTAATCAGTCTTACC
 GTTTTAGCCAGAGCTTCCACCAGATCAGTCTGACGCGCAAGAAGGCCGGTAATTCAGCACATCCACG
 CTTACCCGATCGGTTTGGCTTGCCATGTTTCATGCACATCTGTGATCAGGCTGACACCGAAGGTTTTCTT
 CAGTTCTTCAAAAATACGCAGACCTTCTTCCATACCCGGACCACGGTAGGAATGAAC TGAGGAGCGGTTG
 GCTTTATCCCAGCTGGCTTTAAATACCAAAGGGATGCCAGTTTTTTCAGTTACGGTGACGTGAGGTTTCA
 CAGTGCAGATCGCCAGTCCCGTGACTCCAGTACATTTGCCACCAACAGTACGAATCGTGTATCATT
 AGCAACATTAATCCGTTAACTGAACGACTTCTGTTTTATTTCATCACTCCATCAGTGA AAAAATTATGG
 CGTCATGATCCAGCATCGCCAGCTGTACTTT CAGTACGGCTGCAACCGGATCCTGAGGACATTGTTCAAT
 AAAATAGGTGTAATCATGAGCTGCCAGCTGCGGACATTCAGCTGTTTCATATACCAGCCCGGATCACGG
 ATTTTCATAGGGGTCATCGGGTTTCATCGCCAGCAAACCTCGCTG CAGCGTAAACGCTCCGGCAAACAAC
 GACTTTGTAACATAACAACCTTTAATCAGCTCAATAAACGTGACATTATCTGACGCTGATCTGTCTCTTT
 TGTATTTTACTATCCATCCGGTTAAGTCCACGAGTGTGGCTGCAGCATTAATGACTGTGTCATATA
 TCCCACCTTTCACCGGTAACCGGATCGATTAATACAGGCTTACTTTCCGGGAAAAGCAGCAAGAAGTTAC
 CAGGGAAGCAAACACCACGCACCGGCAATTCATCGAATGAGCCAGATGCATTAACACGATACCCAGCGT
 GAGCGCAAACCAGTGCCTGCATCAGCACCTGATCCAGCAGGCAGTTTTCCGGTGCATAGTAATCTGC
 CAGTCAACCGGAAAATTGCAATTCATGATAGAACGCATGCCAACAGTTGTTCAACGACCCCTTTATCATCAG
 GGGCCGTAATACGGCCGAGCTCTGTACACAAGCGGTAGCCGCTGCAATCCGCTCGATCAATACCGG
 TTTATCTTCTAGATCTTCTCAGCCATCAATCCAGCATGCCAATATCCAGCCGCTTCAAAAATCAGTATCC

TERMINI

SEQUENZIAMENTO

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

wet-lab activity

ASSEMBLAGGIO

- *stoccaggio* delle reads (+ altre informazioni)
- *assemblaggio* delle reads (corte) in *contigs* (lunghe)

dry-lab activity

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

L'INPUT DEGLI ALGORITMI DI *assemblaggio*

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

L'INPUT DEGLI ALGORITMI DI *assemblaggio*

GACTTTGTAACATACAACCTTTAATCACGCTCAATATGACATTATCTGACGCTGATCTGTCTC...

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

L'INPUT DEGLI ALGORITMI DI *assemblaggio*

GACTTTGTAACATAACAACCTTTAATCACGCTCAATATGACATTATCTGACGCTGATCTGTCTC...

GACTTTGTA ACATAACAAC CTTTAATCACG CTCAATATGACAT TATCTGACGCTGA TCTGTCTC...

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

L'INPUT DEGLI ALGORITMI DI *assemblaggio*

GACTTTGTAACATACAACCTTTAATCACGCTCAATATGACATTATCTGACGCTGATCTGTCTC...

CTTTAATCAG
ACATACAAC
CTCAATATGACAT
TATCTGACGCTGA
TCTGTCTC
GACTTTGTA

...

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

L'INPUT DEGLI ALGORITMI DI *assemblaggio*

GACTTTGTAACATACAACCTTTAATCACGCTCAATATGACATTATCTGACGCTGATCTGTCTC...

CTTTAATCAGG
ACATACAAC
CTCAATATGACAT
TATCTGACGCTGA
TCTGTCTC
GACTTTGTA
...

?

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

LA SOLUZIONE: *tante* COPIE

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

LA SOLUZIONE: *tante* COPIE

GACTTTGTA

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

LA SOLUZIONE: *tante* COPIE

GACTTTGTA

TTTGTAACATAC

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

LA SOLUZIONE: *tante* COPIE

GACTTTGTAACATAACAAC
TTTGTAACATAC

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

LA SOLUZIONE: *tante* COPIE

```
GACTTTGTAACATAACAAC  
TTTGTAACATAACAACCTTTAA
```

LA COMBINATORICA DEL PROBLEMA

VINCOLI

- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

LA SOLUZIONE: *tante* COPIE

```
GACTTTGTAACATAACAACCTTTAATCAG  
TTTGTAACATAACAACCTTTAA
```

LA COMBINATORICA DEL PROBLEMA

VINCOLI

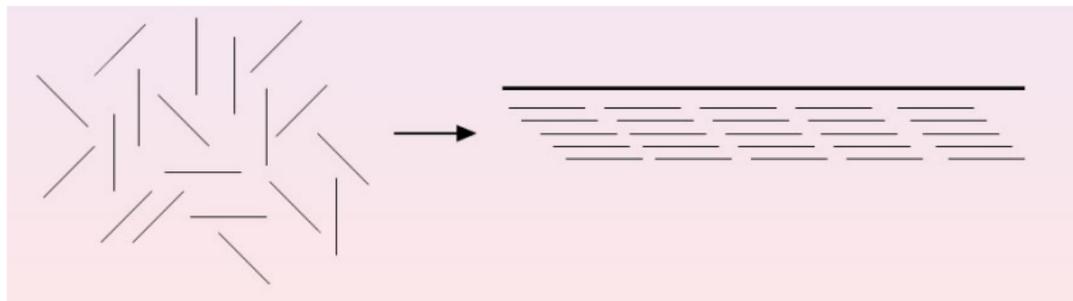
- L'attività di *sequenziamento* produce **reads** di lunghezza limitata
- Il **costo** di un progetto è proporzionale al numero di reads prodotte

LA SOLUZIONE: *tante* COPIE

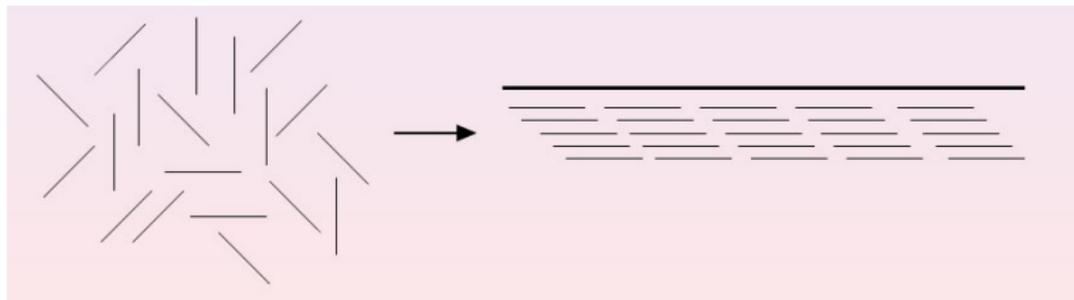
GACTTTGTAACATAACAACCTTTAATCACG . . .

TTTGTAACATAACAACCTTTAATCACGCTCAATA . . .

RIPULIAMO “MATEMATICAMENTE” IL PROBLEMA



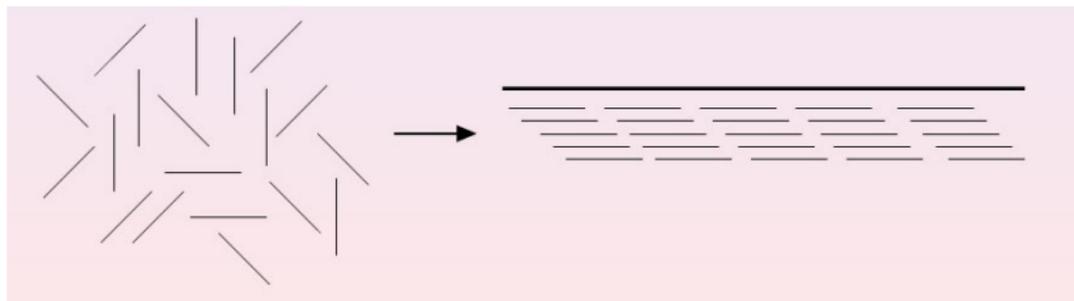
RIPULIAMO “MATEMATICAMENTE” IL PROBLEMA



DEFINIZIONE (SHORTEST SUPERSTRING PROBLEM)

Dato un insieme di stringhe $\{s_1, s_2, \dots, s_n\}$ trova la più corta super-stringa T tale che tutte le s_i siano sotto-stringhe di T .

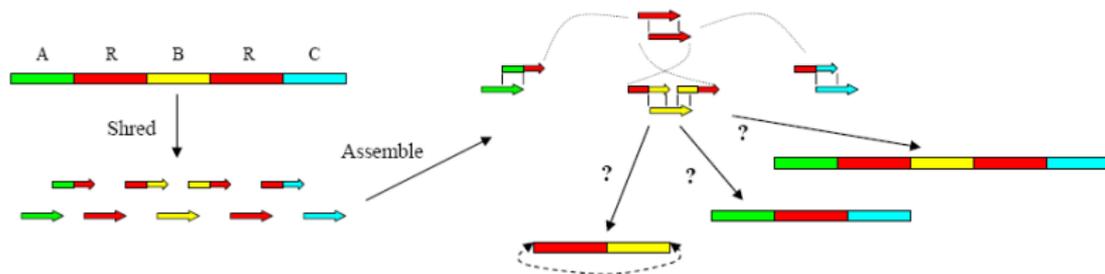
RIPULIAMO “MATEMATICAMENTE” IL PROBLEMA



DUE PROBLEMI:

- 1 \mathcal{NP} -completezza [Gallant et al. 1980]
- 2 Il problema dell'assemblaggio è formulato in modo scorretto!

ELEMENTI RIPETUTI



MORALE

ASSEMBLAGGI DI “QUALITÀ”

- Alte “coperture” (tanti *genomi equivalenti*)
- Tante reads
- **reads lunghe**
- Buon sw (*assembler*)
- Buon hw (parallelo?)
- ...



L'allineamento di sequenze

INPUT

```
...GTTGATTAGCTTATCCCAAAGCAAGGCACTGAAAATGCTAG  
GTGATGTAGCTTAACCCAAGCAAGGCACTAAAATGCCTAG ...
```

L'allineamento di sequenze

INPUT

```
...GTTGATTAGCTTATCCCAAAGCAAGGCACTGAAAATGCTAG  
GTGATGTAGCTTAACCCAAGCAAGGCACTAAAATGCCTAG ...
```

OUTPUT

```
...GTTGAT_TAGCTTATCCCAAAGCAAGGCACTGAAAATG_CTAG  
GT_GATGTAGCTTAACCCAAGCAAGGCACTAAAATGCCTAG...
```

L'allineamento di sequenze

INPUT

...GTTGATTAGCTTATCCCAAAGCAAGGCACTGAAAATGCTAG
 GTGATGTAGCTTAACCCAAGCAAGGCACTAAAAATGCCTAG ...

	G	T	T	G	A	T	T	A	G	C	T	T	A
G	0	1	2	3	4	5	6	7	8	9	10	11	12
T	1	0	1	2	3								
G	2	1	1	1	2								
A	3	2	2	2	1								
T	4	3	2										
G	5	4	3										
T	6	5	4										
A	7	6	5										

OUTPUT

...GTTGAT_ TAGCTTATCCCAAAGCAAGGCACTGAAAATG_ CTAG
 GT_ GATGTAGCTTAACCCA_ GCAAGGCACTAAAAATGCCTAG...

UN ALTRO PROBLEMA ALGORITMICO (SEMPLICE?)

INPUT

T testo su un alfabeto Σ e P pattern su Σ .

UN ALTRO PROBLEMA ALGORITMICO (SEMPLICE?)

INPUT

T testo su un alfabeto Σ e P pattern su Σ .

OUTPUT

Tutte le occorrenze di P in T

UN ALTRO PROBLEMA ALGORITMICO (SEMPLICE?)

INPUT

T testo su un alfabeto Σ e P pattern su Σ .

... è naturale aspettarsi $|P| \cdot |T|$ operazioni di confronto ma si può fare molto meglio!

OUTPUT

Tutte le occorrenze di P in T

Pitagora: “Tutto è numero”

IDEA

Reads (corte) *sono* Numeri (grandi)

FONDAMENTALI

DEFINIZIONE

- $\Sigma_{dna} = \{a, c, g, t\}$ alfabeto;
- P pattern $|P| = m$;
- T testo and $|T| = n$;

T_s è la stringa di m caratteri $T[s, \dots, s + m - 1]$.

STRINGHE IN Σ_{dna} SONO NUMERI IN BASE 4

$$P \rightsquigarrow p \quad T_s \rightsquigarrow t_s$$

FONDAMENTALI

DEFINIZIONE

- $\Sigma_{dna} = \{a, c, g, t\}$ alfabeto;
- P pattern $|P| = m$;
- T testo and $|T| = n$;

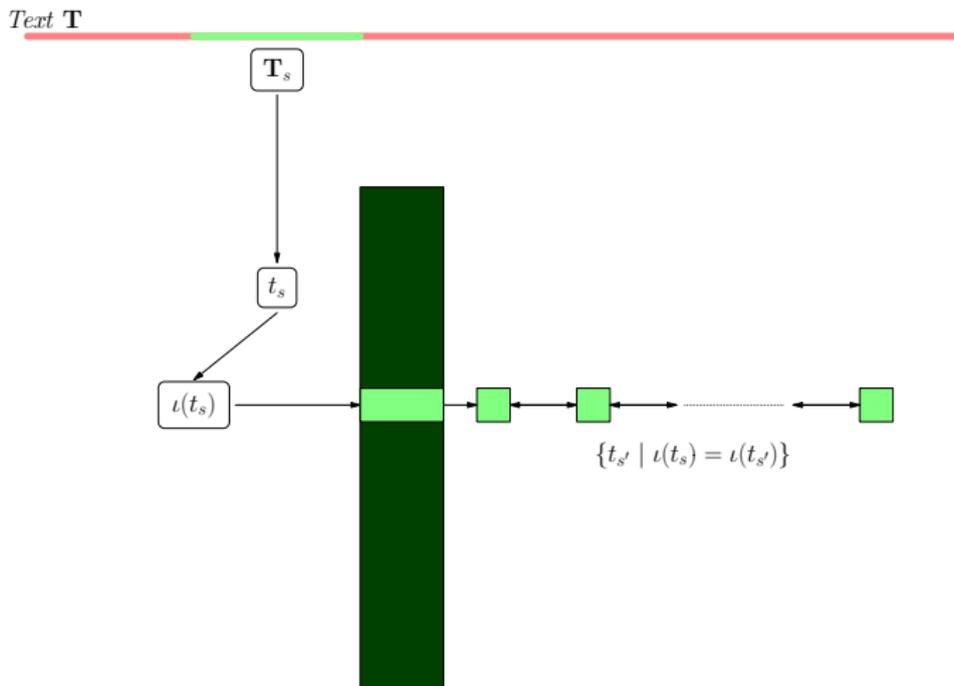
T_s è la stringa di m caratteri $T[s, \dots, s + m - 1]$.

STRINGHE IN Σ_{dna} SONO NUMERI IN BASE 4

$$P \rightsquigarrow p \quad T_s \rightsquigarrow t_s$$

$$a \equiv 0 \quad c \equiv 1 \quad g \equiv 2 \quad t \equiv 3$$

INDICE

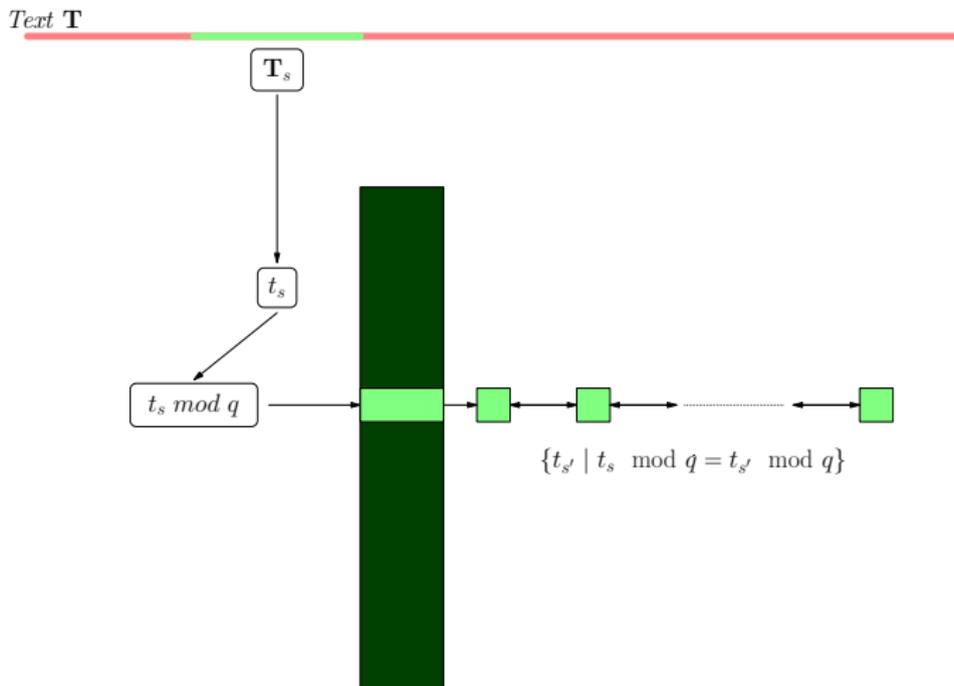


INDICE

PROBLEMA

I numeri diventano è GRANDI

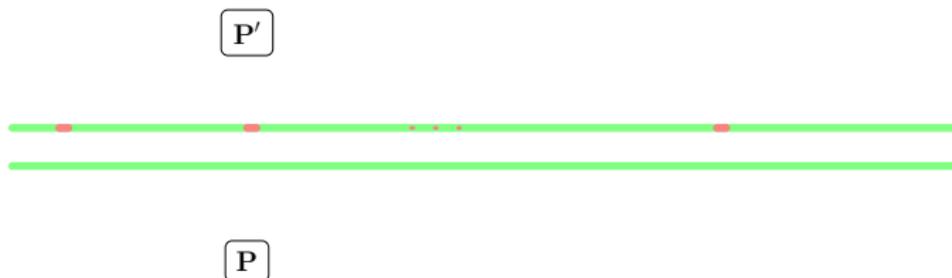
INDICE



... E GLI ERRORI?



... E GLI ERRORI?



... E GLI ERRORI?

$P' - P$

... E GLI ERRORI?



$P' - P$

$$|\{P' \mid d_H(P, P') \leq d\}| = \sum_{d' \leq d} \binom{m}{d'} (\Sigma - 1)^{d'}$$

... E GLI ERRORI?

$P' - P$

$$|\{P' \mid d_H(P, P') \leq d\}| = \sum_{d' \leq d} \binom{m}{d'} (\Sigma - 1)^{d'}$$

$$p' - p \in \mathcal{Z}(d, m)$$

... E GLI ERRORI?

$$\boxed{\mathbf{P}' - \mathbf{P}}$$

$$|\{\mathbf{P}' \mid d_H(\mathbf{P}, \mathbf{P}') \leq d\}| = \sum_{d' \leq d} \binom{m}{d'} (\Sigma - 1)^{d'}$$

$$p' - p \in \mathcal{Z}(d, m)$$

$$(p' - p \bmod q) \in ???$$

RNA

RANDOMIZED NUMERICAL ALIGNER

- opportuno q^*
(\Rightarrow “piccolo” \mathcal{Z})
- diamo *struttura* a \mathcal{Z}
(\Rightarrow implementiamo un test di appartenenza veloce)
- euristiche

RNA

RANDOMIZED NUMERICAL ALIGNER

- opportuno q^*
(\Rightarrow “piccolo” \mathcal{Z})
- diamo *struttura* a \mathcal{Z}
(\Rightarrow implementiamo un test di appartenenza veloce)
- euristiche

COMPLESSITÀ

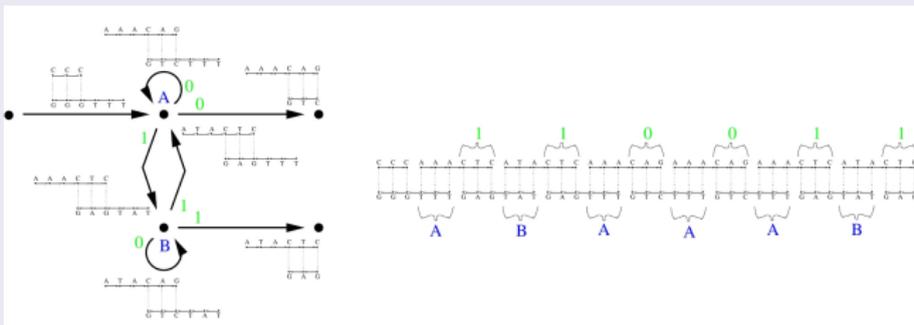
Proporzionale, in media, alla dimensione del pattern.

UNA NOTTE DEL 1993 L. ADLEMAN STAVA LEGGENDO
“THE MOLECULAR BIOLOGY OF THE GENE”...

... SI SEDETTE SUL LETTO E DISSE A SUA MOGLIE: “DIO MIO,
QUESTE COSE POSSONO CALCOLARE “

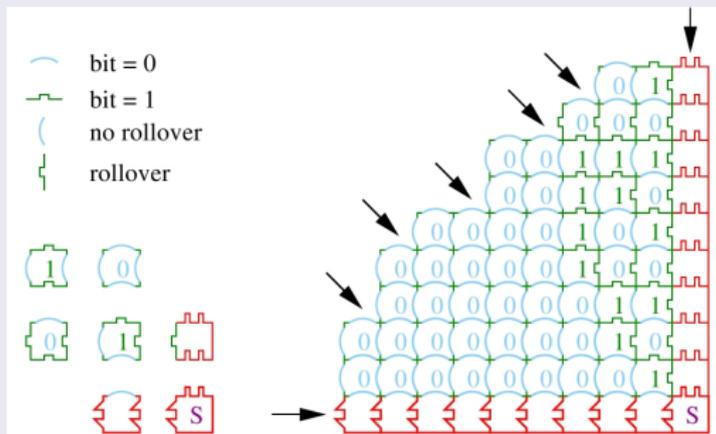
UNA NOTTE DEL 1993 L. ADLEMAN STAVA LEGGENDO “THE MOLECULAR BIOLOGY OF THE GENE”...

... SI SEDETTE SUL LETTO E DISSE A SUA MOGLIE:”DIO MIO,
QUESTE COSE POSSONO CALCOLARE “



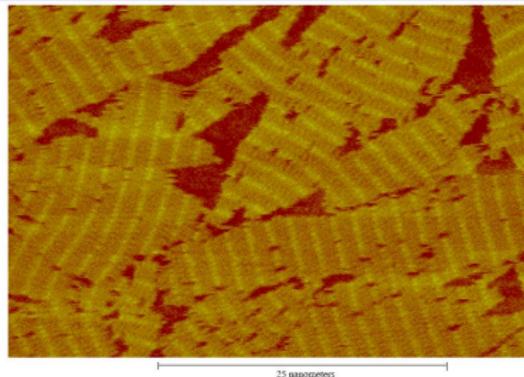
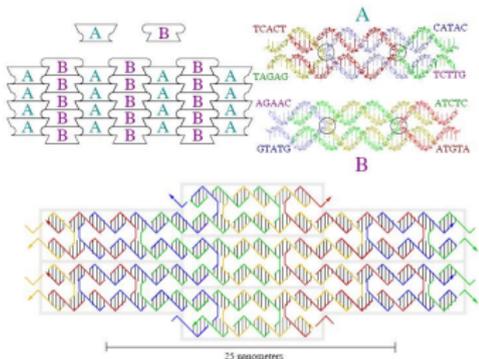
DNA SELF-ASSEMBLY

SI SA DA TEMPO CHE È POSSIBILE COMPUTARE “*in 2D*”



DNA SELF-ASSEMBLY

OGGI SI POSSONO SINTETIZZARE MOLECOLE DI DNA CHE ESEGUONO TALI COMPUTAZIONI!



DNA SELF-ASSEMBLY

PER SAPERNE DI PIÙ

DNA Computing by Self-Assembly

ERIK WINFREE

*Departments of Computer Science and Computation & Neural Systems
California Institute of Technology
Pasadena, California*

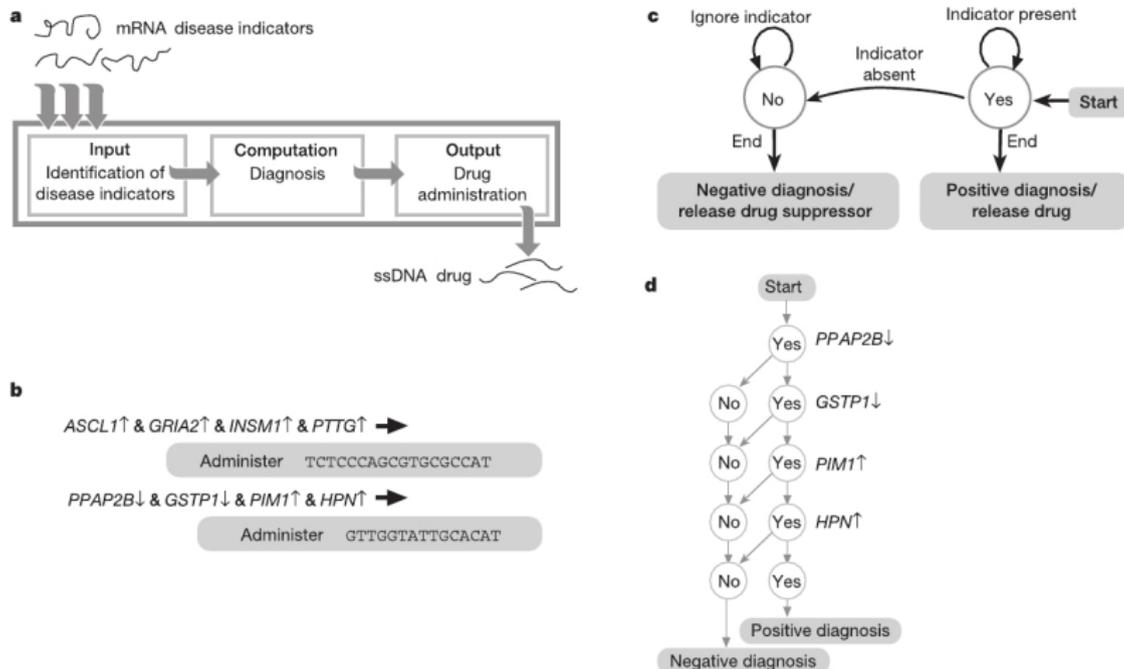
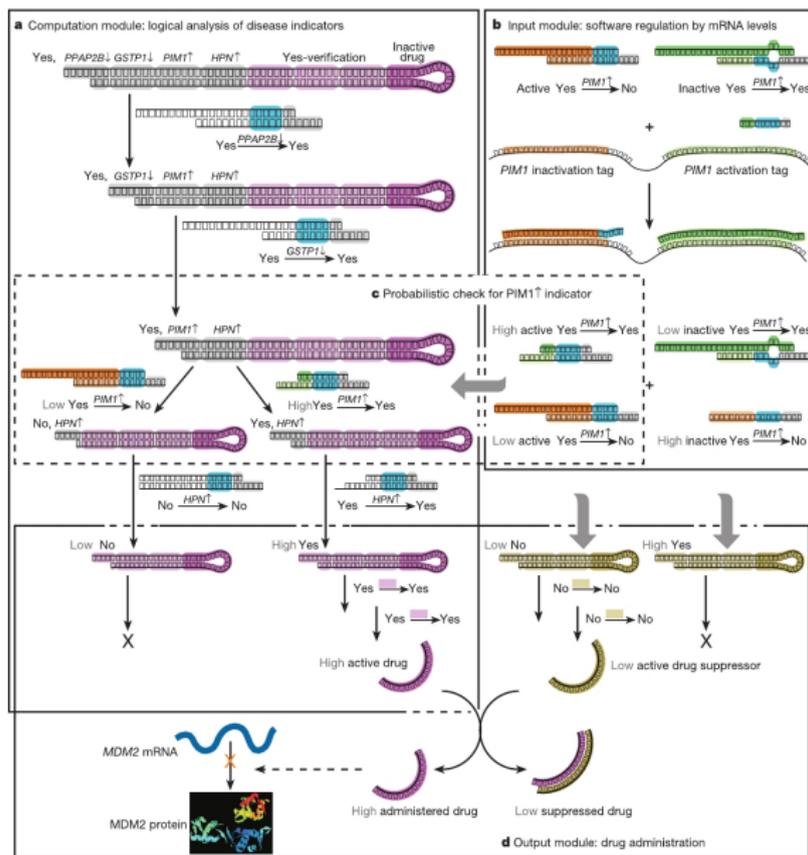


Figure 1 Logical design and logical operation of the molecular computer. **a**, Function and modular organization of the molecular computer. **b**, Example diagnostic rules for simplified models of SCLC¹⁹ and prostate cancer²⁰, indicating overexpression (\uparrow) or underexpression (\downarrow) of a disease-related gene. The first rule states that if genes *ASCL1*, *GRIA2*, *INSM1* and *PTTG1* are overexpressed then administer the ssDNA molecule

TCTCCCAGCGTGCCCAT (oblimersen), purported to be an antisense therapy drug for SCLC²⁷. The second rule states that if the genes *PPAP2B* and *GSTP1* are underexpressed and the genes *PIM1* and *HPN* are overexpressed then administer the ssDNA molecule GTTGGTATTGCACAT, purported to be a drug for prostate cancer²⁸. **c**, Transition diagram of the diagnostic automaton. **d**, The computation that diagnoses prostate cancer.



DOMANDA INTERESSANTE

qual è l'hardware e quale il software?

PER SAPERNE DI PIÙ

An autonomous molecular computer for logical control of gene expression

Yaakov Benenson^{1,2}, Binyamin Gil¹, Uri Ben-Dor¹, Rivka Adar²
& Ehud Shapiro^{1,2}

¹Department of Computer Science and Applied Mathematics and ²Department of Biological Chemistry, Weizmann Institute of Science, Rehovot 76100, Israel

Early biomolecular computer research focused on laboratory-scale, human-operated computers for complex computational problems¹⁻⁷. Recently, simple molecular-scale autonomous programmable computers were demonstrated⁸⁻¹³ allowing both input and output information to be in molecular form. Such computers, using biological molecules as input data and biologically active molecules as outputs, could produce a system for 'logical' control of biological processes. Here we describe an autonomous biomolecular computer that, at least *in vitro*, logically analyses the levels of messenger RNA species, and in response produces a molecule capable of affecting levels of gene expression. The computer operates at a concentration of close to a trillion computers per microlitre and consists of three programmable modules: a computation module, that is, a stochastic molecular automaton¹²⁻¹⁷; an input module, by which specific mRNA levels or point mutations regulate software molecule concentrations, and hence automaton transition probabilities; and an output module, capable of controlled release of a short single-stranded DNA molecule. This approach might be applied *in vivo* to biochemical sensing, genetic engineering and even medical diagnosis and treatment. As a proof of principle we

letters to nature

programmed the computer to identify and analyse mRNA of disease-related genes¹⁸⁻²³ associated with models of small-cell lung cancer and prostate cancer, and to produce a single-stranded DNA molecule modelled after an anticancer drug.

Taking our cue from the terminology of medical treatment, we consider that our molecular computer performs *in vitro* a computational version^{23,24} of 'diagnosis'—the identification of a combination of mRNA molecules at specific levels, which in our example is a highly simplified model of cancer—and 'therapy'—production of a biologically active molecule, which in our case is a drug-like single-stranded (ss)DNA with known anticancer activity (Fig. 1a). The computer operation is governed by a 'diagnostic rule' that encodes medical knowledge in simplified form (Fig. 1b). The left-hand side of the rule consists of a list of molecular indicators for a specific disease, and its right-hand side indicates a molecule to be released, which could be a drug for that disease. For example, the diagnostic rule for prostate cancer states²⁰ that if the genes *PPAP2B* and *GSTP1* are underexpressed and the genes *PIM1* and hepsin (*HPN*) are overexpressed then administer the ssDNA molecule GTTGGTATTGGACATG, which inhibits²⁵ the synthesis of the protein MDM2 by binding to its mRNA. The computer design is flexible in that any sufficiently long RNA molecule can function as a molecular indicator and any short ssDNA molecule, up to at least 21 nucleotides, can be administered.

The computation module is a molecular automaton¹²⁻¹⁵ (Supplementary Fig. S1) that processes such a rule as depicted in Fig. 1c. The automaton has two states: positive (Yes) and negative (No). The computation starts in the positive state and if it ends in that state we call the result 'positive diagnosis', otherwise it is called 'negative diagnosis'. To facilitate rule processing by the automaton, the left-hand side of the diagnostic rule is represented as a string of symbolic indicators, or symbols for short, one for each molecular

CONCLUSIONI

CAMPO STIMOLANTE

- Tanti bei problemi computazionali
- Tecnologie sofisticate
- Forte impatto sulle scienze della Vita

CONCLUSIONI

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.

W. Heisenberg

DEFINIZIONE DI ALGORITMO

INFORMALE (MARKOV, 1960)

Processo di calcolo che soddisfa tre requisiti principali:

- 1 Precisione prescrittiva. *Definitezza*
- 2 Applicabilità a dati iniziali variabili entro limiti prefissati.
Generalità
- 3 Processo orientato ad ottenere un risultato, calcolabile in base ai dati iniziali. *Effettività*

FORMALE (GÖDEL, CHURCH, TURING, ...)

Funzioni ricorsive, λ calcolo. macchine di Turing, ...