Modelli Computazionali per lo Studio di Sistemi di Interesse Biologico (Bioinformatica)

Testo



Anna Tramontano Bioinformatica Ed. Zanichelli

Lezioni *on-line* (http://www.chem.unisa.it/luigicavallo/) (derivate da lezioni tenute dal gruppo BioComputing http://cassandra.bio.uniroma1.it)

Prima lezione.

1. Introduzione alla bioinformatica

2. Evoluzione e informazione

3. Ricerca di geni in genomi di procarioti ed eucarioti

I processi della vita sono dovuti all'azione concertata di (macro) molecole biologiche, principalmente proteine.

Le istruzioni che dirigono la sintesi proteica sono contenute nel genoma.

Dobbiamo quindi

- identificare le istruzioni e decifrarle
- comprendere quando e quanto sono "lette" (regolazione)
- capire quali modifiche subiscono nell'ambiente cellulare
- individuare le loro interazioni

Bioinformatica: scienza multi-disciplinare, al crocevia tra biologia, chimica, matematica, fisica ed informatica, che analizza l'informazione biologica con metodi computazionali al fine di formulare ipotesi sui processi della vita.

Anna Tramontano

BlOinformatica = geni + proteine + informatica (biologia computazionale, biocomputing)

GENE: segmento di DNA che codifica per una proteina specifica e determina UN carattere ereditario

PROTEINA: prodotto di espressione di un gene ed EFFETTORE della funzione biochimica di cui il gene stesso contiene l'INFORMAZIONE Raccolta, archiviazione, organizzazione e interpretazione di dati biologici.

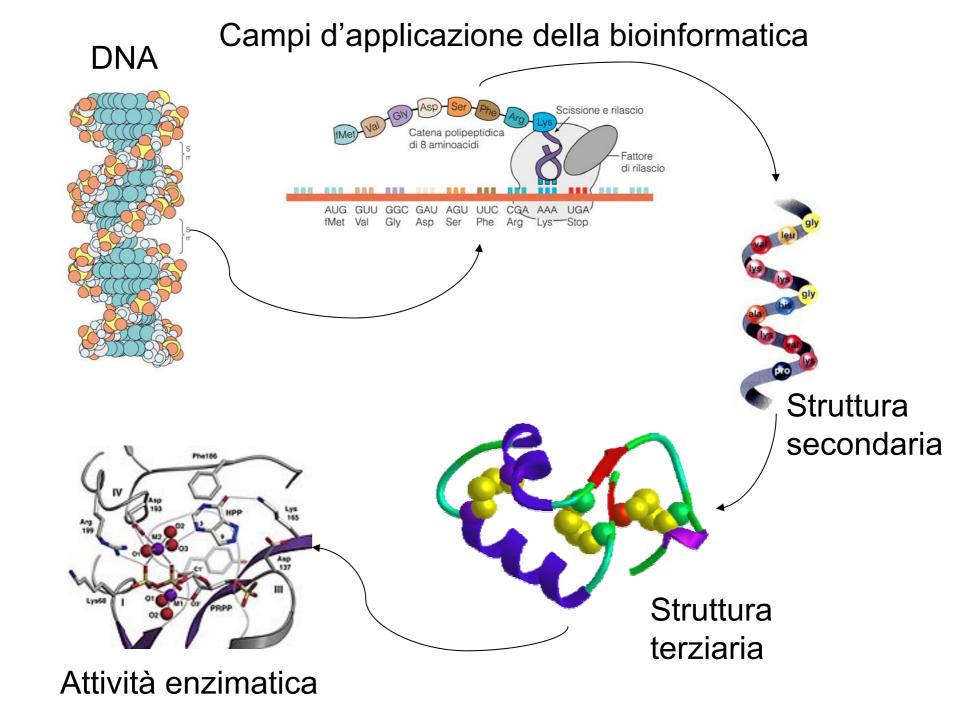
J.M. Thornton

Bioinformatica

La bioinformatica e' cio' che la bioinformatica fa

Eric C. Snowdeal III

- Analisi di sequenze
- Analisi di strutture
- Predizione di strutture
- Disegno di molecole proteiche
- Disegno di inibitori
- Disegno di librerie combinatoriali
- Sviluppo di tools
- Gestione di dati



Cosa ci si aspetta che impariate:

Analizzare una sequenza
Ricercare omologie
Costruire un modello tridimensionale di una proteina
Capire i problemi legati all'analisi di dati genomici e proteomici
Conoscere e usare i principali siti bioinformatici

Conoscere l'affidabilita' dei metodi

Un po' di storia...

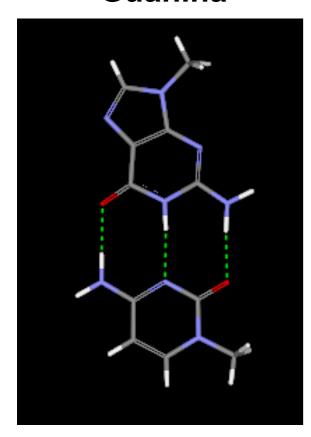
1951	Pauling: alfa e beta	1977	sequenza DNA
1953	Doppia elica DNA	1980	Wutrich
1955	Sequenza insulina	1981	Greer
1959	Struttura mioglobina	1985	FASTP
1960	Anfinsen	1986	Chothia e Lesk
1967	collezione Dayoff	1990	Blast
1968	PAM	1991	Fold recognition
1970	Nedleman and Wunsch	1993	PhD
1977	PDB	1994	CASP
1977	Chou and Fasman		

Un po' di storia...

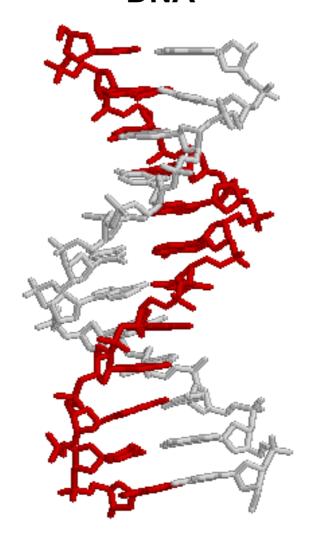
1951	Pauling: alfa e beta	1977	sequenza DNA
1953	Doppia elica DNA	1980	Wutrich
1955	Sequenza insulina	1981	Greer
1959	Struttura mioglobina	1985	FASTP
1960	Anfinsen	1986	Chothia e Lesk
1967	collezione Dayoff	1990	Blast
1968	PAM	1991	Fold recognition
1970	Nedleman and Wunsch	1993	PhD
1977	PDB	1994	CASP
1977	Chou and Fasman		

DNA

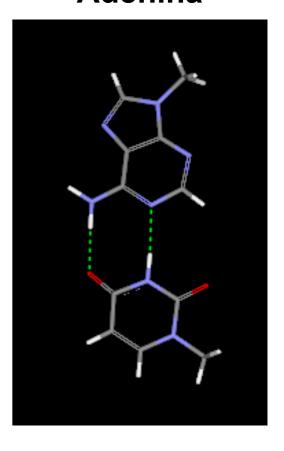
Guanina



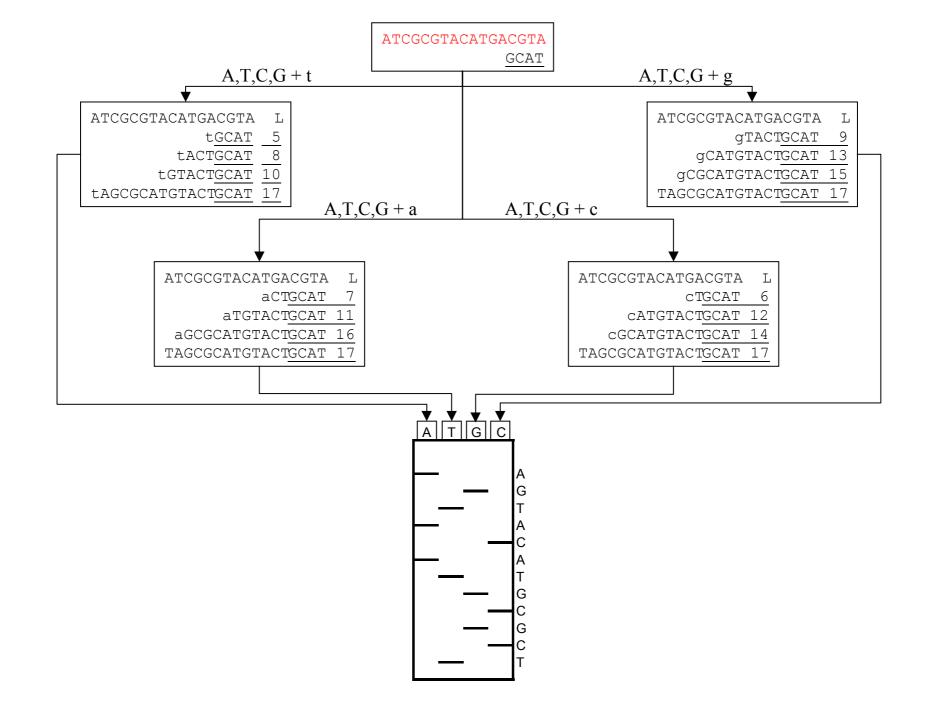
Citosina



Adenina



Timina



ATCGCGTACATGACGTA
GCAT

A,T,C,G + tacg

▼

ATCGCGTACATGACGTA I

tGCAT

tactgcat 8

tGTACTGCAT 10

tAGCGCATGTACTGCAT 17

ATCGCGTACATGACGTA I

gTACTGCAT 9

gCATGTACTGCAT 13

gCGCATGTACTGCAT 15

TAGCGCATGTACTGCAT 17

ATCGCGTACATGACGTA

aCTGCAT 7

aTGTACTGCAT 11

aGCGCATGTACTGCAT 16

TAGCGCATGTACTGCAT 17

ATCGCGTACATGACGTA L

cTGCAT 6

CATGTACTGCAT 12

CGCATGTACTGCAT 14

TAGCGCATGTACTGCAT 17



Il processo e' automatizzabile...



BLUEPRINT OF THE BODY

Overview | Genome guide | Glossary | Related sites | Message board | Story archive | Q&A | Chat Series | Video Archive

Genome announcement 'technological triumph'

Milestone in genetics ushers in new era of discovery, responsibility

June 26, 2000 Web posted at: 12:09 p.m. EDT (1609 GMT)

Urin this story:

Knewledge calchelp treat causes of diseases

Advances could come quickly

RELATED STORIES, SITES **◆**

From staff and wire reports

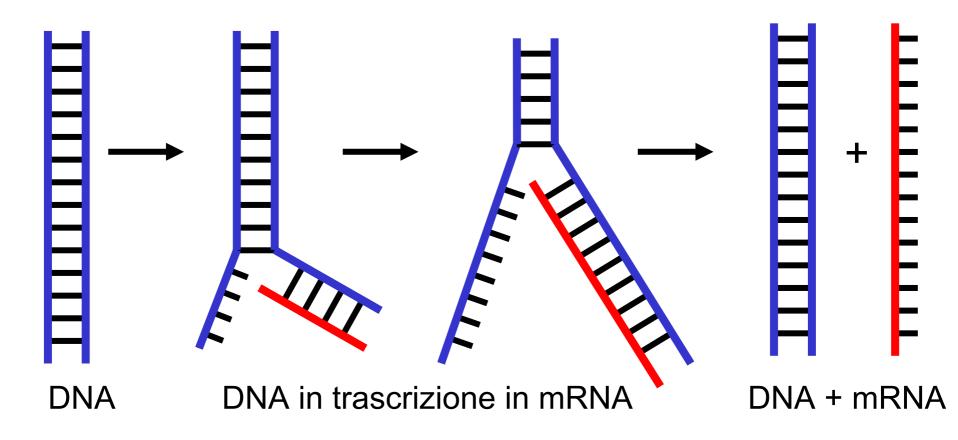
ATLANTA (CNN) — Declaring a new era of medical discovery, U.S. President Bill Clinton and British Prime Minister Tony Blair on Monday praised the efforts of an international team of scientists to decode the genetic makeup of humans.





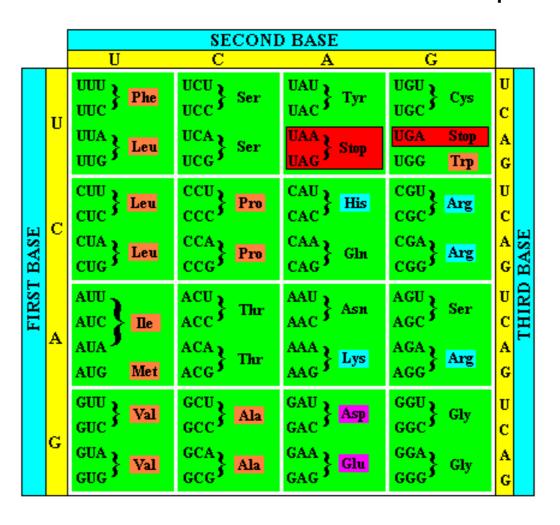
La sintesi proteica richiede la traduzione di sequenze di nucleotidi in sequenze di amminoacidi

 Trascrizione dell'informazione del gene dal DNA allo mRNA da parte di una RNA-polimerasi



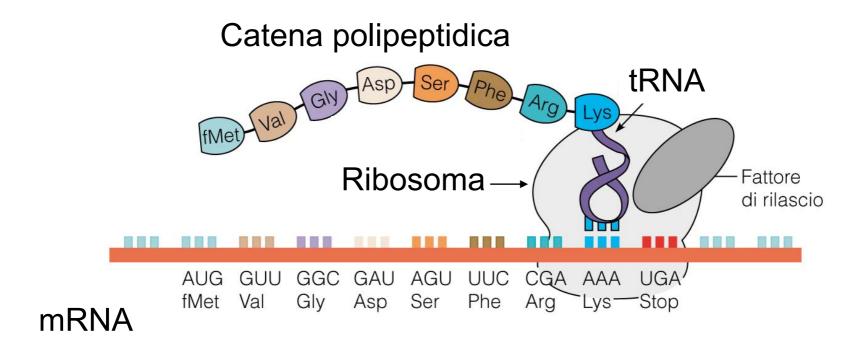
Il codice genetico

La sequenza dei nucleotidi (presi a triplette) dell'mRNA (e del DNA da cui deriva) contiene l'informazione su quali amminoacidi concatenare a formare la proteina



La sintesi proteica richiede la traduzione di sequenze di nucleotidi in sequenze di amminoacidi

2) Traduzione dell'mRNA in proteina da parte dei ribosomi



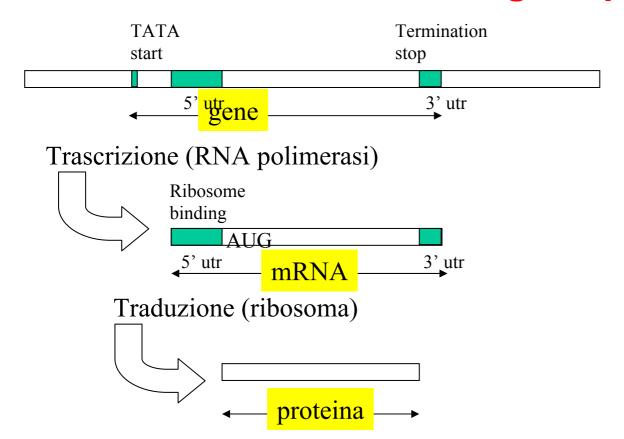
Dove e' il gene?

>cD0826Q1 425-22425 Main

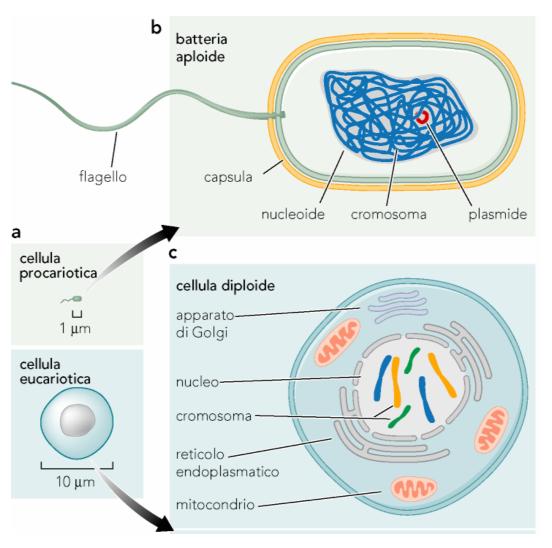
ggcataagaatgatacaatggactttggggacctgagaggaaaggtgggaggggcaagg $\verb|gatactgctcaggtgataggtgcaccaaaatctcacaaatcatcactaaagaacttactc|$ atgtaaccaaatactacctgtaccactataacctacgggggaaaaaagcaacataaccat gaaccaactaataaaaaacaaccttgccttcagtctgcatcctaccctagagacactctc tctgtgtcctcacacttggagctaagcttctgacttttgtctccagtacacccctgagga tcctctcatcacggccatcagaaacctctgtagaaggtcaaatccagtgggttcttgtca tactctqtccttatttttctcctatttactqaatcctccttatcatcctttqaaatctcc tcttaattattatgttctctcatcataccctgagatccctgcatttctgatttttggcac tcttcctggaaaagctcatctaacctgcacctatgcttgatgactctcagttctctggct taaactcctctactgagaccacccatcatacaaaaatgtttacatattatttttccttag ataacttttagatattctaagtgcaatagccccacactgaactcagtctcttctctcagt caggctgtcttctctcattaccctttttaatgaatggaatcaagatgtttgcattgggtt tacatttcatttagattggaggaataattttaagagttttattgtataacatggactata gttgctaacaatgtattgttgaaaattgctaaaagggtggattttaagtgttctcaccac aaaaaataagtatgtgaggtgagccataagttctttagcttgatgtagccggtccatgat gtacatacatttcaaaacaacatattatacatgataaatataaatatttttgtcaatca ggccagtctaggaataagagttatctgggagttttctaagtcggatgccaccgacatcac tcaccaataatccctttaatgtcaatcaaattaagtcctcttcttccatcattttactcc tatgcccatttcctcactctttgttcaggcactattagtcttgcctcttgaaccaacttc tttcactcatgctgcccactgttgccgtagtgatcttcctaaattgcaaatgcgccatca ctctcctqcttaaaatccttcaatqattccttatqacttccaqqacaqaqtaqccactcc $\verb|cttgctgctgttccacatccaaagctggctccattcatactgaagcagctgaagttcttc|\\$ agatatgtcattgccacactgggcccacacttttgaacctgcttcctcctgtgtgagaag tggcttctgccctgttttcggactgcctacattgaagccatctgttccccaggaagcctt ccctgatgccttgacagcatcttgtgcctgcccatatctgcacttatccatctggg cctqctqttqtcttqtcacttqtqttctcttctqtqaactqtaaacatcaqqaqqacaaq acctatqtcttacttttatttqaatatttaqcatctaacaatqttcqacatataqtaqqc ttttgatactatttttttactatgacattgtagtatatgttaatatccagtaggacatag gacagtgtggaaagccaggctgggactagggatgcacttaccttaggtgcaaaatttagg aggataccaaaagaactcagtaataaaagtcaatcatattttaatgaaatatcttaagaa atctaaattaatggaaaatatataatgaacaaaatgtcaaaagagaactattcaaagaaa atggagaagcagaggcagaagaattagtagaatatactggcacataagccaaggaggt ${\tt aaa} gattttccaggaaggaagtagagttggagttcagaagttcaacagaagttcatttcag$ aaatcttaccttggttttgaaatcctttcagagagcagttttacataatgtgagcaatta tttctccttcatccccatcattccagaattgagcttcttctctggcttcagaaatgtggc ggggttttgggggtgaaattaattgactttagggaactccttgaatgctaagttctgttca cctqqaqqaccaqaqqqcacaqaqatqaccacctaqcttctqcctqqqacctaaacaq $\tt ggcagagaaataggaggatcaggtataaagggagcagggaagatgggtctgggcttacag$



Struttura di un gene procariota

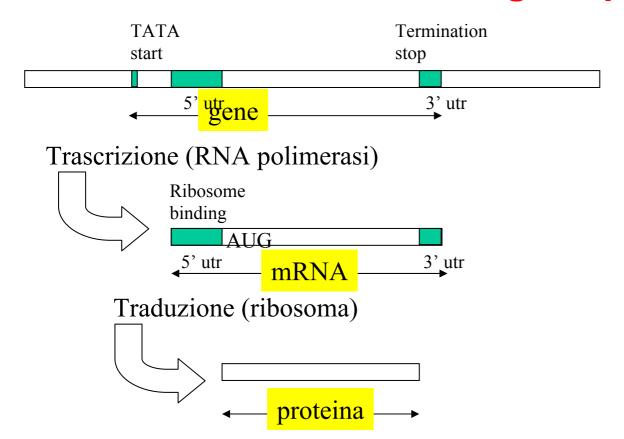


Cellula procariotica

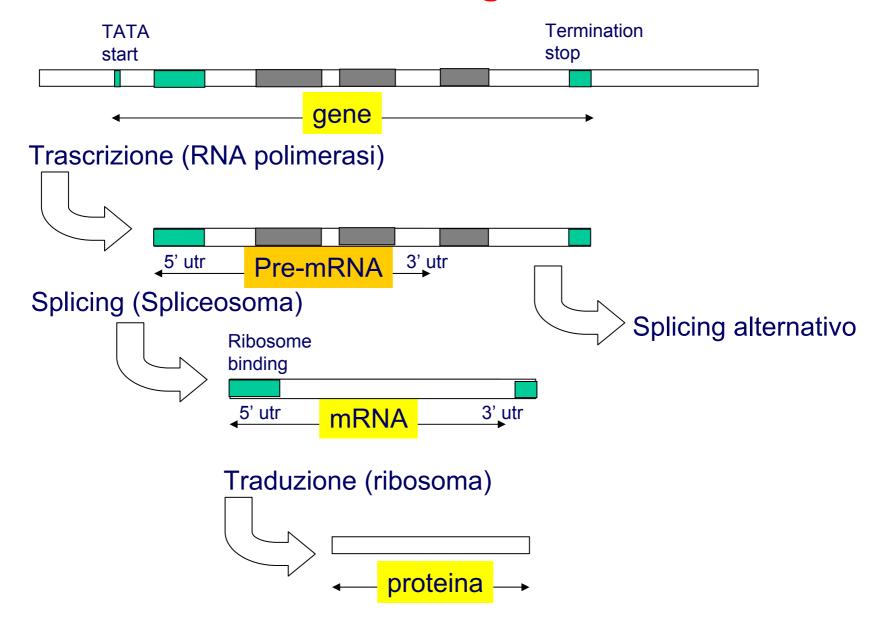


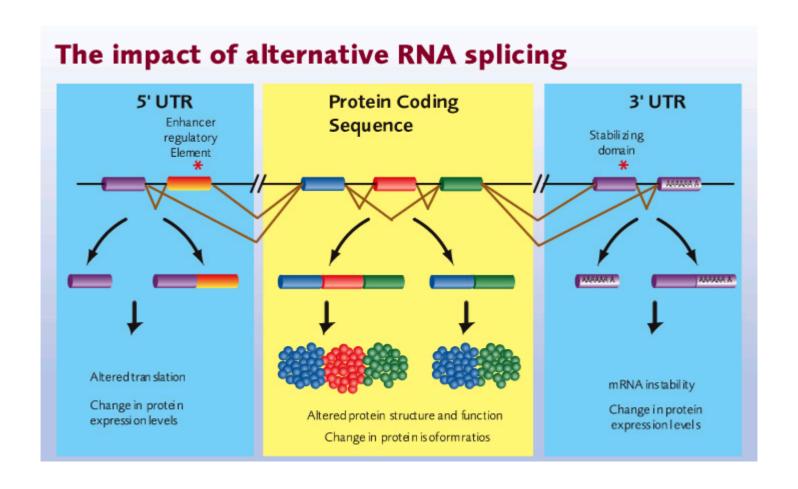
Cellula eucariotica

Struttura di un gene procariota

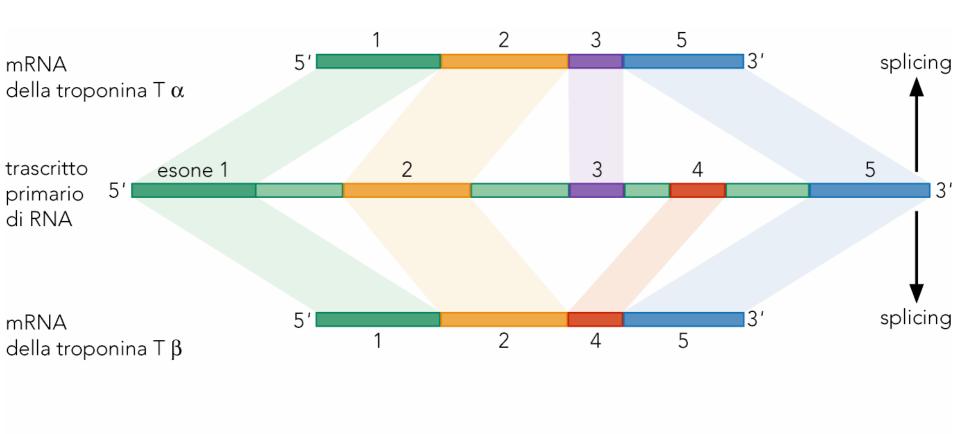


Struttura di un gene eucariota

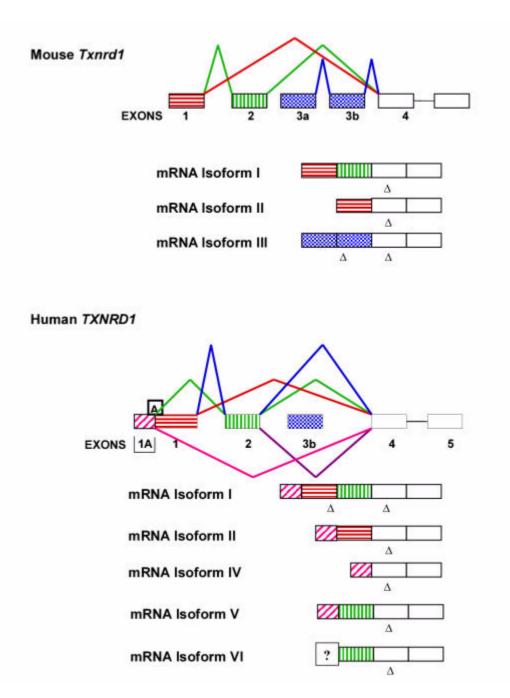


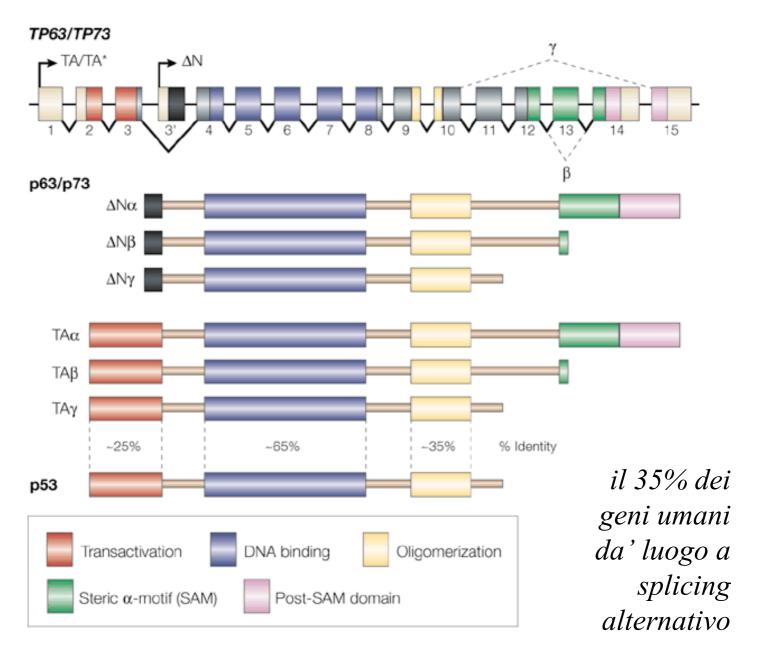


Un esempio: le immunoglobuline possono presentare un dominio che le tiene ancorate alla membrana cellulare, oppure essere solubili

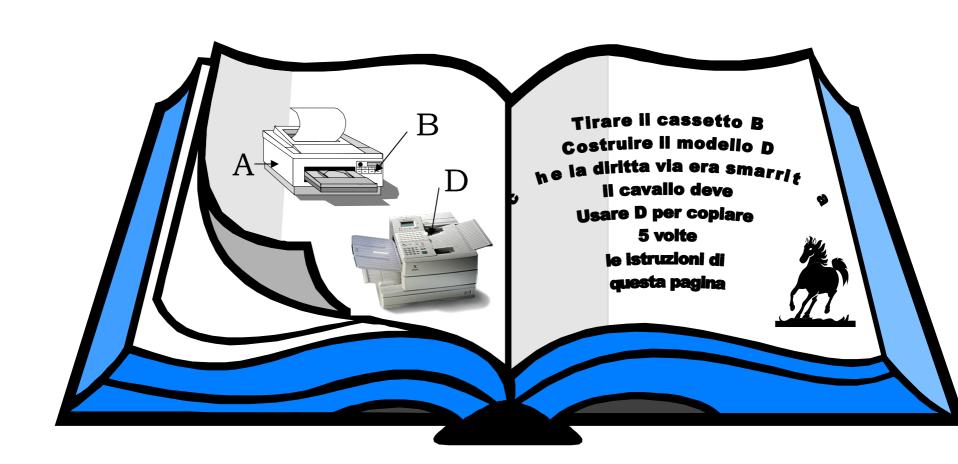


L'organizzazione del gene della tioredossina reduttasi 1 umana e da topo e' altamente conservata ed entrambi I geni danno luogo a splicing alternativo









Un manuale d'uso:

- non include quasi mai l'informazione che cercate
- nei rari casi in cui la include, e' incomprensibile

Ma le cose potrebbero andare peggio...

Se le istruzioni fossero anche divise in 27 paragrafi diversi interrotti da lunghe pagine di informazioni irrilevanti, avremmo qualcosa di molto simile al gene del retinoblastoma umano (e a molti altri geni)

Complessita' dei geni

Geni codificanti per 100Kb

E. coli	87
S. cerevisiae	52
C. elegans	22
H. sapiens	5

•	Numero medio		Lunghezza media	
	di introni per Kb		geni(Kb)	mRNA(Kb)
	S. cerevisiae	0	1.6	1.6
	C. elegans	3	4.0	3.0
	Drosophila	3	11.3	2.7
	Gallus	8	13.9	2.4
	 H. sapiens 	6	16.6	2.2

La lunghezza approssimativa del DNA del genoma umano e' 3 miliardi di basi.

Questa sequenza di nucleotidi (A, G, T, C) <u>non</u> è casuale !!!

Teoria dell'informazione (Claude Shannon)

L'*informazione* e' una *misura universale* dell'ordine e puo' essere applicata a qualunque struttura o sistema.

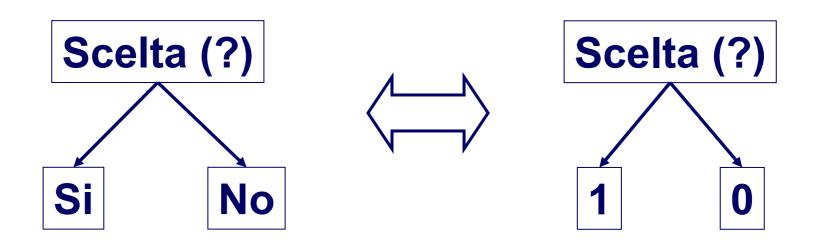
L'ordine si riferisce alla disposizione strutturale del sistema.



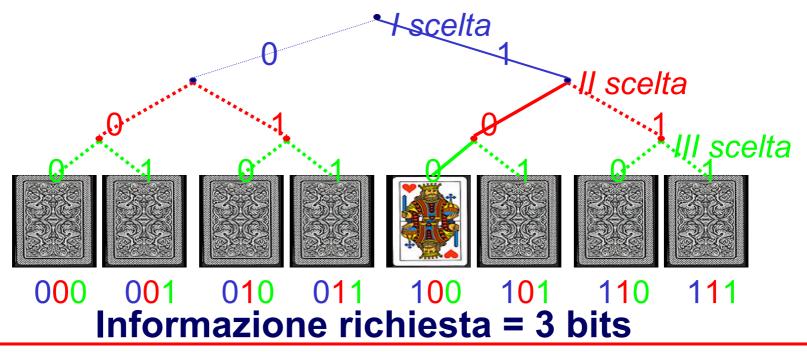


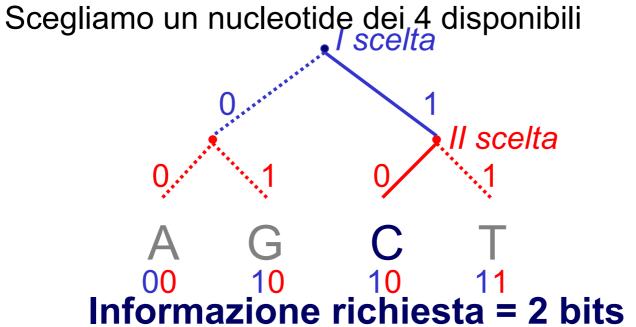
MSKGPAVGIDLGTTYSCVGVFQHG
KVEIIANDQGNRTTPSYVAFTDTE
RLIGDAAKNQVAMNPTNTVFDAKR
LIGRRFDDAVVQSDMKHWPFMVVN
DAGRPKVQVEYKGETKSFYPEEVS
SMVLTKMKEIAEAYLGKTVTNAVV
TVPAYFNDSQRQATKDAGTIAGLN
VLRIINEPTAAAIAYGLDKKVGAE
RNVLIFDLGGGTFDVSILTIEDGI
FEVKSTAGDTHLGGEDFDNRMVNH
FIAEFKRKHKKDISENKRAVRRLR

L'informazione quantifica le istruzioni necessarie a produrre una determinata forma di organizzazione e puo' essere raggiunta in termini di scelte binarie ed espressa in bit.



Scegliamo una carta da una pacchetto di 8...





Quanta informazione contiene una sequenza dinucleotidica, es. GC?

$$N_{(tot)} = y^x = 4^2$$

N(tot) = n di stati possibili y = n di possibili scelte per ogni posizione (A,C, G, T) x = n di posizioni disponibili (1, 2)

La probabilita' che la sequenza nucleotidica GC si generi spontaneamente e' quindi $1/4^2 = 1/16 = 0.0625$

?? Quanto e' ordinato il genoma umano ??

La lunghezza approssimativa del DNA del genoma umano e' 3*10⁹.

Trascurando le mutazioni somatiche spontanee, si puo' dire APPROSSIMATAMENTE che le molecole di DNA di ogni individuo mostrano la stessa sequenza e calcolarne la capacita' di informazione:

$$N_{\text{(tot)}} = 4^{3.000.000.000}$$
 $x = 3*10^9$
 $y = 4$ (A, G, T, C)

Il numero di possibilita' stimato N(tot) e' maggiore del numero stimato di particelle presenti nell'universo!!!

Ma la natura ne sceglie **SOLTANTO UNA**...

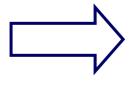
- Il contenuto informativo del genoma degli organismi appartenenti alle varie specie e' altissimo.
- L'informazione riguardante i sistemi biologici in natura si accumula in maniera graduale attraverso i processi di variazione casuale del genotipo e selezione naturale (Charles Darwin, Origin of Species)

Mutazione genotipica neutrale o deleteria per il fenotipo



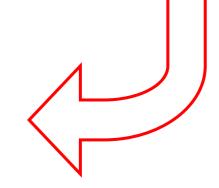
Selezione negativa

Mutazione genotipica vantaggiosa per il fenotipo

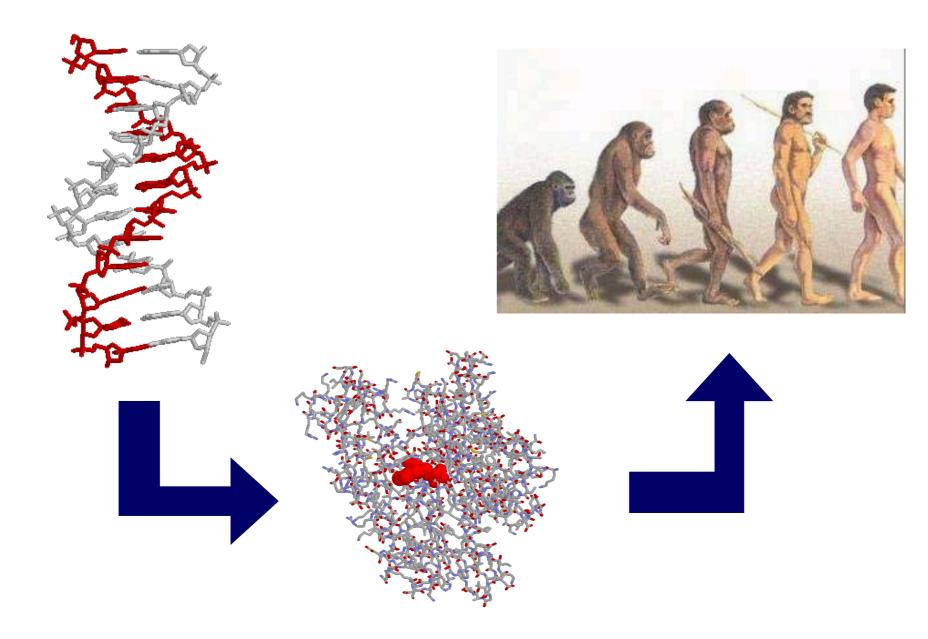


Selezione positiva

NUOVA FUNZIONE BIOLOGICA



EVOLUZIONE MOLECOLARE



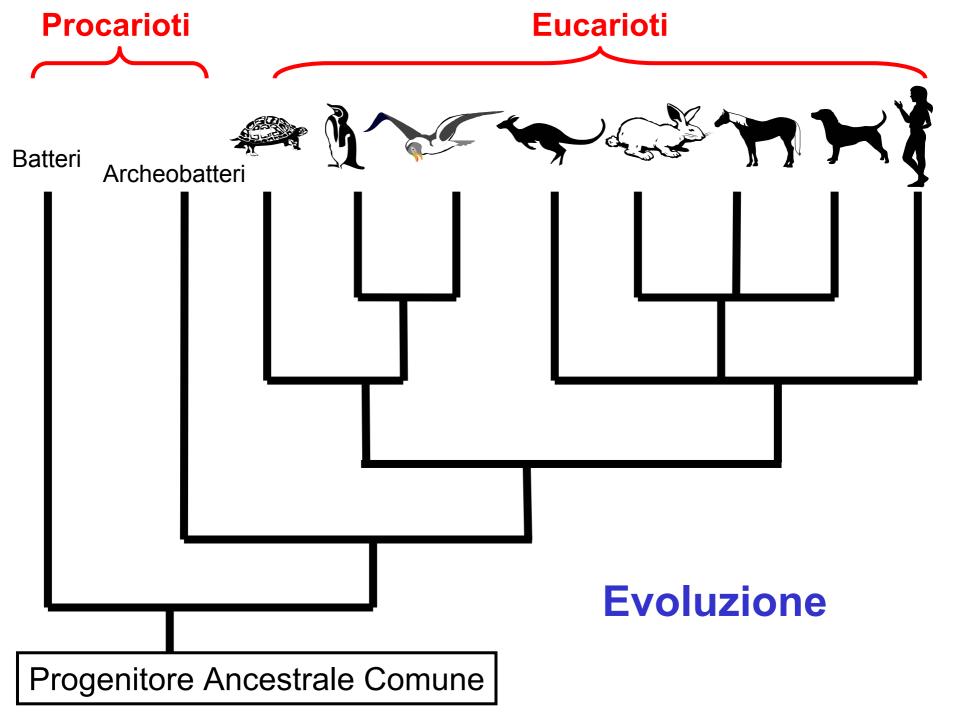
I principi dell'evoluzione

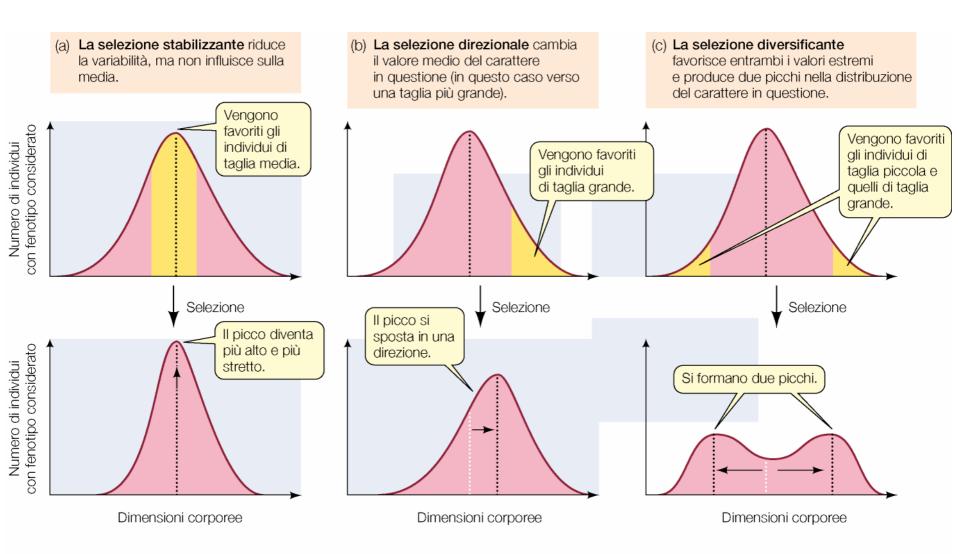
Tutte le specie viventi si sono evolute da altre specie

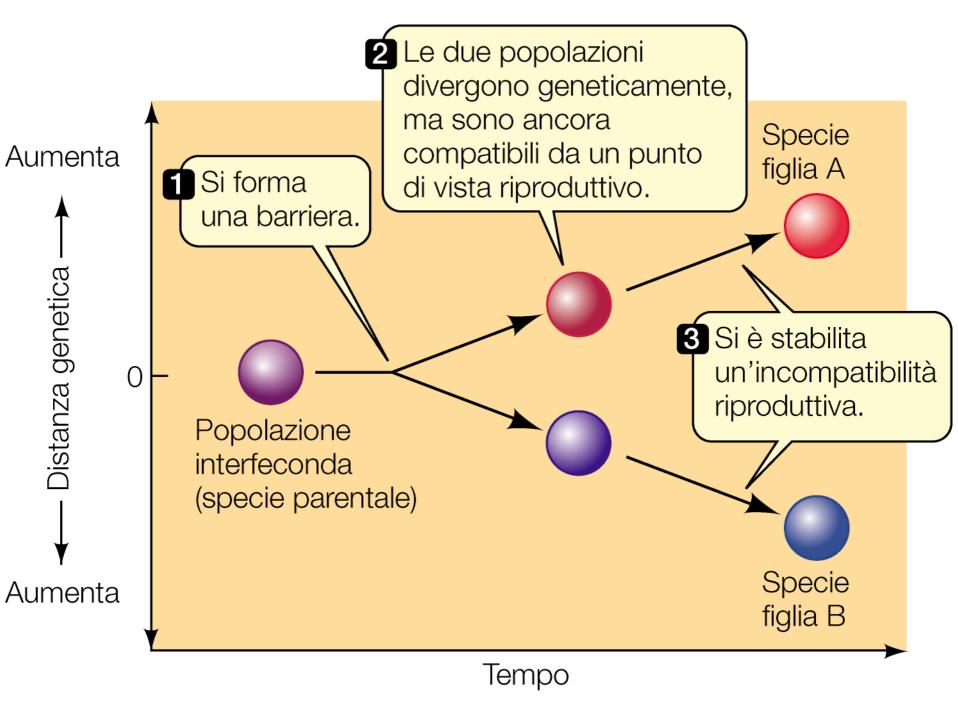
Tutte le specie viventi sono legate le une alle altre a vari gradi attraverso progenitori comuni

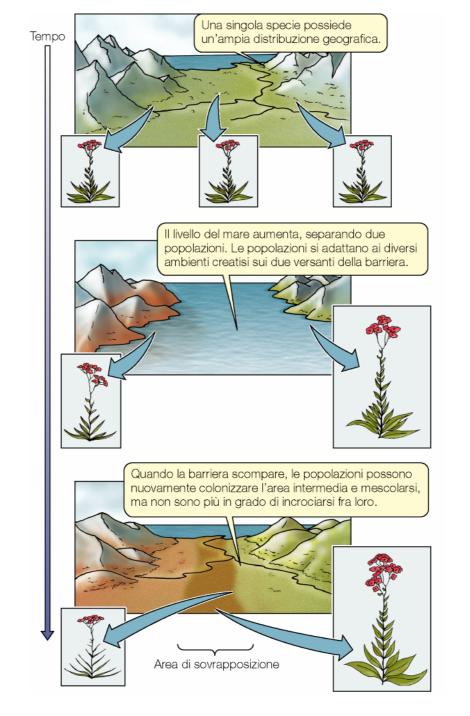
Tutte le forme di vita sulla terra hanno una origine comune. E' esistita una forma di vita originale che ha dato luogo a tutte le forme successive (*LUCA: Last Universal Common Ancestor*)

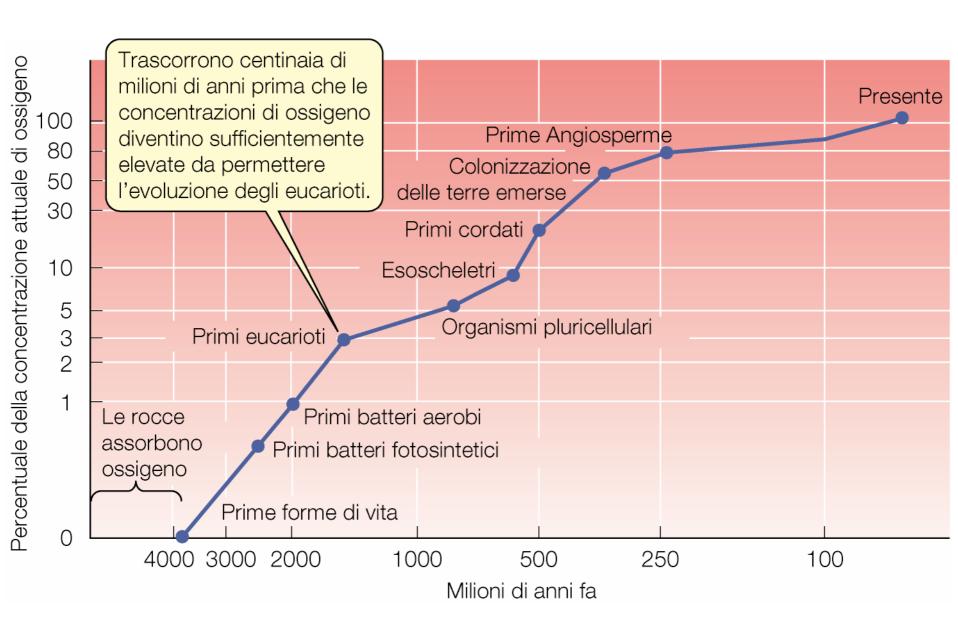
Il processo attraverso cui una specie evolve in un'altra coinvolge mutazioni casuali, le mutazioni che risultano in vantaggio di sopravvivenza si diffondono e persistono piu' di quelle neutre o svantaggiose





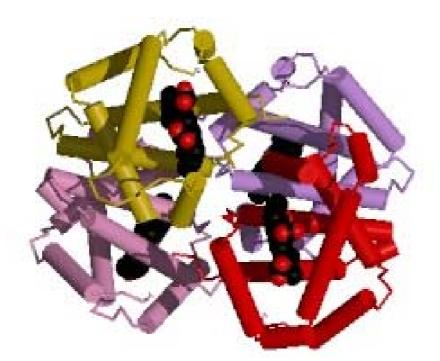






Un esempio di evoluzione a livello molecolare:

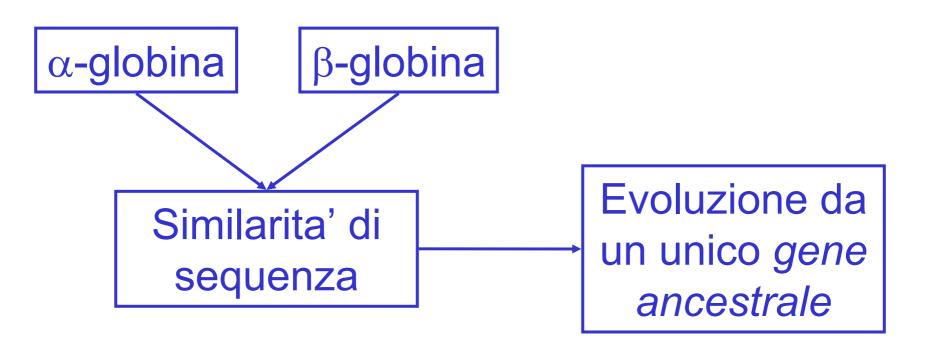
la famiglia delle globine



L'emoglobina e' composta da 4 catene polipeptidiche, chiamate globine, necessarie al trasporto delle molecole di O₂ a due a due uguali.

Un esempio: la famiglia delle globine

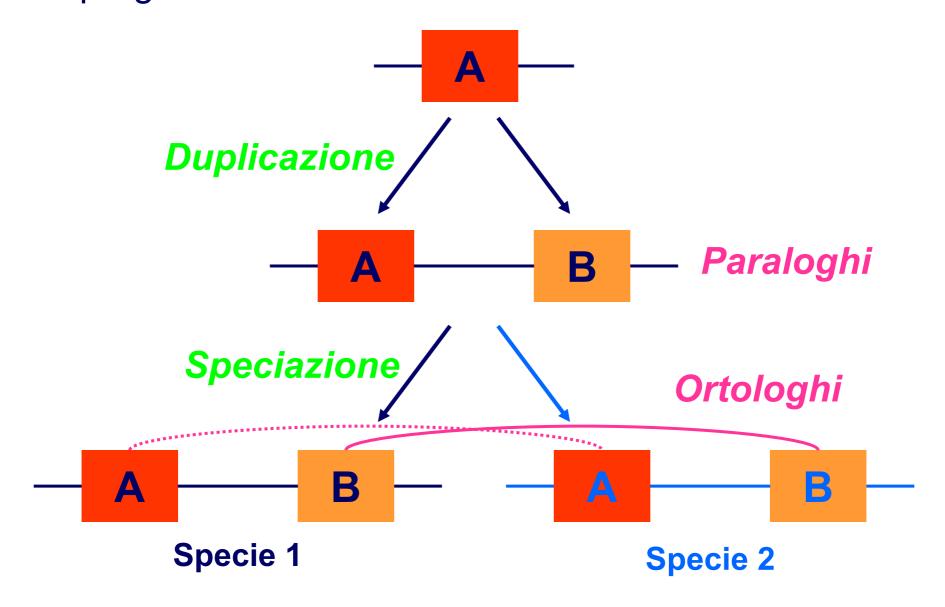
Le globine emoglobiniche possono essere di tipo α o di tipo β .



Similarita' e omologia

- Due o piu' proteine possono essere definite omologhe se derivano da un progenitore comune
- L'omologia tra proteine non puo' essere direttamente osservata ma si deduce dalla loro similarita' in sequenza o funzione.
- Due sequenze sono simili se possono essere allineate in modo che molti ammino acidi corrispondenti sono identici o simili.

Due proteine si definiscono omologhe se derivano da un progenitore comune:

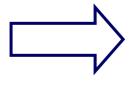


Mutazione genotipica neutrale o deleteria per il fenotipo



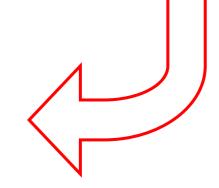
Selezione negativa

Mutazione genotipica vantaggiosa per il fenotipo



Selezione positiva

NUOVA FUNZIONE BIOLOGICA



Cosa causa le mutazioni genetiche?

Mutazioni spontanee

Danneggiamenti nel DNA avvengono continuamente

Errori "editoriali" durante la replicazione

Mutazioni chimiche

Agenti chimici dall'ambiente (mutagenesi, carcinogenesi)

Modificazione di basi

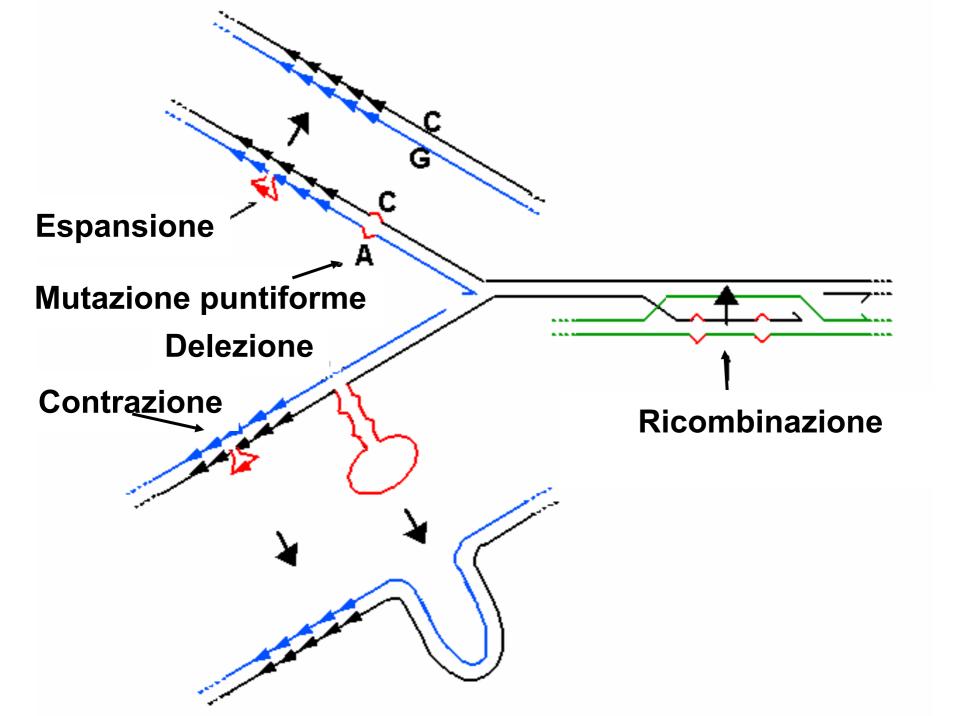
Inserzioni tra basi

UV, radiazioni ionizzanti

Cross-linking tra coppie di basi

Aperture di anelli

Rottura di filamenti di DNA



con chi vai nel bus

von chi vai nel bus

coc hiv ain elb us

Classificazione delle sostituzioni nucleotidiche

Transizione: sostituzione di una purina con una purina (A ⇔ G) o di una pirimidina con una pirimidina (C ⇔ T)

Trasversione: sostituzione di una purina con una pirimidina (A o G \Leftrightarrow C o T) o viceversa C o T \Leftrightarrow A o G)

Sinonima: non si verifica variazione dell'amminoacido codificato (TCT ⇔ TCC, entrambe codificano la L)

Non-sinonima: si verifica variazione dell'amminoacido codificato (CTT ⇔ CCT, L ⇔ P)

Errori di copia di materiale genetico (DNA/RNA)

Sistema	Tasso d'errore stimato (Mut/N _(Pos))	_
Reazione chimica	0.05-0.1 (5-10/100)	
RNA virus (influenza, HIV)	10-3-10-5	RNA-polimerasi
Procarioti (<i>E. coli</i>)	10-10-11	DNA- polimerasi &
Eucarioti (<i>H. sapiens</i>)	3*10⁻8 →	meccanismi riparazione

La *variabilta' intraspecie* e' responsabile della *sopravvivenza* della specie stessa.

Maggiore e' la variabilita' intraspecie, maggiore e' la probabilita' che si verifichino mutazioni positive e che la specie riesca ad adattarsi a nuove condizioni ambientali e quindi a sopravvivere piu' a lungo.



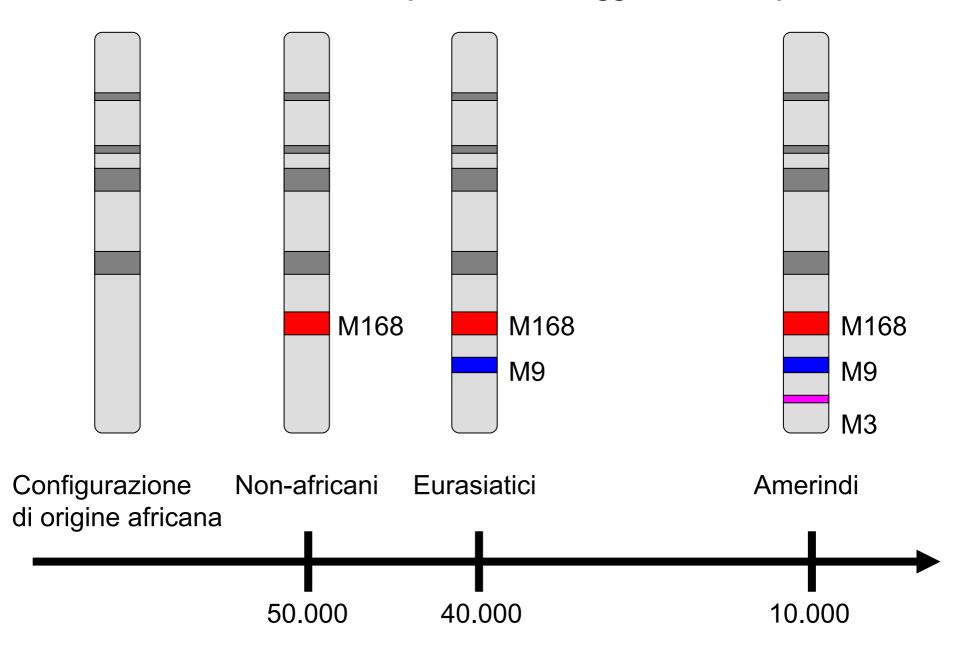
Gli RNA-virus (influenza, HIV, polio) sono tra gli organismi geneticamente piu' variabili e per questo sono difficilmente trattabili farmacologicamente

- I risultati delle variazioni genotipiche forniscono la traccia molecolare dell'evoluzione;
- Quanto piu' due specie sono evolutivamente vicine tanto piu' simili sono le corrispondenti sequenze genomiche ed i loro prodotti di espressione (cioe' le proteine);
- Laddove le sequenze hanno subito variazioni cosi' grandi da non essere piu' riconoscibili, le corrispondenti strutture 3D proteiche possono aver conservato una similarita' informativa.



Relazioni filogenetiche tra organismi

Il cromosoma Y permette di viaggiare nel tempo





Variabilita' inter-specie uomo/scimpanze'

= ~1 %

Variabilita' intraspecie scimpanze' = ~0,4 %

Variabilita' intraspecie uomo = ~0,1 %

Popolazioni geneticamente isolate in Europa

lapponi islandesi finlandesi gallesi baschi



ggcataagaatgatacaatggactttggggacctgagaggaaaggtgggagggggcaagggatactgctca ggtgataggtgcaccaaaatctcacaaatcatcactaaagaacttactcatgtaaccaaatactacctgta ccactataacctacgggggaaaaaagcaacataaccatgaaccaactaataaaaaacaaccttgccttcag tctgcatcctaccctagagacactctctctgtgtcctcacacttggagctaagcttctgacttttgtctcc agtacacccctgaggatcctctcatcacggccatcagaaacctctgtagaaggtcaaatccagtgggttct gtccttatttttctcctatttactgaatcct ctttgaaatctcctcttaattattatgttc tctcatcataccctgagatccctgcattt tcttcctggaaaagctcatctaacctgc acctatgcttgatgactctcagttctct ctgagaccacccatcatacaaaaatgt ttacatattatttttccttagataacti attc caatagccccacactgaactcagtctc ttctctcagtcaggctgtcttctctcattacccttt tggaatcaagatgtttgcattgggttg gggagatgttggt a aggatacatccatttcatt gatacatttcaaa g tacatttcattt agattggaggaataa ttt a g tta tt agt togith that atgtattgttgaa aaaataagtatgtgaggtgagccataagttct aattgctaaaagggtggattttaagtgttctc? ttagcttgatgtagccggtccatgatgtacata baaaacaacatattatacatgataaatataaat aatttttgtcaatcaaaataatttagaaaagt == pacttacacacacacacaaaagagatgattg cattggccagtctaggaataagagttatctgg tctaagtcggatgccaccgacatcactcaccaa taatccctttaatgtcaatcaaattaagtcct catcattttactcctatgcccatttcctcact ctttgttcaggcactattagtcttgcctcttgauccauctttcttcactcatgctgcccactgttgccgta gtgatcttcctaaattgcaaatgcgccatcactctcctgcttaaaatccttcaatgattccttatgacttc tcctttttttttttttttcttgctgctgttccacatccaaagctggctccattcatactgaagcagctgaagttcttcag atatgtcattgccacactgggcccacacttttgaacctgcttcctcctgtgtgagaagtggcttctgccct gttttcggactgcctacattgaagccatctgttccccaggaagccttccctgatgccttgacagcagcatc ttgtgcctgccccatatctgcacttatccatctgggcctgctgttgtcttgtcacttgtgttctcttctgt gaactgtaaacatcaggaggacaagacctatgtcttacttttatttgaatatttagcatctaacaatgttc gacatatagtaggcttttgatactatttttttactatgacattgtagtatatgttaatatccagtaggaca

Metodi per la ricerca di geni

Obiettivo principale della ricerca di geni e' riuscire a distinguere nel genoma tra DNA-genico e DNAnon genico e a localizzare i geni.

- 1. Lunghezza degli ORFs (Open Reading Frames)
- 2. Ricerca di *segnali* nel gene (es. promotori)
- 3. Differenze nella composizione nucleotidica
- 4. Ricerca di *omologia* (similarita' di sequenza)

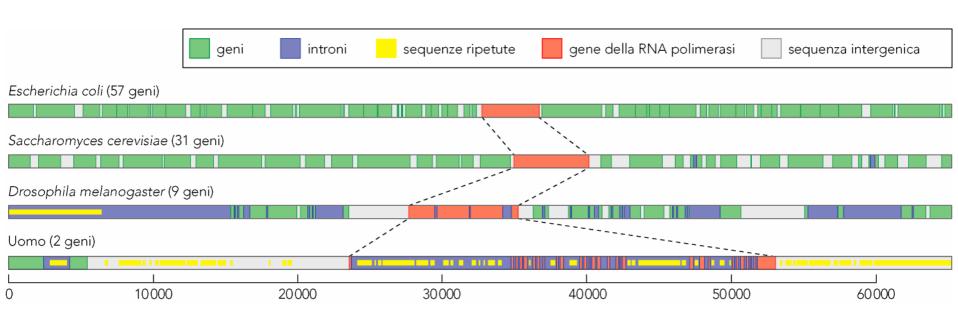
Complessita' dei geni

• Geni codificanti per 100Kb

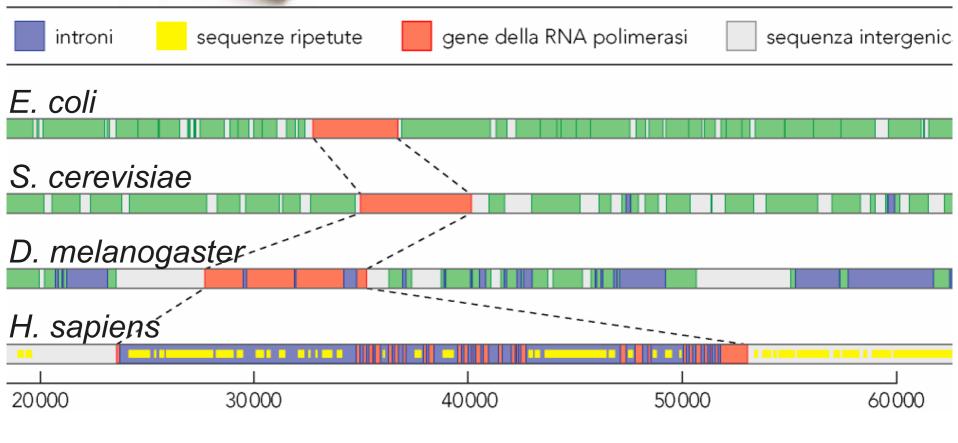
E. coli	87
S. cerevisiae	52
C. elegans	22
 H. sapiens 	5

•	N medio di introni		gene Kb	mRNA Kb
	S. cerevisiae	0	1.6	1.6
	C. elegans	3	4.0	3.0
	Drosophila	3	11.3	2.7
	Gallus	8	13.9	2.4
	 H. sapiens 	6	16.6	2.2

Guardiamo in dettaglio una regione di ~ 65 Kb contenente il gene della RNA polimerasi in 4 diversi genomi: *E. coli, S. cerevisiae, D. melanogaster e H. sapiens*



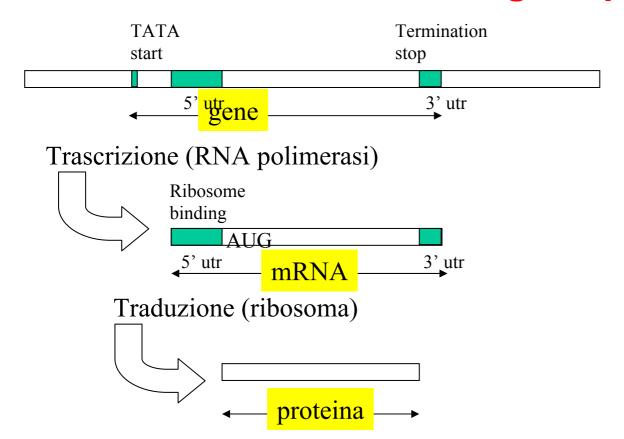
Guardiamo in dettaglio una regione di ~ 65 Kb contenente il gene della RNA polimerasi in 4 diversi genomi: *E. coli, S. cerevisiae, D. melanogaster e H. sapiens*



Genomi procariotici

- Piccole dimensioni (< 10 Mb)
- Alta densita' genica (>85% del genoma e' coding)
- Scarsa ridondanza nucleotidica
- · Introni (quasi) assenti
- Scarsa complessita' dell'architettura genica e genomica
- Facilmente sequenziabili

Struttura di un gene procariota



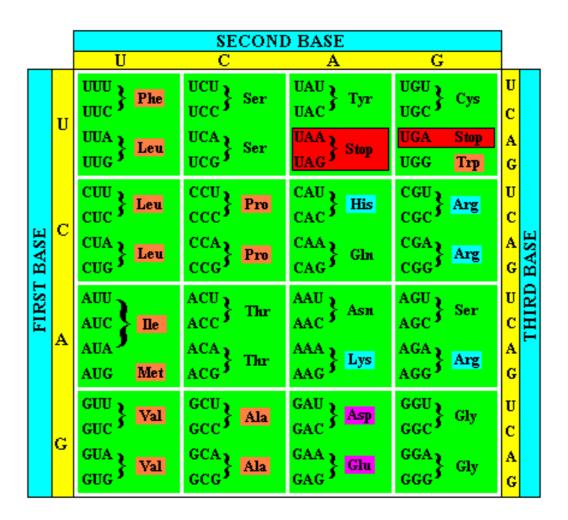
ORFs lunghi nei procarioti

Se assumiamo che i codoni sono distribuiti uniformemente ci aspettiamo uno stop codon ogni 3/64 codoni.

Le proteine sono in media lunghe circa 1000 bp e ogni regione codificante deve contenere solo uno stop codon.

La probabilità casuale di trovare un ORF lungo N codoni preceduto da un ATG e seguito da un codone di terminazione è:

$$P=(1/64) \times (61/64)^{N} \times (3/64)$$



ORFs lunghi nei procarioti

Se assumiamo che i codoni sono distribuiti uniformemente ci aspettiamo uno stop codon ogni 3/64 codoni.

Le proteine sono in media lunghe circa 1000 bp e ogni regione codificante deve contenere solo uno stop codon.

La <u>probabilità casuale</u> di trovare un ORF lungo N codoni preceduto da un ATG e seguito da un codone di terminazione è:

$$P=(1/64) \times (61/64)^{N} \times (3/64)$$

Per N=10

$$P=(1/64) \times (61/64)^{N} \times (3/64) = 0.045$$

Per N=100

$$P=(1/64) \times (61/64)^{N} \times (3/64) = 0.0006$$

Genomi eucariotici

- Dimensioni piu' grandi: da 13 Mb per i funghi piu' semplici a 10.000 Mb per alcune piante superiori
- Architettura genica e genomica complessa in proporzioni alla complessita' dell'organismo. DNA coding: 70% S. cerevisiae, 25% D. melanogaster, 1-3% vertebrati e piante superiori
- Forte presenza di introni
- Forte ridondanza nucleotidica

Geni eucariotici

I geni eucarioti sono piu' complessi:

```
Promotore - 5'UTR - Esone - Introne - Esone - ... 3'UTR - PolyA

Sito di poli-
adenilazione
```

Il segnale per un introne e' di sole 2 bp (AG....GU)

In media un gene di vertebrati e' ~ 30Kb, la regione codificante solo 1-3Kb. In media contiene 6 esoni (~ 150 bp ognuno)

La regione 5'UTR e' lunga ~ 750bp

La regione 3'UTR e' lunga ~450bp

Metodi basati su similarita'

- Confronto con il database di EST (cDNA).
- Confronto della sequenza genomica tradotta con i database di proteine.
- Confronto della sequenza genomica con sequenze genomiche omologhe
- Confronto della sequenza genomica con sequenze di cDNA
- Spliced alignment.

AAAAA AAAAA AAAAA AAAAA AAAAA AAAAA

BANCA DATI EST

(Expressed Sequence Tag)

Seleziona e ottieni la sequenza di tutti gli mRNA contenenti polyA

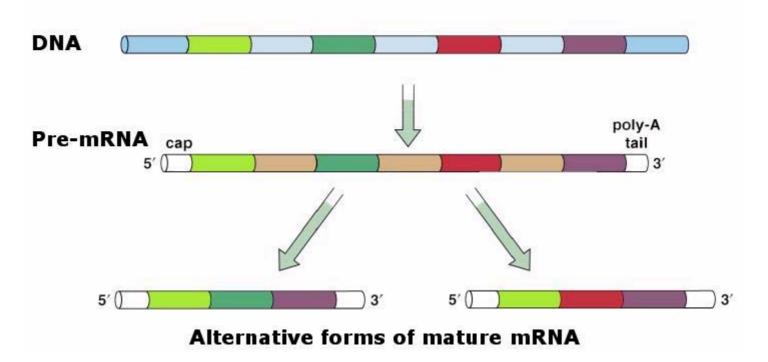
Spliced Alignment

Dati una sequenza genomica ed un insieme di esoni possibili (cioe' tutti i possibili frammenti fra i siti donatori e accettori di splicing):

- costruzione di tutte le possibili catene di esoni
- traduzione del corrispondente mRNA nella corrispondente sequenza proteica
- confronto con una serie di sequenze proteiche fino a che non si rileva una proteina simile.

Alternative RNA splicing

RNA structure can also be rearranged by alternative splicing.



Differenze di composizione nucleotidica

Le frequenze di ricorrenza dei codoni (cioe' le frequenze con cui i codoni si trovano nelle regioni codificanti) possono costituire un indice di *genicita*' del DNA.

In che rapporto di frequenza si trovano alanine e triptofani nelle sequenze di proteine note?

Quante volte troviamo A o T nella terza posizione di un codone? E quante volte G o C?

Differenze di composizione nucleotidica

Le frequenze di ricorrenza dei codoni (cioe' le frequenze con cui i codoni si trovano nelle regioni codificanti) possono costituire un indice di *genicita*' del DNA.

Esempio 1: L, A e W (leucina, alanina e triptofano) sono codificati rispettivamente da 6, 4 e 1 codoni diversi e quindi, nella traduzione di un DNA casuale questi amminoacidi si dovrebbero trovare nel rapporto 6:4:1. Tuttavia, in una proteina di solito si osserva il rapporto 6:5:1 → DNA-coding contiene i codoni in rapporto 6:5:1

Quante volte troviamo A o T nella terza posizione di un codone? E quante volte G o C?

Differenze di composizione nucleotidica

Le frequenze di ricorrenza dei codoni (cioe' le frequenze con cui i codoni si trovano nelle regioni codificanti) possono costituire un indice di *genicita*' del DNA.

Esempio 1: L, A e W (leucina, alanina e triptofano) sono codificati rispettivamente da 6, 4 e 1 codoni diversi e quindi, nella traduzione di un DNA casuale questi amminoacidi si dovrebbero trovare nel rapporto 6:4:1. Tuttavia, in una proteina di solito si osserva il rapporto 6:5:1 → DNA-coding contiene i codoni in rapporto 6:5:1

Esempio 2: in un gene, A o T sono nella terza posizione di un codone il 90% delle volte, G e C solo il 10% (le proporzioni esatte dipendono dalla specie), mentre in un DNA casuale tale rapporto sara' 50%-50%.

Ricerca di segnali nel gene

I **segnali** genici sono unita' discrete con una sequenza consenso riconoscibile (promotori, siti di poliadenlilazione, siti donatori o accettori di splicing)

Gli algoritmi utilizzati per ricercare segnali genici sono chiamati sensori di segnali (*signal sensors*).

I sensori di segnali possono ricercare il segnale isolatamente o in maniera contestuale (cioe' relativamente ad altri segnali locali)

Ricerca di promotori

Le sequenze codificanti sono fiancheggiate da regioni promotore, che non hanno una sequenza specifica.

Si possono definire sequenze consenso a cui tutti i promotori somigliano.

Per esempio, in *E. Coli* dall'analisi di 263 promotori e' stato derivato il consenso:

Di solito si utilizzano misure statistiche:

	1	2	3	4	5	6		
Α	0.1	0.1	0.1	0.6	0.1	0.5		
Т	0.1 0.7	0.7	0.1	0.1	0.1	0.3		
G	0.1	0.1	0.6	0.1	0.1	0.1		
С	0.1	0.1	0.2	0.2	0.7	0.1		
	Т	т	G	A	C	A		
	Sequenza consenso (consensus)							

Come possiamo utilizzare un consensus derivato da casi noti per predire nuovi casi?

Sequenza consensus per un segnale es. promotori, sito di splicing, ...

	1	2	3	4	5	6
Α	0.30	0.50	0.10	0.45	0.65	0.32
T	0.20	0.15	0.05	0.20	0.05	0.25
С	0.45	0.20	0.80	0.30	0.25	0.40
G	0.05	0.15	0.05	0.05	0.05	0.03

La probabilità che la sequenza GGTCAC AACCCAGCG appartenga alla regione usata per derivare la matrice precedente è il prodotto delle probabilità cioè :

se consideriamo la prima G come corrispondente alla prima posizione

 $0.05 * 0.15 * 0.05 * 0.30 * 0.65 * 0.40 = 2.92 * 10^{-5}$

Sequenza consensus per un segnale es. promotori, sito di splicing, ...

	1	2	3	4	5	6
Α	0.30	0.50	0.10	0.45	0.65	0.32
Т	0.20	0.15	0.05	0.20	0.05	0.25
С	0.45	0.20	0.80	0.30	0.25	0.40
G	0.05	0.15	0.05	0.05	0.05	0.03

La probabilità che la sequenza GGTCACAACCCAGCG appartenga alla regione usata per derivare la matrice precedente è il prodotto delle probabilità cioè :

se consideriamo la prima G come corrispondente alla prima posizione

$$0.05 * 0.15 * 0.05 * 0.30 * 0.65 * 0.40 = 2.92 * 10^{-5}$$

se invece partiamo dalla seconda G:

$$0.05 * 0.15 * 0.80 * 0.45 * 0.25 * 0.32 = 2.59 * 10^{-5}$$
 e così via.

La sequenza AACCCA ha una probabilità di

$$0.30 * 0.50 * 0.80 * 0.30 * 0.25 * 0.32 = 2.88 * 10^{-3}$$

TUTTE da confrontare con la probabilita' casuale: 0.25 * 0.25 * 0.25 * 0.25 * 0.25 * 0.25 * 0.25 * 0.25 * 0.25

<u>Dividiamo ciascun valore nella matrice sito</u> <u>specifica precedente per 0.25 (valore casuale), ne facciamo il log in base 2 e otteniamo:</u>

	1	2	3	4	5	6
Α	0.26	1.00	-1.32	0.85	1.38	0.36
Т	-0.32	-0.74	-2.32	-0.32	-2.32	0.00
С	0.85	-0.32	1.68	0.26	0.00	0.68
G	-2.32	-0.74	-2.32	-2.32	-2.32	-3.06
Caso 1 = -3.06	G	G	Т	С	Α	С
	-2.32	-0.74	-2.32	+0.26	+1.38	+0.68
Caso 2 = -0.17	G	Т	С	Α	С	Α
	-2.32	-0.74	+1.68	+0.85	+0.00	+0.36
Caso 3 = +3.56	Α	Α	С	С	С	Α
	+0.26	+1.00	+1.68	+0.26	+0.00	+0.36

Problema dello 0!!!

PSEUDOCOUNTS

Se
$$F=0.25 \rightarrow F/0.25=1 \rightarrow In(F/0.25)=0$$

Se F<0.25
$$\rightarrow$$
 F/0.25<1 \rightarrow In(F/0.25) < 0 neg

Se F>0.25
$$\rightarrow$$
 F/0.25>1 \rightarrow In(F/0.25) > 0 pos

dove F è la frequenza di un evento, ad es. di osservare un dato nucleotide in una data posizione

e 1/2 e' la probabilita' casuale di osservare quell'evento

Se F=Pcaso
$$\rightarrow$$
 F/ Pcaso =1 \rightarrow In(F/ Pcaso) = 0

cioe' se la frequenza con cui osserviamo qualcosa e' uguale a quella che ci aspettiamo per caso il valore corrispondente in una matrice log-odd e' "0"

Se lanciamo un dado molte volte e otteniamo "6"

1/6 delle volte



Se F< Pcaso \rightarrow F/ Pcaso < 1 \rightarrow In(F/ Pcaso) < 0

cioe' se la frequenza con cui osserviamo qualcosa e' piu' bassa di quella che ci aspettiamo per caso il valore corrispondente in una matrice log-odd e' negativo

Se lanciamo un dado molte volte e otteniamo "6"

1/10 delle volte

Il dado e' truccato La possibilita' di numeri alti e' sfavorita rispetto al caso

Se F> Pcaso \rightarrow F/ Pcaso > 1 \rightarrow In(F/ Pcaso) > 0

cioe' se la frequenza con cui osserviamo qualcosa e' piu' alta di quella che ci aspettiamo per caso il valore corrispondente in una matrice log-odd e' positivo

Se lanciamo un dado molte volte e otteniamo "6"

1/2 delle volte



Il dado e' truccato

La possibilita' di numeri alti e' favorita rispetto al caso

Prima lezione.

1. Introduzione alla bioinformatica

2. Evoluzione e informazione

3. Ricerca di geni in genomi di procarioti ed eucarioti