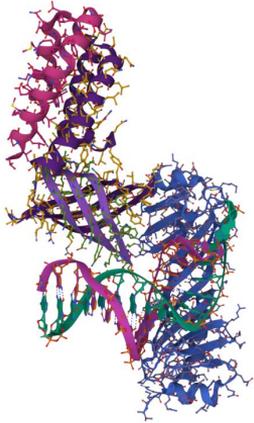


DEF: Approccio computazionale allo studio di sistemi biologici (a livello molecolare) => studio *in silico* (contrapposto allo studio *in vitro* o *in vivo*)

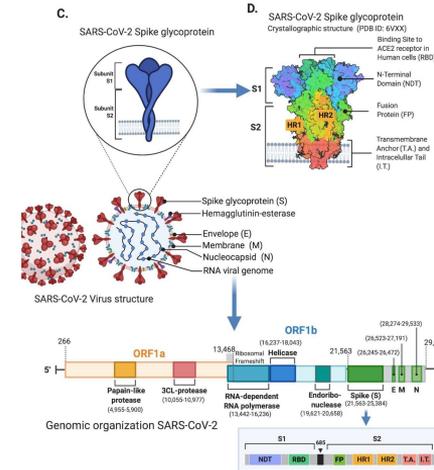
Cosa si può studiare mediante approccio bioinformatico



- Relazioni **evolutive** tra molecole o frammenti di molecole
- Relazioni **strutturali** tra molecole
- Relazione **tra sequenza e struttura**
- Interazione tra molecole
- **Network** di interazione (INTERATTOMA); di regolazione; di metabolica (METABOLOMA)
- STRUTTURA DEI GENOMI; RELAZIONE TRA GENOMI (Comparative Genomics)
- Modificazioni epigenetiche dei genomi (es metilazione: metiloma)

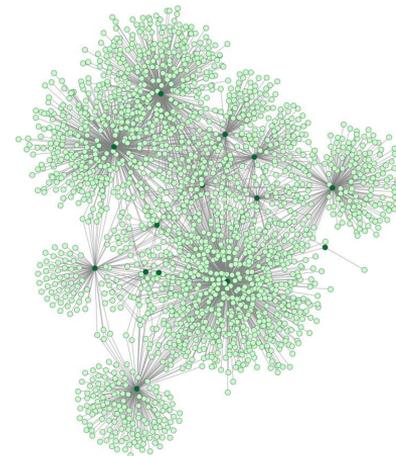
AREE DELLA BIOINFORMATICA

- **GENOMICA (Computational Genomics)** Study of entire genomes. Huge amount of data, fast algorithms, limited to sequence.
- **BIOINFORMATICA STRUTTURALE** Study of the folding process of bio-molecules. Less structural data than sequence data available, step toward function.
- **SYSTEMS BIOLOGY** Study of complex interactions in biological systems. High level of representation.



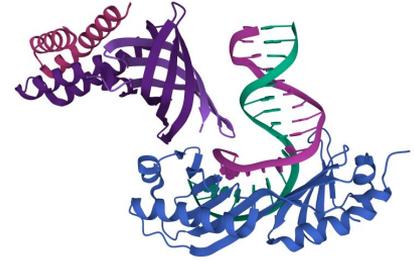
Schematic diagram of the genomic structure of the 29.3 kb(nCoV) gene and domain structure of the 1273aa spike glycoprotein.

<http://dx.doi.org/10.1136/jclinpath-2020-206658>



The protein-protein interaction network obtained from IMEx. Dark green nodes represent SARS-CoV-2 proteins, whereas light green nodes are human proteins.

Computational genomics



La genomica computazionale si focalizza sulla comprensione dei principi secondo cui in DNA regola la biologia di qualsiasi specie a livello molecolare

Al giorno d'oggi abbiamo a disposizione una enorme quantità di dati biologici

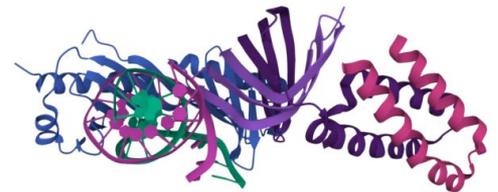
=> gli studi computazionali sono diventati uno strumento determinante nello studio della biologia

L'uso di algoritmi e tecniche di **machine learning** permettono di scoprire segnali biologici in grandi genomi, ricostruire network cellulari e scoprire meccanismi di evoluzione del genoma.

(da: MIT EECS dpt)

RISORSE DEL WEB

- BANCHE DATI (databases)
- GENOME BROWSERS



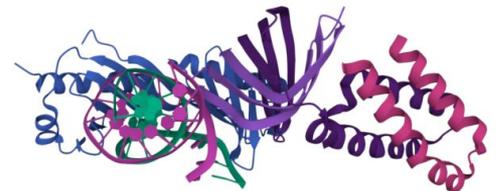
BANCHE DATI

preistoria



In 1955 English biochemist [Frederick Sanger](#)^{en} sequenced the amino acids of insulin, the first of any protein.

Sanger's work “revealed that a protein has a definite constant, genetically determined sequence—and yet a sequence with no general rule for its assembly. Therefore it had to have a code” (Judson, *Eighth Day of Creation*, 188).



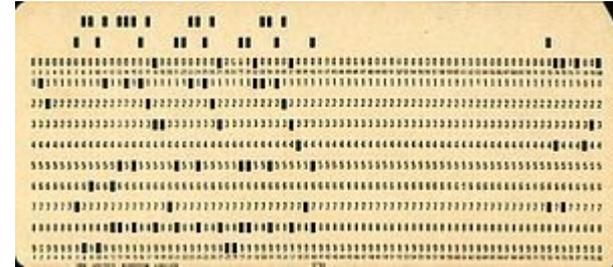
BANCHE DATI



Dr. Margaret Oakley Dayhoff, pioniera della **bioinformatica**, nel **1965** ha pubblicato 'Atlas of protein sequence and structure'. Da questo progetto ha preso vita il **database** Protein Information Resource (PIR) che nel 2002 diventa Uniprot (**Universal Protein Resource**)



Margaret Oakley Dayhoff in the computer room at the National Biomedical Research Foundation, in front of the **punch card** reader.





UniProtKB

Advanced

Search

BLAST Align Retrieve/ID mapping Peptide search SPARQL

Help Contact



The new UniProt website is here! [Take me to UniProt BETA](#)

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.

UniProtKB

UniProt Knowledgebase

Swiss-Prot
(566,996)



Manually annotated and reviewed.

Records with information extracted from literature and curator-evaluated computational analysis.

TrEMBL

UniRef

Sequence clusters



UniParc

Sequence archive



Proteomes

Proteome sets

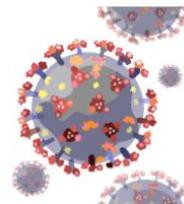


Supporting data

Literature citations

Taxonomy

Subcellular locations



New UniProt portal for the latest SARS-CoV-2 coronavirus protein entries and receptors, updated independent of the general UniProt release cycle.

[View SARS-CoV-2 Proteins and Receptors](#)

News

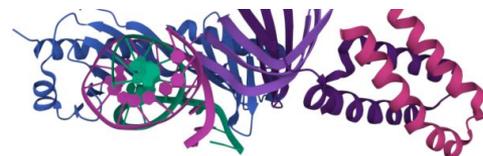


Forthcoming changes

Planned changes for UniProt

UniProt release 2022_01

A phospholipase for clear vision | Cross-references to MANE-Select



GENOME BROWSERS

DEF: GENOME BROWSER: interfaccia grafica per dati genomici di un database biologico.

I genome browser permettono di visualizzare ed analizzare interi genomi.

I dati rappresentati sono 'annotati' ovvero portano annotazioni che permettono di predire le proteine espresse da un gene; la regolazione; la variazione; effettuare analisi comparative e molto altro

Si possono trovare molti genome browsers. Tra questi, FREE ed accessibili online e tra i più utilizzati possiamo citare:

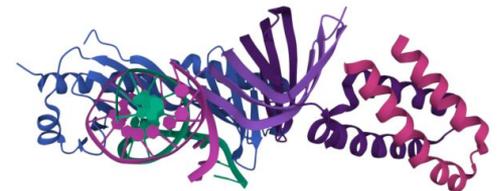
[UCSC Genome Browser](#),

[NCBI's Genome Data Viewer](#).

[Ensembl Genome Browser](#) and



NB: Quelli elencati coprono molti genomi, ma per alcune specie esistono genome browsers dedicati.





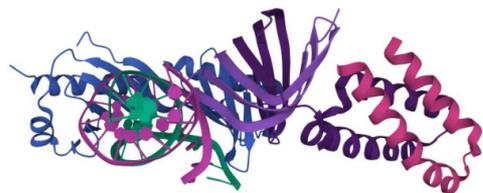
Ensembl hosts several genomic databases and resources for comparative genomics, variation, gene regulation, and epigenetics. Browse through over eighty genomes, from **alpacas** through **zebrafish**. Ensembl also provides tools for identifying phenotypic effects of genetic variants and comparing your data with known genomic variation.



The **UCSC Genome Browser** provides an interactive interface for navigating several sources of genomic information including sequence variation, epigenetic modifications, and transcription factor binding sites. In addition to the genomes of several model organisms, the UCSC Genome Browser also has annotated genomes of many other species, including those of naked mole-rats, the Ebola virus, and over ten species of fruit flies!

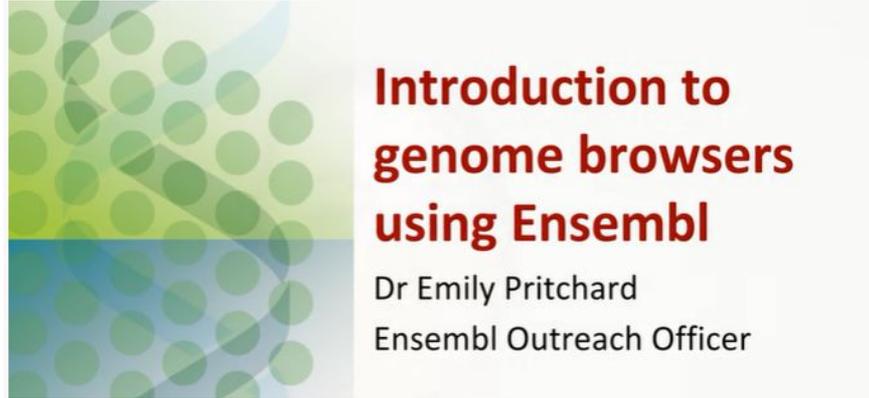


The **Variation Viewer** at NCBI allows you to search the human genome and visualize known genomic variants. The viewer integrates data from several genetic databases, including Gene, dbSNP, and dbVar. The viewer also includes information on medical genetics and gene expression. You can even use the Variation Viewer to compare your own data to what is already known about genetic variation and gene expression.



The Ensembl browser

<https://youtu.be/42qZyXSH0Cc>

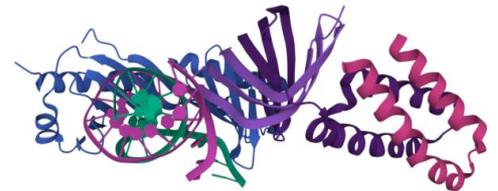


Cos' è un genome browser

Cos'è un genome assembly e perché si aggiorna

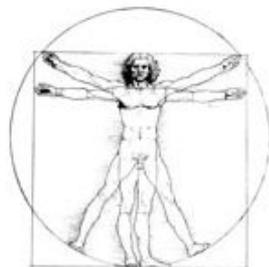
Cosa sono le annotazioni (snp, deletion, insertion...)

Cos'è il regulatory build



Genome Assemblies

The GRC has built tools to facilitate the curation of genome assemblies based on the sequence overlaps of long, high quality sequences (clones and PCR products, not short sequence reads). The GRC currently supports production of assemblies for human, mouse or zebrafish. If your assembly data fits this model and you are interested in using these tools, please [contact us](#). [Subscribe](#) to the grc-announce email list to receive email notification for all GRC assembly updates.



Human

The human genome assembly was produced as part of the [Human Genome Project \(HGP\)](#). The previous assembly (NCBI36) was the last one produced by the HGP and was described in 2004 ([PMID: 15496913](#)); this was the starting point for the GRC. The assembly is based largely on assembling overlapping clone sequences.

Human assembly information

Current major assembly	GRCh38
Regions with alternate loci	178
Assembly N50	67,794,873 bp
Remaining gaps	875
Patch release version	p14
Patches released	FIX: 164 , NOVEL: 90

[More human assembly statistics...](#)

ANNOTAZIONE

Es. gene insulina umana

e!Ensembl BLAST/BLAT | VEP | Tools | BioMart | Downloads | Help & Docs | Blog

Human (GRCh38.p13) ▼

Location: 11:2,159,779-2,161,221 Gene: INS Jobs ▼

Gene-based displays

- Summary
 - Splice variants
 - Transcript comparison
 - Gene alleles
- Sequence
 - Secondary Structure
- Comparative Genomics
 - Genomic alignments
 - Gene tree
 - Gene gain/loss tree
 - Orthologues
 - Paralogues
 - Ensembl protein families
- Ontologies
 - GO: Cellular component
 - GO: Molecular function
 - GO: Biological process
- Phenotypes
- Genetic Variation
 - Variant table
 - Variant image
 - Structural variants
- Gene expression
 - Pathway
 - Regulation
 - External references

Gene: INS ENSG00000254647

Description insulin [Source:HGNC Symbol;Acc:HGNC:6081]

Gene Synonyms IDDM1, IDDM2

Location [Chromosome 11: 2,159,779-2,161,221](#) reverse strand.
GRCh38:CM000673.2

About this gene This gene has 5 transcripts ([splice variants](#)), [247 orthologues](#), [3 paralogues](#) and is associated with [6 phenotypes](#).

Transcripts [Hide transcript table](#)

Show/hide columns (1 hidden)								Filter				
Transcript ID	Name	bp	Protein	Biotype	CCDS	UniProt Match	RefSeq Match	Flags				
ENST00000381330.5	INS-202	465	110aa	P Protein coding	CCDS7729	P01308-1	NM_000207.3	MANE Select v0.95	Ensembl Canonical	GENCODE basic	APPRIS P1	TSL:1
ENST00000397262.5	INS-203	639	110aa	M Protein coding	CCDS7729	P01308-1	-		GENCODE basic	APPRIS P1	TSL:1	
ENST00000250971.7	INS-201	503	110aa	P Protein coding	CCDS7729	P01308-1	-		GENCODE basic	APPRIS P1	TSL:1	
ENST00000421783.1	INS-204	464	92aa	P Protein coding	-	C9JNR5	-		TSL:2	CDS 3' incomplete		
ENST00000512523.1	INS-205	297	98aa	P Protein coding	-	A6XGL2	-		GENCODE basic	TSL:1		

Summary

TRANSCRIPT NAMES AND COLOURS (protein coding)

- A **red** transcript comes from either the Ensembl automatic annotation pipeline or manual curation by the VEGA/Havana project.
- A **gold**, or **merged**, transcript is identical between Ensembl automated annotation and VEGA/Havana manual curation. Only human, mouse, and zebrafish will have gold transcripts. This transcript can be thought of as stable (unlikely to change), and is coloured gold. It is assigned a number beginning with 0.
- A **blue**, **pink** or **grey** transcript is non-coding. See the 'NON-CODING TRANSCRIPTS' section below for more.

Codons Alternating codons Alternating codons

Exons An exon Another exon

Variants 3 prime UTR 5 prime UTR Coding sequence Frameshift Inframe deletion Missense Start lost Stop gained

Synonymous

Other UTR

Markup loaded

• Variants are filtered by consequence type

```

121 YR R NR M RK MYR B K M *V RRR M M
   A C T T G G C T T G G C T C G T G A G C A T C T G G G G T G A G C C C A G G G G C C C A A G G C A G G C A C C T
180

181 KHY ***S*Y* N HY D Y * K YYYYY* B DR Y
   G G C T T C A C C C T G C C T C A G C G T G C G F G T C C C C A G A T C A C T G T G C C T T C G C C T G G C T C
240
   . . . . . A T G G C C C
   . . . . . - M - A - -
2

241 * RRR NY YSBY B Y HYYKY Y S S Y SYVR YHD
   T G T G A T S A T S G C C T G T G C C C G G T G G C G G T G G C C T G G S A C C T G A C C C C C C G G
300
   8 T G T G A T G C G C C T C T G C C C C T G T G G C C T G C T G G C C C T T G S G A C C T G A C C C A G C G
67
   3 L - W - M - R - L - L - P - L - L - A - L - L - A - L - W - G - P - D - P - A -
22

301 HM VS ** *YV S SR YV RV ** *YV Y K
   C A G C C T T T G T G A A C C A G A C G T G C G G C T A C A C C T G G T G G A A G G T G M C A C C T A G T G T
360
   68 C A G C C T T T G T G A A C C A C A C C T G T G C G G C T C A C A C C T G T G G A A G C T C T C A C T A G T G T
127
   23 A - - A - F - V - N - Q - H - L - L - C - G - S - H - L - V - E - A - L - Y - L - V -
42

361 YR*R**YRSKY*B***B**YRYVMSSRM M YRYVYR Y D H * K
   G C G G G A A C A G G G T T C T T C T A C A C C C A A G C C C G G G A G G C A G A G A C C T G C A G
420
   128 G C G G G A A C A G G G T T C T T C T A C A C C C A A G C C C G G G A G G C A G A G A C C T G C A G
187
   43 C - - G - E - R - G - F - F - Y - T - P - X - T - T - R - R - E - A - E - D - L - Q -
62

481 Y S KYWR NM*V*RYR*RRD*Y D B W SM MYY** YY RR
   T G S G G T G G T G G A G C G C G C G G C C C T G G T G C A G C A G C C T C A G C C C T T G C C T G G
480
   182 T G G G C A G T G G A G C T G G C G G G G G C C C T G G T G C A G C A G C C T G C A G C C C T T G C C T G G
247
   63 V - - G - Q - V - - E - L - L - G - G - G - P - G - A - G - S - L - Q - P - L - A - L -
82

541 R HN RYN KB B K RR ** *V SSD Y S Y
   A G S G G T C C T G C A G A A G C G T G C C A T T T G G A C A A T C G T A C C A C A T C T G C T C C C T T
540
   248 A G S G G T C C C T G C A G A A G C G T G C C A T T G G A A C A A T G C T A C C A C A T C T G C T C C C T C T
307
   83 E - - G - S - L - Q - K - R - G - I - V - E - Q - C - C - T - S - I - C - S - L -
102

601 RY Y*B R***R***M WYRY YR M K Y Y H*VH YRMHD H
   A C A G C T G G A G A T G A C T C A C T A C C A G C C C A G S C A G C C C A G A G C C G C C G C G C C
600
   308 A C C A G C T G G A G A C T A C T G C A C T A G S . . . . .
333
   103 Y - - Q - L - E - N - Y - - C - N - * . . . . .
110

601 Y VYBRN**** RR B B Y R
   C T G T S C A C C G A G A G A T G G A T A A A G C C C C T G A A C C A G C
639

```

snp (single nucleotide polymorphism)

```

.....
121 YR R NR M RK MYR B K M *V RRR M M 180
   A C T T G G C T T G G C T C G T G A G C A T C T G G G G T G A G C C C A G G G G C C C A A G G C A G G C A C C T
.....
181 KHY ***S*Y* N HY D Y * K YYYYY* B DR Y 240
   G G C T T C A C C C T G C C T C A G C G T G C G F G T C C C C A G A T C A C T G T G C C T T C G C C T G G C T C
.....
   . . . . . A T G G C C C
   . . . . . - M - A - -
2

241 * RRR NY YSBY B Y HYYKY Y S S Y SYVR YHD 300
   T G T G A T S A T S G C C T G T G C C C G G T G G C G G T G G C C T G G S A C C T G A C C C C C C G G
.....
   8 T G T G A T G C G C C T C T G C C C C T G T G G C C T G C T G G C C C T T G S G A C C T G A C C C A G C G
67
   3 L - W - M - R - L - L - P - L - L - A - L - L - A - L - W - G - P - D - P - A -
22

301 HM VS ** *YV S SR YV RV ** *YV Y K 360
   C A G C C T T T G T G A A C C A C A C G T G C G G C T A C A C C T G G T G G A A G G T G M C A C C T A G T G T
.....
   68 C A G C C T T T G T G A A C C A C A C C T G T G C G G C T C A C A C C T G T G G A A G C T C T C A C T A G T G T
127
   23 A - - A - F - V - N - Q - H - L - L - C - G - S - H - L - V - E - A - L - Y - L - V -
42

361 YR*R**YRSKY*B***B**YRYVMSSRM M YRYVYR Y D H * K 420
   G C G G G A A C A G G G T T C T T C T A C A C C C A A G C C C G G G A G G C A G A G A C C T G C A G
.....
   128 G C G G G A A C A G G G T T C T T C T A C A C C C A A G C C C G G G A G G C A G A G A C C T G C A G
187
   43 C - - G - E - R - G - F - F - Y - T - P - X - T - T - R - R - E - A - E - D - L - Q -
62

```

rs75124361 SNP

Most severe consequence **start lost** | [See all predicted consequences](#)

Alleles **T/C** | Ancestral: T | Highest population MAF: < 0.01

Change tolerance **CADD: C:23.5 | GERP: 1.70**

Location **Chromosome 11:2160971** (forward strand) | **VCF: 11 2160971 rs75124361 T C**

Evidence status **Ex AC** **gnom AD**

HGVs names **This variant has 29 HGVS names - [Show](#)**

Synonyms **ClinGen Allele Registry [CA5818208](#)** (C)

Original source **Variants (including SNPs and indels) imported from dbSNP (release 154) | [View in dbSNP](#)**

About this variant **This variant overlaps 18 transcripts.**



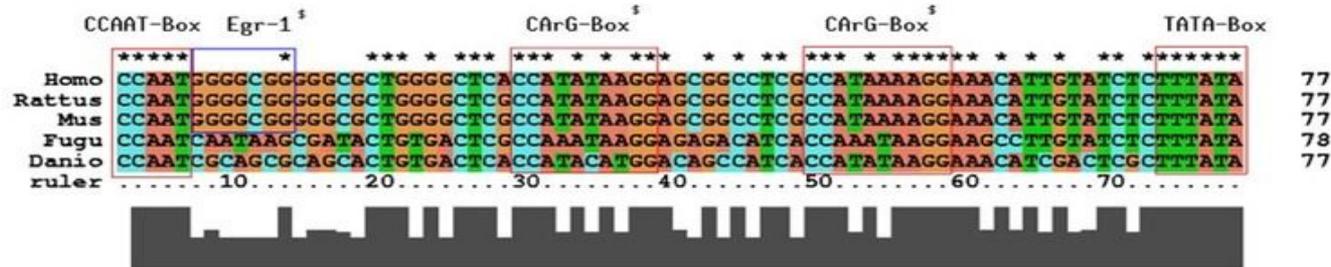
ALIGNMENT



Similarity profiles

Organism	 CHIMP	 MOUSE	 CHICKEN	 FRUIT FLY
Gene Conservation with Humans (%)	99.5	88	75	60

Researchers can learn a great deal about the structure and function of human genes by examining their counterparts in **model organisms**.



ALLINEAMENTO DI SEQUENZE

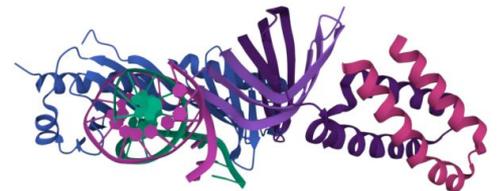
Un allineamento di più o 2 sequenze nucleotidiche o proteiche si basa sull'assunzione che esista una relazione di tipo evolutivo tra queste.

Le tecniche di confronto di sequenze mediante allineamento e gli algoritmi di ricerca nei database biologici sono tool fondamentali della bioinformatica

Il confronto tra sequenze aiuta nella comprensione dell'informazione contenuta in queste e della loro funzione

Gli allineamenti di sequenza vengono usati ad esempio per

- Determinare la funzione di una nuova sequenza genica
- Determinare la relazione evolutiva tra geni, proteine e specie intere
- Predire la struttura e la funzione di nuove proteine (step iniziale dell' Homology modeling)

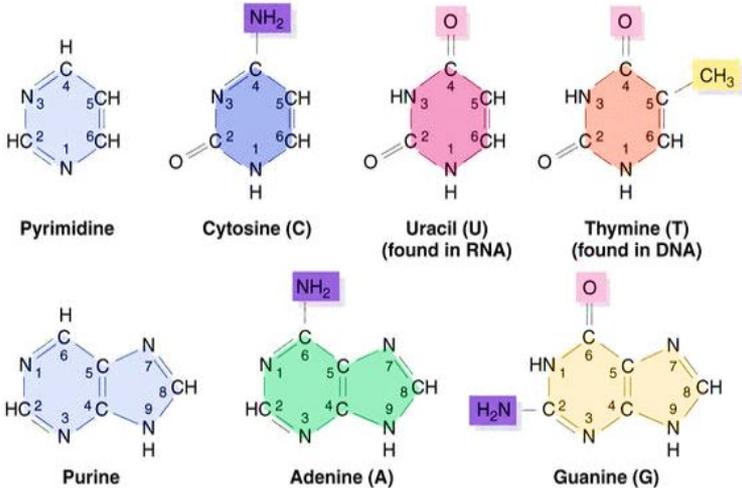


SEQUENZE

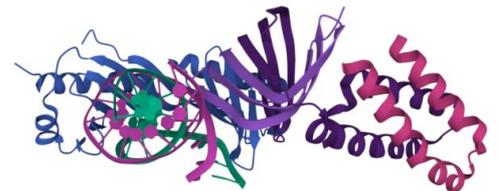
una sequenza nucleotidica o proteica è una successione di caratteri di tipo testo

DNA/RNA: tipicamente 4 caratteri (A-C-G-T/U)

proteine: 20 caratteri (uno per ogni aminoacido: si usa la notazione ad 1 lettera)



A	Alanine
R	Arginine
N	Asparagine
D	Aspartic acid
C	Cysteine
E	Glutamic acid
Q	Glutamine
G	Glycine
H	Histidine
I	Isoleucine
L	Leucine
K	Lysine
M	Methionine
F	Phenylalanine
P	Proline
S	Serine
T	Threonine
W	Tryptophan
Y	Tyrosine
V	Valine



FORMATO FASTA

Il formato più utilizzato per rappresentare una sequenza è il FASTA (suffisso del file: .fas; .fasta, è un file di tipo testo, quindi leggibile da qualsiasi programma di elaborazione testi)

Una sequenza in formato fasta è formata da 2 parti:

una riga di intestazione

una serie di righe di dati di sequenza vera e propria

La riga di intestazione inizia sempre con il carattere > seguito da una descrizione

La sequenza vera e propria inizia nella riga successiva

NB IMPORTANTE IL FONT UTILIZZATO!!!! (vd Courier New)

CARATTERI AMMESSI

A	adenosine	C	cytidine	G	guanine
T	thymidine	N	A/G/C/T (any)	U	uridine
K	G/T (keto)	S	G/C (strong)	Y	T/C (pyrimidine)
M	A/C (amino)	W	A/T (weak)	R	G/A (purine)
B	G/T/C	D	G/A/T	H	A/C/T
V	G/C/A	-	gap of indeterminate length		

A	alanine	P	proline
B	aspartate/asparagine	Q	glutamine
C	cystine	R	arginine
D	aspartate	S	serine
E	glutamate	T	threonine
F	phenylalanine	U	selenocysteine
G	glycine	V	valine
H	histidine	W	tryptophan
I	isoleucine	Y	tyrosine
K	lysine	Z	glutamate/glutamine
L	leucine	X	any
M	methionine	*	translation stop
N	asparagine	-	gap of indeterminate length

NB

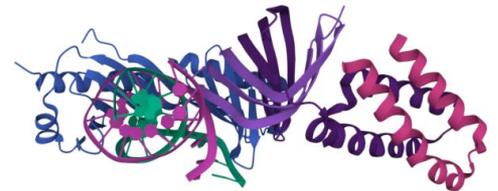
Quando si lavora con le sequenze è IMPORTANTE IL FONT UTILIZZATO!!!!

(ogni carattere testo deve occupare lo stesso spazio per poter lavorare)

Courier A-B-C-D-E-F-G-H-I-L-*

Questrial A-B-C-D-E-F-G-H-I-L-*

Arial A-B-C-D-E-F-G-H-I-L-*



ESEMPI di sequenze

> Sequences (2+)

Sequence status: Complete.

Sequence processing: The displayed sequence is further processed into a mature form.

This entry describes 2 isoforms¹ produced by **alternative splicing**. [Align](#) [Add to basket](#)

This entry has 2 described isoforms and 2 potential isoforms that are computationally mapped. [Show all](#) [Align All](#)

Isoform 1 (Identifier: **P01308-1**) [UniParc] [FASTA](#) [Add to basket](#)

This isoform has been chosen as the canonical¹ sequence. *All positional information in this entry refers to it. This is also the sequence that appears in the downloadable versions of the entry.* [Sequence data in FASTA format](#)

[Hide](#)

```
      10      20      30      40      50
MALWMRLLP LALLALWGPD PAAAFVNHQL CGSHLVEALY LVCGERGFFY
      60      70      80      90     100
TPKTRREAED LQVGQVELGG GPGAGSLQPL ALEGLQKRG IVEQCCTSI
      110
SLYQLENYCN
```

Length: 110

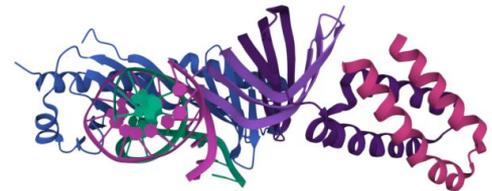
Mass (Da): 11,981

Last modified: July 21, 1986 - v1

Checksum: C2C3B23B85E520E5

```
>sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
MALWMRLLP LALLALWGPDPAAAFVNHQLCGSHLVEALYLVCGERGFFYTPKTRREAED
LQVGQVELGGGPGAGSLQPLALEGLQKRGIVEQCCTSIQSLYQLENYCN
```

```
>
ATGGCCCTGTGGATGCGCCTCCTGCCCTGTGCGCGCTGCTGGCCCTCTGGGGACCTGAC
CCAGCCGCAGCCTTTGTGAACCAACACCTGTGCGGCTCACACCTGGTGAAGCTCTCTAC
CTAGTGTGCGGGGAACGAGGCTTCTTCTACACACCCAAGACCCGCGGGAGGCAGAGGAC
CTGCAGG
```



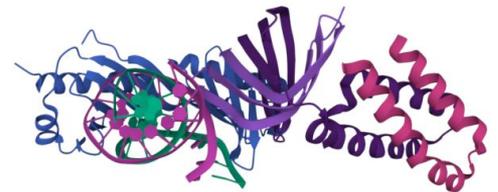
ALLINEAMENTO BINARIO

L'allineamento più semplice è l'**allineamento binario**:

confronto di 2 caratteri appartenenti a 2 sequenze diverse.

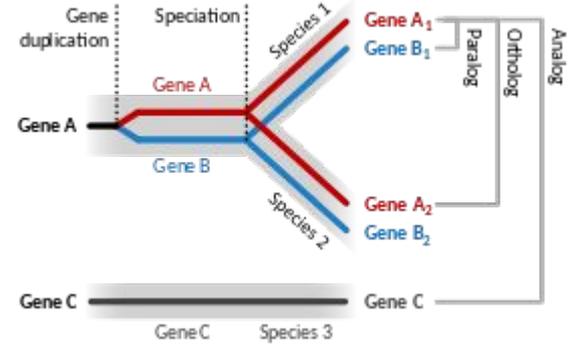
L'allineamento di nucleotidi o aminoacidi riflette la loro relazione evolutiva, ovvero la presenza di un antenato comune.

La similarità viene quantificata mediante un valore razionale (percentuale di omologia)



2 segmenti di DNA possono avere un antenato comune per

- SPECIAZIONE (ortologi)
- DUPLICAZIONE (paraloghi)

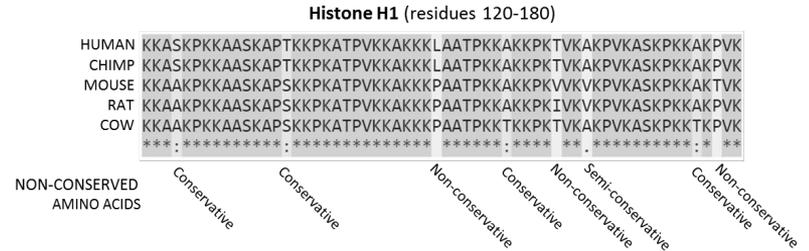


Il grado di similarità tra 2 sequenze viene misurato mediante la percentuale di 'omologia'

PROTEINE: distinguiamo

percentuale di residui identici (percentuale di identità)

percentuale di residui con proprietà fisicochimiche simili (percentuale di similarità)



By Thomas Shafee - Own work, CC BY 4.0,
<https://commons.wikimedia.org/w/index.php?curid=37188728>



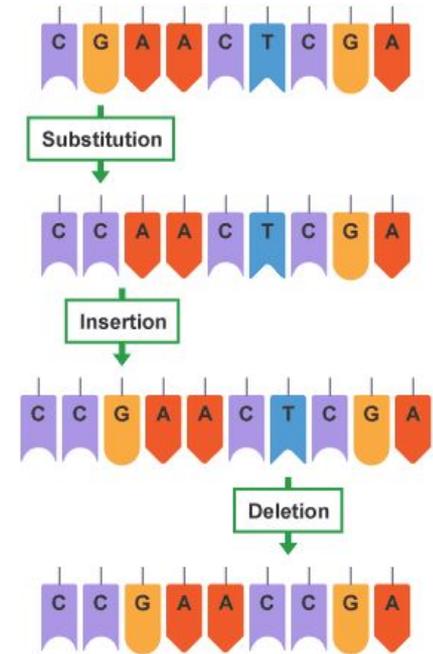
In una data posizione all'interno della sequenza possono trovarsi 3 tipi di variazioni:

una SOSTITUZIONE (di un carattere con un altro)

un'INSERZIONE (aggiunta di 1 o più caratteri)

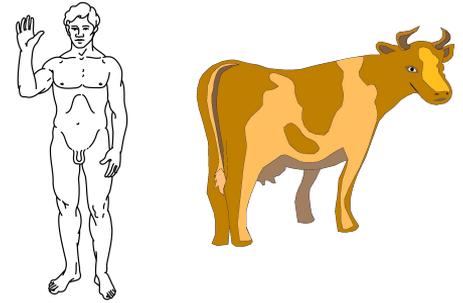
una DELEZIONE (perdita di uno o più caratteri)

In natura, inserzioni e delezioni sono significativamente meno frequenti delle sostituzioni!



ALLINEAMENTO DI SEQUENZE:

ALIGN (es. in UNIPROT): allineamento di 2 o più sequenze proteiche
<https://www.uniprot.org/align/>



ES. allineamento tra:

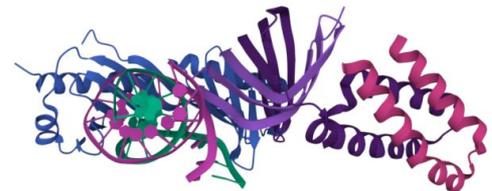
Insulina bovina (insulin bos taurus): uniprot P01317

Insulina umana (insulin homo sapiens: uniprot P01308)

Job status: COMPLETED

```
P01317 INS_BOVIN      1  MALWTRLRPLLALLALWPPPPARAFVFNQHLGSHLVEALYLVCGERGFFYTPKARREVEG      60
P01308 INS_HUMAN     1  MALWMRLLPLLALLALWGPDPAAAFVFNQHLGSHLVEALYLVCGERGFFYTPKTRREAED      60
      ****  *  *****  *  **  *****:****.*

P01317 INS_BOVIN     61  PQVGALELAGGPGAGGL-----EGPPQKRGIVEQCCASVCSLYQLENYCN      105
P01308 INS_HUMAN     61  LQVGGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN      110
      ***  .**.******.*  **  *****:*****
```



ANNOTATION

Highlight

Annotation

- Peptide
- Beta strand
- Disulfide bond
- Signal peptide
- Helix
- Turn
- Propeptide
- Natural variant

Amino acid properties

- Similarity
- Hydrophobic
- Neaative

```
P01317 INS_BOVIN      1 MALWTRLRPLLALLALWPPPPARA FVNQHLCGSHLVEALYLVCGERGFFYTPKARREVEG      60
P01308 INS_HUMAN     1 MALWMRLRPLLALLALWGPDPAAA FVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAD      60
                        ***** * ***** * * *****;***.*.

P01317 INS_BOVIN     61 PQVGALELAGGPGAGGL-----EGPPQKRGIVEQCCASVCSLYQLENYCN      105
P01308 INS_HUMAN     61 LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSI CSLYQLENYCN      110
                        *** :*.*****.* ** *****;*:*****
```

You may add additional sequences to this alignment (in FASTA format)

☰ Add sequence and align

PTM / Processingⁱ

Molecule processing

Feature key	Position(s)	Description
Signal peptide ⁱ	1 – 24	📄 1 Publication ▾
Peptide ⁱ (PRO_0000015819)	25 – 54	Insulin B chain
Propeptide ⁱ (PRO_0000015820)	57 – 87	C peptide
Peptide ⁱ (PRO_0000015821)	90 – 110	Insulin A chain