

Docente: **Giorgio Valentini**

UNIVERSITÀ DEGLI
STUDI DI MILANO



Bioinformatica

A.A. 2016-2017 semestre II

2

Banche dati biologiche

Banche dati biologiche e bioinformatica:

Come mai esistono numerosi strumenti web dedicati in modo specifico alla manipolazione ed elaborazione di dati di tipo biomolecolare e biotecnologico?

Lo sviluppo di questi strumenti è iniziato a metà degli anni ottanta del secolo scorso per rispondere all'esigenza di rendere possibile l'estrazione di informazioni da collezioni di schede testuali che descrivevano molecole presenti nelle cellule di diversi organismi.

Lo sviluppo di questi strumenti di estrazione e di metodi per confrontare biomolecole ha portato alla nascita di una disciplina che prende il nome di bioinformatica.

Bioinformatica è l'**applicazione** di strumenti propri delle scienze dell'informazione (es. algoritmi, intelligenza artificiale, databases) a problemi di interesse biologico, biotecnologico e biomedico.

Banche dati biologiche e bioinformatica:

Attualmente la bioinformatica suscita grande interesse perché:

- La creazione di nuove biotecnologie ha permesso di ridurre il costo ed il tempo necessario per l'acquisizione di informazioni sulle biomolecole. → **esistono banche dati contenenti informazioni riguardanti milioni di biomolecole,**
 - La dimensione dei problemi biologici è sufficiente a motivare lo sviluppo di algoritmi efficienti
 - I problemi sono accessibili (elevata quantità di dati pubblici e letteratura inerente) ed interessanti
 - Le scienze biologiche si avvalgono sempre più spesso di strumenti computazionali
-

Bioinformatica: Scienze dell'Informazione o Biologia? (I)

Gli sviluppi delle scienze biomediche, (in particolare per quanto riguarda la biologia molecolare) si verificano ad un ritmo tale da porre seri problemi:

La nostra capacità tecnologica di acquisire nuovi dati (spesso in quantità elevate) rende impossibile la loro analisi **in assenza di strumenti efficienti.**

Bioinformatica:

Scienze dell'Informazione o Biologia? (II)

Cosa hanno **in comune** scienze biologiche e scienze dell'informazione?

La **biologia**, ed in particolare la biologia molecolare, si occupa dei fenomeni che avvengono nei viventi a livello di atomi e molecole. L'unità di base dei viventi è la **cellula**. La costruzione di una cellula richiede la lettura e la manipolazione di informazioni ...

Scienze dell'informazione è la disciplina che si occupa di calcolo (in senso generico) e delle sue applicazioni. Essa si basa sullo studio sistematico di fattibilità, struttura e automatizzazione di metodi che permettono l'acquisizione, accesso e manipolazione dell'informazione (intesa in senso generico).

Bioinformatica:

Scienze dell'Informazione o Biologia? (III)

Cosa hanno **in comune** scienze biologiche e scienze dell'informazione?

Scienze dell'informazione

Dati di input



Lettura

Elaborazione



Dati di output

Biologia

DNA



copiatura (trascrizione)

traduzione



PROTEINA

Bioinformatica: Scienze dell'Informazione o Biologia? (IV)

Cosa hanno **in comune** scienze biologiche e scienze dell'informazione?

Scienze dell'informazione

Biologia

STUDIANO ENTRAMBE:

- **Flussi di informazione** (biologia nei viventi, scienze dell'informazione all'interno di processi di calcolo)
- **Manipolazione dell'informazione** (biologia nei viventi, scienze dell'informazione all'interno di processi di calcolo)
- **Organizzazione e rappresentazione dell'informazione** (biologia nei viventi, scienze dell'informazione all'interno di processi di calcolo)

L'INFORMAZIONE NEI VIVENTI

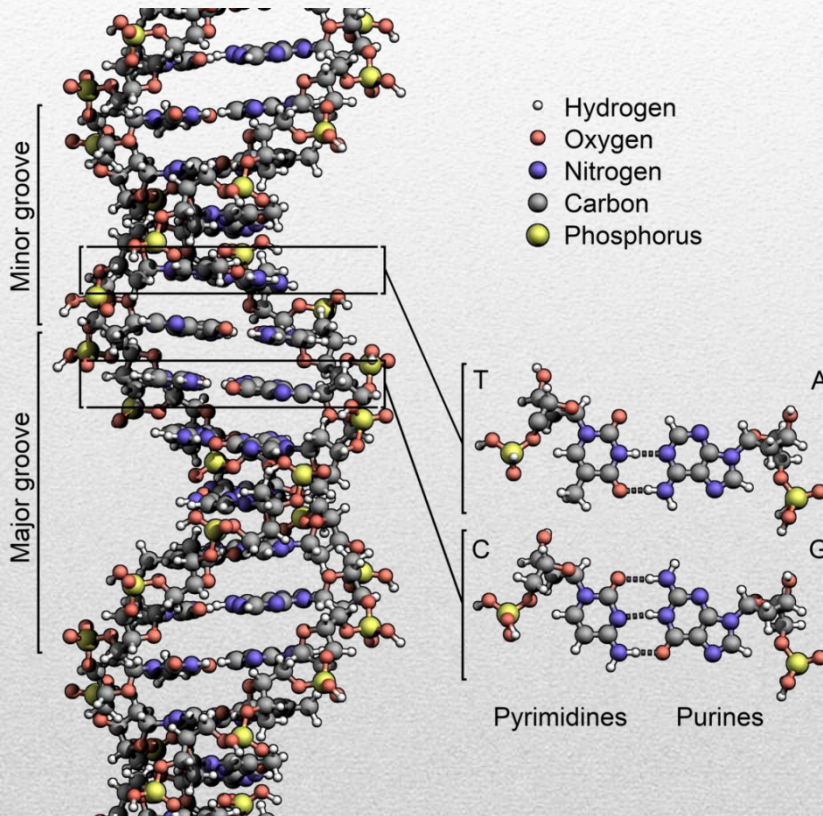
I viventi sono formati da un numero estremamente ampio di molecole. Ognuna di esse è funzionalmente unica e costruita rispettando criteri ben definiti.

Se ogni molecola è costruita seguendo un «progetto» questo implica che, nei viventi, **deve** esistere un luogo il cui ruolo è quello di immagazzinare e rendere disponibili al momento del bisogno le informazioni relative ai «progetti» delle molecole.

A livello biomolecolare come sono realizzate l'**organizzazione** e la manipolazione dell'informazione?

Depositario informazione: DNA

- **Doppia elica**: ognuna delle due eliche può essere ricostruita a partire dall'informazione presente nell'altra (banca dati con sistema di «backup» incorporato)
- **Informazione**: ogni elica è una lunga catena formata da una **sequenza** di 4 elementi : adenina (**A**), Timina (**T**), citosina (**C**) e guanina (**G**). Essi vengono detti **nucleotidi**.

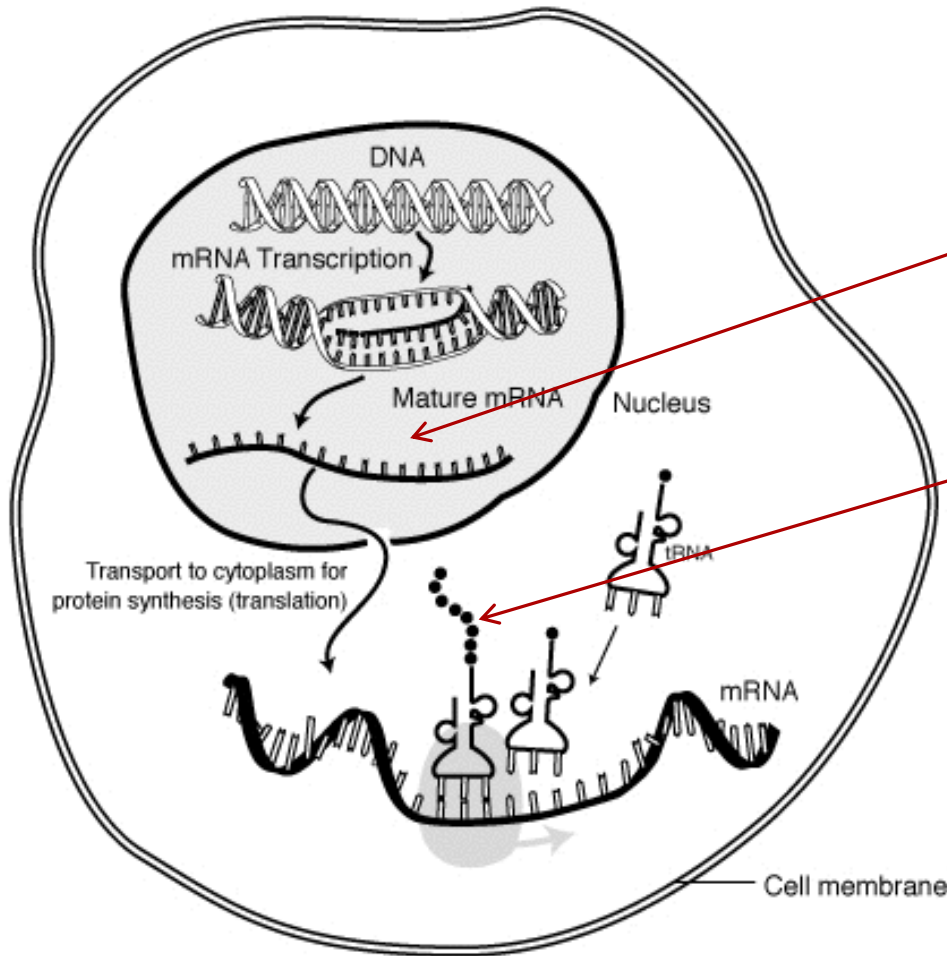


...AGCGGAGGAGCATGCGGATTAGGCTTCGGATCGGAT...

...TCGCCTCCTCGTACGCCTAATCCGAAGCCTAGCCTA...

Lunghezza DNA umano? ... più di 3 miliardi di caratteri.

A livello biomolecolare come sono realizzate l'organizzazione e la **manipolazione** dell'informazione?

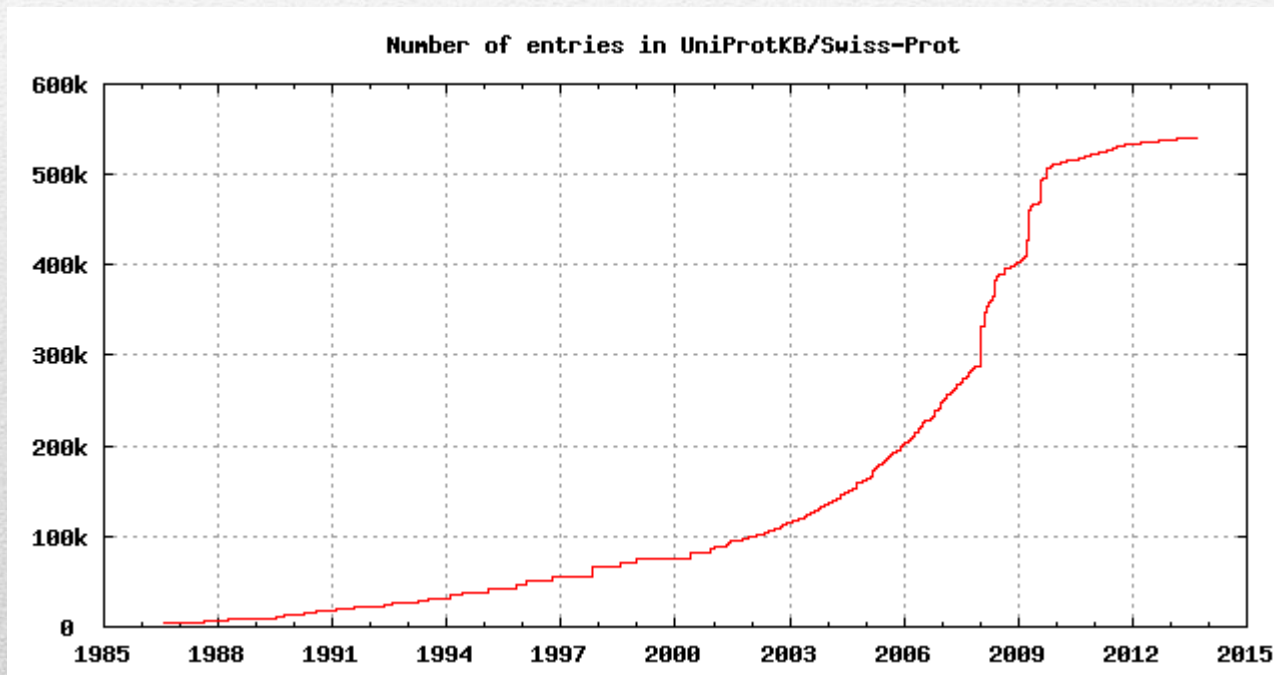


Manipolazione informazione: es. sintesi proteina

- **Copia informazione:** un tratto di DNA (**gene**) viene copiato in una molecola (**RNA**) in grado di trasportare l'informazione altrove.
- **Traduzione:** l'RNA che ha portato l'informazione altrove (**messaggero**) viene tradotto: viene letto a gruppi di tre lettere, ad ogni tripletta corrisponde uno di 20 aminoacidi. mRNA viene letto e viene sintetizzata una catena di aminoacidi (la **proteina** codificata dal messaggero)

Cosa c'entra tutto questo con la lezione di oggi?

Tutte le biosequenze (DNA, geni, mRNA, proteine,...) vengono immagazzinate in apposite banche dati. La banca dati di riferimento per le sequenze di proteine (**Uniprot** : <http://www.uniprot.org/>) contiene, al momento, quasi 600.000 proteine.



Dato il numero estremamente elevato di sequenze non è possibile accedere a queste informazioni senza utilizzare strumenti dedicati. Oggi ci occupiamo di questi strumenti.

BANCHE DATI BIOLOGICHE

Le prime banche dati biologiche sono state create nel 1982. In quel periodo i calcolatori erano poco potenti. La possibilità di scambiarsi informazioni utilizzando internet non era diffusa come al giorno d'oggi.

Tuttavia questa possibilità esisteva già nelle università. In quel periodo serviva molto tempo e molto denaro per ottenere informazioni su una biomolecola.

Ogni volta che il progetto dedicato alla caratterizzazione di una biomolecola terminava, tutti i risultati venivano inseriti in semplici file di testo e inviati ad una banca dati che rendeva queste «schede informative» pubblicamente disponibili.

BANCHE DATI BIOLOGICHE

Quindi le prime banche dati biologiche erano semplici pagine web che fornivano l'accesso a collezioni di file di testo.

Ogni file di testo aveva un nome che corrispondeva ad un codice identificativo della molecola descritta al suo interno. Ad esempio se il file conteneva informazioni su una proteina avente codice **PR001** allora il file di testo contenente la sua scheda aveva nome **PR001.txt**.

La pratica di assegnare un identificativo ad ogni molecola inserita in una banca dati biologica è utilizzata ancora adesso. L'identificativo viene detto **accession number** ed è una parola contenente lettere, numeri e simboli (spazi vuoti non ammessi)

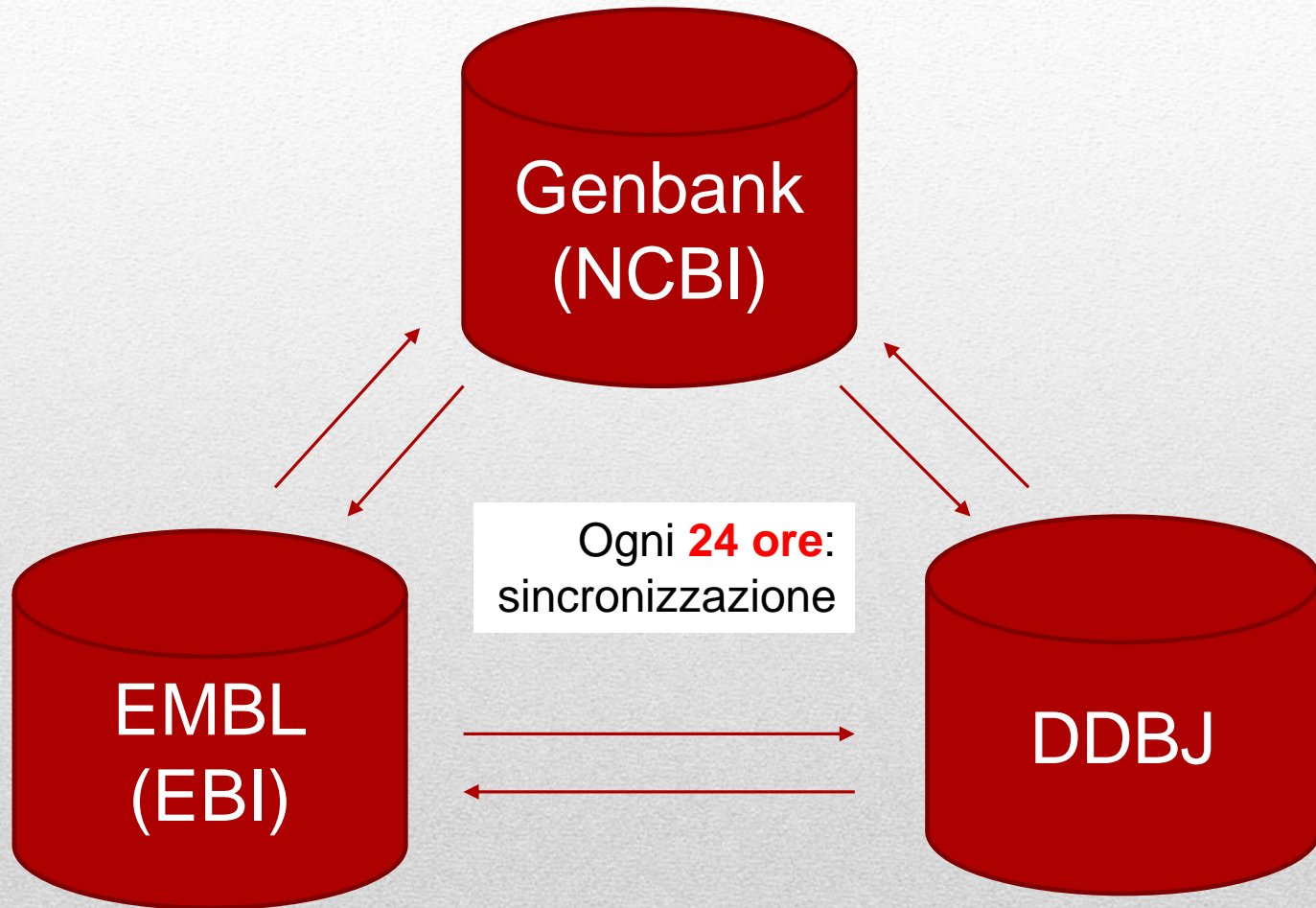
BANCHE DATI BIOLOGICHE

Esistono diversi tipi di banche dati biologiche. Esse si possono classificare in base a vari criteri. Uno di questi riguarda la **qualità delle informazioni** in esse contenute.

Da questo punto di vista possiamo suddividere la banche dati bio in 2 classi principali:

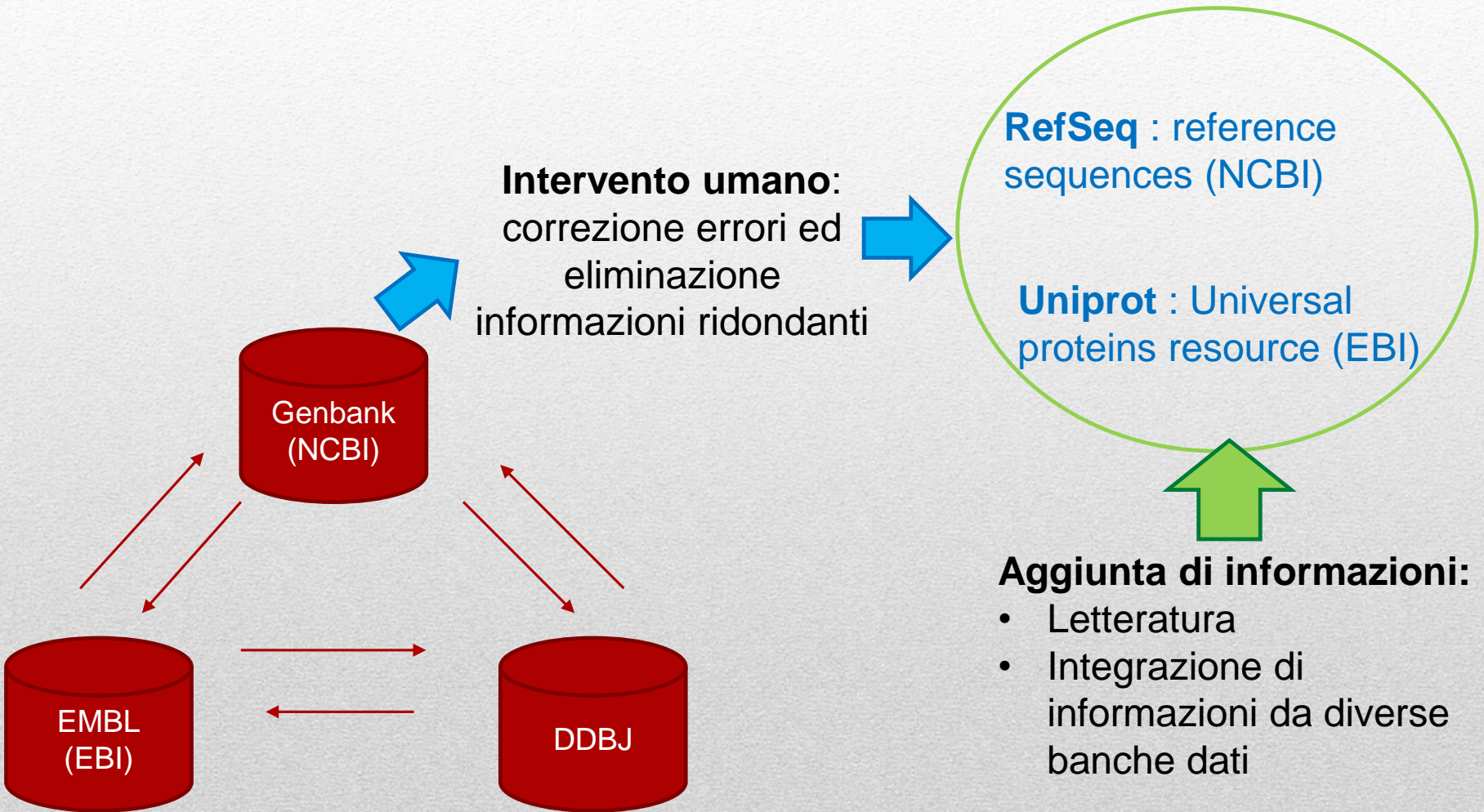
- **Banche dati PRIMARIE** : dette anche collettori primari. Il loro ruolo è quello di raccogliere, giornalmente, tutte le informazioni che riguardano biomolecole prodotte in tutti i laboratori del mondo e renderle disponibili.
 - **Banche dati SECONDARIE**: dato che l'informazione contenuta nei collettori primari è sporca e ridondante esse esaminano i dati dei collettori, correggono eventuali errori, includono informazioni aggiuntive e rendono disponibili i risultati di questo processo di raffinamento.
-

BANCHE DATI PRIMARIE



... va avanti così dagli anni '80 ... secondo voi quanti dati contengono ?

BANCHE DATI SECONDARIE



... va avanti così dagli anni '80 ... secondo voi quanti dati contengono ?

BANCHE DATI PRIMARIE

Queste banche dati contengono, letteralmente, miliardi di schede. Sarebbe impossibile trovare quello di cui abbiamo bisogno in assenza di strumenti che permettano di cercare le informazioni a cui siamo interessati.

Questo ci aiuta a capire il motivo per cui le banche dati biologiche non sono **mai** costituite solamente dalla collezione di dati che contengono ma anche da un insieme di strumenti progettati per rendere possibile estrazione e manipolazione delle informazioni in esse contenute.

... va avanti così dagli anni '80 ... secondo voi quanti dati contengono ?



Banca dati biologica: definizione

Obiettivi:

1. Disseminare dati ed informazioni biologiche
2. Strutturare l'informazione in modo che essa sia leggibile/modificabile da parte di un calcolatore

Una banca dati biologica **DEVE** avere **almeno uno strumento specifico** per la ricerca ed estrazione dei dati.

Pagine web, libri, articoli scientifici, tabelle, file di testo, e fogli di calcolo **non possono** essere considerati banche dati biologiche.

Liste pubbliche di banche dati biologiche

- **Wikipedia (lista di banche dati biologiche)**
https://en.wikipedia.org/wiki/List_of_biological_databases
 - **Nucleic Acids Research Database Listing**
<http://nar.oupjournals.org/cgi/content/full/30/1/1/DC1> (esempio di pubblicazione in cui è presente una lista di database biologici “storica” ... articolo del 2002)
 - Sono un buon punto di partenza per farsi un’idea sul numero e varietà delle banche dati biologiche.
 - **Più di 500** banche dati esistenti sono state catalogate fino ad oggi. E sono costantemente in crescita.
-

Esempio di accesso a banca dati biologica:

Supponiamo di conoscere il “nome” di un gene: **INDY** (ebbene sì ... ogni gene ha un suo nome). E di voler cercare informazioni su di esso in una banca dati.

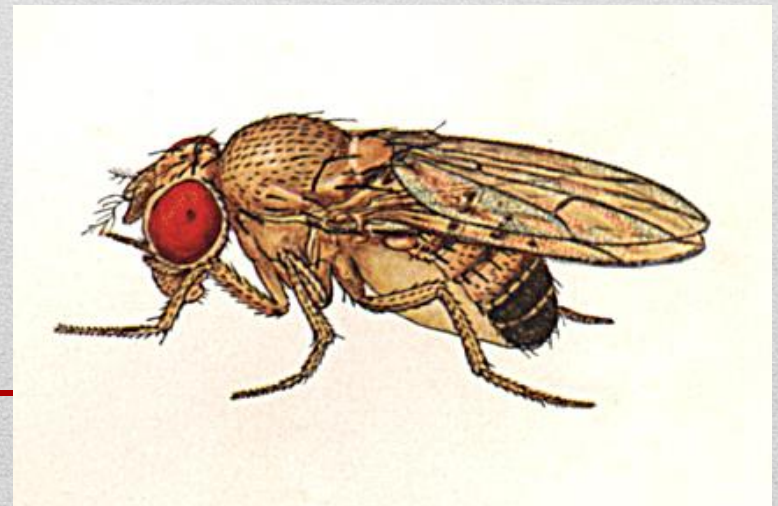
Da dove iniziamo?

Prima di iniziare ...

L'effetto delle mutazioni nei geni viene studiato utilizzando organismi facili da manipolare in laboratorio.

In uno di questi organismi, il moscerino della frutta (nome scientifico: *Drosophila melanogaster*) è stato identificato un gene che, se mutato, raddoppia la **durata della vita media dei moscerini**. A questo gene è stato dato il nome di **INDY**:

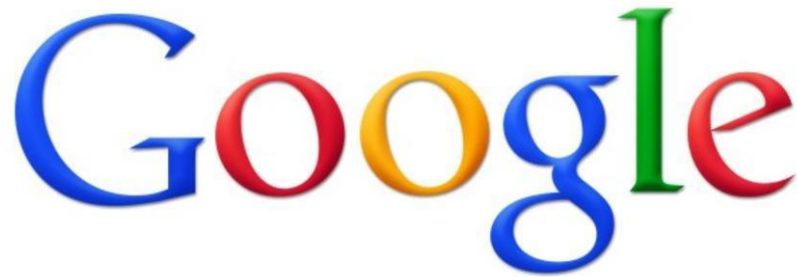
“I’m Not Dead Yet”



Esempio di accesso a banca dati biologica:

Ricerca di informazioni sul gene: **INDY**

Prima prova: usiamo strumenti “classici” per cercare informazioni ...

The image shows the Google logo in its characteristic multi-colored font (blue, red, yellow, blue, green, red) on a white background. The logo is centered within a white rectangular box.

Google

Ricerca di informazioni sul gene: INDY

← → 📶 🔒 https://www.google.it/search?q=INDY&btnG=Cerca&hl=it&gbv=1&site=imghp

+Tu **Ricerca** Immagini Maps Play YouTube News Gmail Altro ▾

Google 🔍

Ricerca Circa 43.200.000 risultati

Web [The Official Site of IndyCar News, Drivers, Schedule & Shop ...](#) 📶
[www.indycar.com/](#) - Copia cache - Simili

Immagini Other Schedules. IZOD IndyCar Series Schedule · Firestone **Indy** Lights Schedule · Pro Mazda Championship Schedule · Cooper Tires USF2000 Schedule.

Maps [Live Timing & Scoring - Schedule - Drivers - Stats](#)

Shopping

Notizie [Indy](#) 📶
[www.indyproject.org/](#) - Copia cache - Simili

Altro An open source internet component suite comprised of popular internet protocols that is included in both Delphi 6 and Kylix. Both client and server ...
[Indy Sockets - Support - CLR - About](#)


Qualsiasi Paese
Pagine da: Italia


Qualsiasi lingua
Pagine in italiano

Qualsiasi data
Ultima ora
Ultime 24 ore
Ultima settimana
Ultimo mese
Ultimo anno

Tutti i risultati
Verbatim

Notizie relative a INDY

 [Houston: formula **Indy**, lo schianto di Franchitti](#)
[La Repubblica](#) - 1 giorno fa
Brutto incidente per il tre volte campione di Indianapolis Dario Franchitti sulla pista di Houston: all'ultimo giro il contatto con Takuma Sato e...

 [Indy car: terribile schianto per Franchitti - VIDEO](#)
[Adnkronos/IGN](#) [AGI - Agenzia Giornalistica Italia](#) - 18 ore fa

[Maserati **Indy** - Wikipedia](#) 📶
[it.wikipedia.org/wiki/Maserati_Indy](#) - Copia cache - Simili
La **Indy** è un modello di autovettura Maserati costruita dal 1969 al 1975. Disegnata da Virgino Vairo, fu presentata dalla Vignale al Salone dell'automobile di ...

[Formula **Indy**, auto fuori pista: feriti 13 spettatori - Video - Corriere TV](#) 📶
[video.corriere.it/.../8d8dfbf6-2f33-11e3-bfe9-e2443a6320c1](#)
22 ore fa ... Formula **Indy**, auto fuori pista: feriti 13 spettatori: Brutto incidente sul circuito di Houston di Indycar. All'ultimo giro della gara di domenica 6 ...

[Indy Week](#) 📶
[www.indyweek.com/](#) - Copia cache - Simili
Progressive news, culture and commentary for Raleigh, Cary, Durham and Chapel Hill, North Carolina.

Non ci siamo ... ci sono troppe cose in internet che si chiamano INDY ... e il gene non è uno dei primi risultati riportati.

Inoltre anche se trovassimo informazioni relative al gene troveremmo molti collegamenti (uno per ogni banca dati che contiene informazioni sul gene ...)

Ricerca di informazioni sul gene: INDY

RICERCA PER PAROLA CHIAVE

Apriamo il web browser e colleghiamoci alla divisione

→ **NUCLEOTIDE** delle banche dati gestite da NCBI:

<http://www.ncbi.nlm.nih.gov/nucleotide>

Scegliendo la “sezione” Nucleotide di Genbank otterremo solo risultati riguardanti molecole composte da nucleotidi (DNA o RNA). Non otterremo risultati riguardanti le schede delle proteine.

Tipo di sequenze contenute: DNA e RNA ... *NON* Proteine

Ricerca di informazioni sul gene: INDY

The image shows a screenshot of the NCBI Nucleotide search interface. At the top, there is a search bar with a dropdown menu set to 'Nucleotide' and a 'Search' button. Below the search bar, there is a red banner with a notice about government funding. The main content area features a dark blue header with the word 'Nucleotide' and a paragraph describing the database. Below this, there are three columns of links: 'Using Nucleotide', 'Nucleotide Tools', and 'Other Resources'. A blue callout box on the left contains the text 'MODALITA' DI RICERCA : RICERCA PER PAROLA CHIAVE'. Two red arrows point from the callout box to the search input field and the 'Search' button.

www.ncbi.nlm.nih.gov/nucleotide

NCBI Resources How To Sign in to NCBI

Nucleotide Nucleotide Search Limits Advanced Help

The information on this web site remains accessible; but, due to the lapse in government funding, the information may not be up to date and we may not be able to respond to inquiries until appropriations are enacted. For updates regarding funding, please see USA.gov.

Nucleotide

The Nucleotide database is a collection of sequences from several sources, including GenBank, RefSeq, TPA and PDB. Genome, gene and transcript sequence data provide the foundation for biomedical research and discovery.

Using Nucleotide

- [Quick Start Guide](#)
- [FAQ](#)
- [Help](#)
- [GenBank FTP](#)
- [RefSeq FTP](#)

Nucleotide Tools

- [Submit to GenBank](#)
- [LinkOut](#)
- [E-Utilities](#)
- [BLAST](#)
- [Batch Entrez](#)

Other Resources

- [GenBank Home](#)
- [RefSeq Home](#)
- [Gene Home](#)
- [SRA Home](#)
- [INSDC](#)

Scrivete **qui** il nome del gene e, poi, premete **search**

gene: INDY, ricerca per **parola chiave** NCBI Nucleotide

NCBI Resources How To Sign in to NCBI

Nucleotide Nucleotide INDY Search Save search Limits Advanced Help

Display Settings: Summary, 20 per page, Sorted by Default order **Send to:** **Filter your results:**

Found 81 nucleotide sequences. Nucleotide (53) EST (28)

Results: 1 to 20 of 53 << First < Prev Page 1 of 3 Next > Last >>

[Drosophila melanogaster INDY transporter protein \(Indy\) mRNA, complete cds](#)
1. 2,602 bp linear mRNA
Accession: AF509505.1 GI: 27127245
[GenBank](#) [FASTA](#) [Graphics](#) [Related](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant A, mRNA](#)
2. 2,581 bp linear mRNA
Accession: NM_001169994.2 GI: 442633232
[GenBank](#) [FASTA](#) [Graphics](#) [Related](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant A, mRNA](#)
3. 2,484 bp linear mRNA
Accession: NM_079426.4 GI: 442633232
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant C, mRNA](#)
4. 2,600 bp linear mRNA

Otteniamo lista di «schede» informative che contengono la parola INDY. Ogni elemento è un link che porta ad una singola scheda (entry)

All (53)
Bacteria (0)
[INSDC \(GenBank\) \(33\)](#)
[mRNA \(15\)](#)
[RefSeq \(20\)](#)
[Manage Filters](#)

▼ Top Organisms [\[Tree\]](#)
Drosophila melanogaster (27)
Mus musculus (6)
Homo sapiens (4)
synthetic construct (4)
Oryctolagus cuniculus (3)
All other taxa (9)
[More...](#)

Find related data
Database:

gene: INDY, ricerca per **parola chiave** NCBI Nucleotide

NCBI Resources How To Sign in to NCBI

Nucleotide Nucleotide INDY Search Save search Limits Advanced Help

Display Settings: Summary, 20 per page, Sorted by Default order **Send to:** **Filter your results:**

Found 81 nucleotide sequences. Nucleotide (53) EST (28)

Results: 1 to 20 of 53 << First < Prev Page 1 of 3 Next > Last >>

[Drosophila melanogaster INDY transporter protein \(Indy\) mRNA, complete cds](#)
1. 2,602 bp linear mRNA
Accession: AF509505.1 GI: 27127245
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant A, mRNA](#)
2. 2,581 bp linear mRNA
Accession: NM_079426.4 GI: 442633232
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant C, mRNA](#)
3. 2,484 bp linear mRNA
Accession: NM_079426.4 GI: 442633232
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant D, mRNA](#)
4. 2,600 bp linear mRNA
Accession: NM_079426.4 GI: 442633232
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

All (53)
Bacteria (0)
[INSDC \(GenBank\) \(33\)](#)
[mRNA \(15\)](#)
[RefSeq \(20\)](#)
[Manage Filters](#)

▼ Top Organisms [Tree](#)
Drosophila melanogaster (27)
Mus musculus (6)
Homo sapiens (4)
synthetic construct (4)
Oryctolagus cuniculus (3)
All other taxa (9)
[More...](#)

Find related data
Database:

Abbiamo ottenuto **53** entries da Nucleotide e **28** entries da EST (espressed sequence tags) collezione di sequenze **PARZIALI** (dati meno affidabili di quelli provenienti da Nucleotide)

gene: INDY, ricerca per **parola chiave** NCBI Nucleotide

NCBI Resources How To Sign in to NCBI

Nucleotide Nucleotide INDY Search Save search Limits Advanced Help

Display Settings: Summary, 20 per page, Sorted by Default order **Send to:** **Filter your results:**

Found 81 nucleotide sequences. Nucleotide (53) EST (28)

Results: 1 to 20 of 53 << First < Prev Page 1 of 3 Next > Last >>

[Drosophila melanogaster INDY transporter protein \(Indy\) mRNA, complete cds](#)
1. 2,602 bp linear mRNA
Accession: AF509505.1 GI: 27127245
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant A, mRNA](#)
3. 2,484 bp linear mRNA
Accession: NM_079426.4 GI: 442633232
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant C, mRNA](#)
4. 2,600 bp linear mRNA

[mRNA](#)

Top Organisms [Tree]

- Drosophila melanogaster (27)
- Mus musculus (6)
- Homo sapiens (4)
- synthetic construct (4)
- Oryctolagus cuniculus (3)
- All other taxa (9)
- More...

Find related data Database:

E' disponibile il conteggio delle sequenze estratte in base all'organismo da cui derivano. La lista può essere molto lunga. Per visualizzare la lista completa fate click su **More...**

gene: INDY, ricerca per parola chiave

- [Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant D, mRNA](#)
2. 2,581 bp linear mRNA
Accession: NM_001169994.2 GI: 442633233
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)
- [Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant A, mRNA](#)
3. 2,484 bp linear mRNA
Accession: NM_079426.4 GI: 442633232
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)
- [Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant C, mRNA](#)
4. 2,600 bp linear mRNA
Accession: NM_168779.2 GI: 442633231
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)
- [Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant B, mRNA](#)
5. 2,572 bp linear mRNA
Accession: NM_168778.2 GI: 442633230
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)
- [Biston betularia I'm not dead yet \(indy\) mRNA, partial cds](#)

- ▼ Top Organisms [Tree]
- Drosophila melanogaster (27)
 - Mus musculus (6)
 - Homo sapiens (4)
 - synthetic construct (4)
 - Oryctolagus cuniculus (3)
 - All other taxa (9)
 - More...

Find related data

Database:

- Select
- Select
- Nucleotide
- Assembly
- BioProject
- BioSample
- BioSystems
- Clone
- dbVar
- Gene
- Genome
- GEO Profiles
- HomoloGene
- EST
- GSS
- OMIM
- PubChem BioAssay
- PubChem Compound
- PubChem Substance
- PMC
- PopSet

Al di sotto della lista dedicata agli organismi di provenienza delle sequenze c'è uno strumento che permette di identificare dati correlati alle sequenze ma **PRESENTI IN ALTRE BANCHE DATI**. Per ora non usiamo questo strumento...

- [Drosophila mauritiana strain G105 I am not dead yet \(Indy\) gene, partial sequence](#)
8. 782 bp linear DNA
Accession: EF388947.1 GI: 126429573
[GenBank](#) [FASTA](#) [Graphics](#) [PopSet](#) [Related Sequences](#)
- [Drosophila sechellia strain S9 I am not dead yet \(Indy\) gene, partial sequence](#)
9. 780 bp linear DNA
Accession: EF388946.1 GI: 126429572

INDY (53)

Nucleotide

See more...

Nucleotide

Nucleotide

INDY

Search

Save search

Limits

Advanced

Help

The information on this web site remains accessible; but, due to the lapse in government funding, the information may not be up to date, and the agency may not be able to respond to inquiries until appropriations are enacted. For updates regarding government operating status see USA.gov.

Display Settings: Summary, 20 per page, Sorted by Default order

Send to: Filter your results:

Found 81 nucleotide sequences. Nucleotide (53) EST (28)

Results: 1 to 20 of 53

<< First < Prev Page 1 of 3 Next > Last >>

Drosophila melanogaster INDY transporter protein (Indy) mRNA, complete cds

1. 2,602 bp linear mRNA

Accession: AF509505.1 GI: 27127245

GenBank FASTA Graphics Related Sequences

All (53)

Bacteria (0)

INSDC (GenBank) (33)

mRNA (15)

RefSeq (20)

Manage Filters

Ora cerchiamo di filtrare i risultati. Vogliamo ottenere solo le sequenze di un certo tipo di molecola: RNA messaggero (mRNA). Fate click su Advanced

3. 2,484 bp linear mRNA

Accession: NM_079426.4 GI: 442633232

GenBank FASTA Graphics Related Sequences

Drosophila melanogaster I'm not dead yet (Indy), transcript variant C, mRNA

4. 2,600 bp linear mRNA

Accession: NM_168779.2 GI: 442633231

GenBank FASTA Graphics Related Sequences

Top Organisms [Tree]

Drosophila melanogaster (27)

Mus musculus (6)

Homo sapiens (4)

synthetic construct (4)

Oryctolagus cuniculus (3)

Drosophila pseudoobscura (3)

Macaca fascicularis (2)

Drosophila pseudoobscura

pseudoobscura (2)

Gallus gallus (1)

Drosophila mauritiana (1)

Drosophila sechellia (1)

Biston betularia (1)

Less...

NCBI Resources How To Sign in to NCBI

Nucleotide Nucleotide INDY Search Help

Advanced

The information on this web site remains accessible, but, due to the lapse in government funding, the information may not be up to date, and the agency may not be able to respond to inquiries until appropriations are enacted. For updates regarding government operating status see USA.gov.

Limits

Published in the last
Any Date

Modified in the last
Any Date

Segmented Sequences
Any

Source database
Any

Molecule
mRNA

Exclude
 STSs
 working draft
 TPA
 patents

Reset Search

1: Selezionate mRNA dalla lista disponibile nella sezione **Molecule**

2: Premete il pulsante Search

The information on this web site remains accessible; but, due to the lapse in government funding, the information may not be up to date, and the agency may not be able to respond to inquiries until appropriations are enacted. For updates regarding government operating status see USA.gov.

Display Settings: Summary, 20 per page, Sorted by Default order

Limits Activated: Molecule: mRNA [Change](#) | [Remove](#)

Questa scritta ci ricorda che in **Limits** ci sono elementi attivi

Results: 15

[Drosophila melanogaster INDY transporter protein \(Indy\) mRNA, complete cds](#)

1. 2,602 bp linear mRNA

Accession: AF509505.1 GI: 27127245

[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant D, mRNA](#)

2. 2,581 bp linear mRNA

Accession: NM_001169994.2 GI: 442633233

[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant A, mRNA](#)

3. 2,484 bp linear mRNA

Accession: NM_079426.4 GI: 442633232

[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant C, mRNA](#)

4. 2,600 bp linear mRNA

Accession: NM_168779.2 GI: 442633231

[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant B, mRNA](#)

5. 2,572 bp linear mRNA

Accession: NM_168778.2 GI: 442633230

[INDY \(GenBank\) \(5\)](#)

[mRNA \(15\)](#)

[RefSeq \(10\)](#)

[Manage Filters](#)

Top Organisms [Tree]

Drosophila melanogaster (6)

Homo sapiens (4)

Macaca fascicularis (2)

Mus musculus (1)

Oryctolagus cuniculus (1)

All other taxa (1)

More...

Il numero di risultati è diminuito. Tutte le sequenze ottenute sono di tipo: mRNA

Find related data

Database:

Select

The information on this web site remains accessible; but, due to the lapse in government funding, the information may not be up to date, and the agency may not be able to respond to inquiries until appropriations are enacted. For updates regarding government operating status see [USA.gov](#).

Display Settings: ☑ Summary, 20 per page, Sorted by Default order

Send to: ☑ **Filter your results:**

⚠ **Limits Activated:** Molecule: mRNA [Change](#) | [Remove](#)

All (15)

Bacteria (0)

[INSDC \(GenBank\) \(5\)](#)

[mRNA \(15\)](#)

[RefSeq \(10\)](#)

[Manage Filters](#)

Results: 15

[Drosophila melanogaster INDY transporter protein \(Indy\) mRNA, complete cds](#)

Nonostante il filtro i risultati ottenuti corrispondono a sequenze appartenenti a più organismi. Siamo interessati solo alle sequenze di **Drosophila melanogaster** (il moscerino della frutta ... un noto organismo utilizzato in laboratorio).

Selezionamo l'organismo da questa lista

▼ **Top Organisms [Tree]**

Drosophila melanogaster (6)

Homo sapiens (4)

Macaca fascicularis (2)

Mus musculus (1)

Oryctolagus cuniculus (1)

All other taxa (1)

[More...](#)

Accession: NM_079426.4 GI: 442633232

[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant C, mRNA](#)

4. 2,600 bp linear mRNA

Accession: NM_168779.2 GI: 442633231

[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant B, mRNA](#)

5. 2,572 bp linear mRNA

Accession: NM_168778.2 GI: 442633230

Analyze these sequences

[Run BLAST](#)

Find related data

Database:

Select

The information on this web site remains accessible; but, due to the lapse in government funding, the information may not be up to date, and the agency may not be able to respond to inquiries until appropriations are enacted. For updates regarding government operating status see [USA.gov](https://www.usa.gov).

[Display Settings:](#) Summary, 20 per page, Sorted by Default order

[Send to:](#) **Filter your results:**

Limits Activated: Molecule: mRNA [Change](#) | [Remove](#)

All (6)

Bacteria (0)

[INSDC \(GenBank\) \(2\)](#)

[mRNA \(6\)](#)

[RefSeq \(4\)](#)

[Manage Filters](#)

Results: 6

[Drosophila melanogaster INDY transporter protein \(Indy\) mRNA, complete cds](#)

1. 2,600 bp linear mRNA

Accession: NM_168779.2

[Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant D, mRNA](#)

Accession: NM_168779.2

[Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant A, mRNA](#)

Accession: NM_168779.2

[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant C, mRNA](#)

4. 2,600 bp linear mRNA

Accession: NM_168779.2 GI: 442633231

[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant B, mRNA](#)

5. 2,572 bp linear mRNA

Accession: NM_168778.2 GI: 442633230

[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Drosophila melanogaster LD24274 full insert cDNA](#)

6. 2,488 bp linear mRNA

Accession: AY102686.1 GI: 20976879

[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)

[Display Settings:](#) Summary, 20 per page, Sorted by Default order

[Send to:](#)

[Turn Off](#) [Clear](#)

I risultati ottenuti derivano da più banche dati. Vogliamo che i risultati provengano da **una sola** banca dati. Selezioniamo **RefSeq** da questa lista

Ora, tutti i risultati ottenuti corrispondono a sequenze che appartengono ad un solo organismo.

Analyze these sequences

[Run BLAST](#)

Find related data

Database:

Select

[Find items](#)

[Search](#)

[See more...](#)

Recent activity

gene: INDY, ricerca per **parola chiave**

ESERCIZIO 1 :

Al momento abbiamo ottenuto la lista di sequenze presenti in NCBI Nucleotide che :

- Contengono **INDY** nella loro scheda descrittiva
- Corrispondono a molecole di mRNA
- Corrispondono a sequenze del moscerino della frutta

Rispondete a queste domande:

- **Quante sono** le sequenze ottenute?
 - Quante di esse sono presenti nella banca dati secondaria **RefSeq**?
-

gene: INDY, ricerca per **parola chiave**

ESERCIZIO 2 :

Limitate i risultati ottenuti alle sole sequenze che appartengono alla banca dati secondaria **RefSeq**.

Suggerimento: cercate qualcosa che vi permetta di restringere il numero di banche dati sulle quali viene effettuata l'estrazione.

Domanda (a cui rispondere dopo aver risolto l'esercizio) :

- Quante sequenze avete ottenuto?
-

Modalità di visualizzazione dei risultati

Chi è stato attento ha notato che, nelle slide precedenti, a volte descrivevamo i risultati riferendoci ad essi come «sequenze». Questo sembra suggerire che per ogni risultato ottenuto il link che punta ad esso restituisca la sequenza della molecola.

Sarà davvero così?

Selezionate il link che punta a :

[Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant A, mRNA](#)

Modalità di visualizzazione dei risultati : ENTRY

Drosophila melanogaster I'm not dead yet (Indy), transcript variant A, mRNA

NCBI Reference Sequence: NM_079426.4

[FASTA](#) [Graphics](#)

[Go to:](#)

LOCUS NM_079426 2484 bp mRNA linear INV 16-JAN-2013

DEFINITION Drosophila melanogaster I'm not dead yet (Indy), transcript variant A, mRNA.

ACCESSION NM_079426

VERSION NM_079426.4 GI:442633232

KEYWORDS RefSeq.

SOURCE Drosophila melanogaster (fruit fly)

ORGANISM [Drosophila melanogaster](#)

Eukaryota; Metazoa; Ecdysozoa; Arthropoda; Hexapoda; Insecta; Pterygota; Neoptera; Endopterygota; Diptera; Brachycera; Muscomorpha; Ephydroidea; Drosophilidae; Drosophila; Sophophora.

REFERENCE 1 (bases 1 to 2484)

AUTHORS Hoskins,R.A., Carlson,J.W., Kennedy,C., Acevedo,D., Evans-Holm,M., Frise,E., Wan,K.H., Park,S., Mendez-Lago,M., Rossi,F., Villasante,A., Dimitri,P., Karpen,G.H. and Celniker,S.E.

TITLE Sequence finishing and mapping of Drosophila melanogaster heterochromatin

JOURNAL Science 316 (5831), 1625-1628 (2007)

PUBMED 17569867

Titolo , banca dati (RefSeq) e identificativo nella banca dati

Lunghezza molecola, tipo molecola, data ultimo aggiornamento

Identificativo sequenza (ACCESSION), organismo di provenienza

Lista pubblicazioni che parlano di questa sequenza (elenco può essere lungo ...)

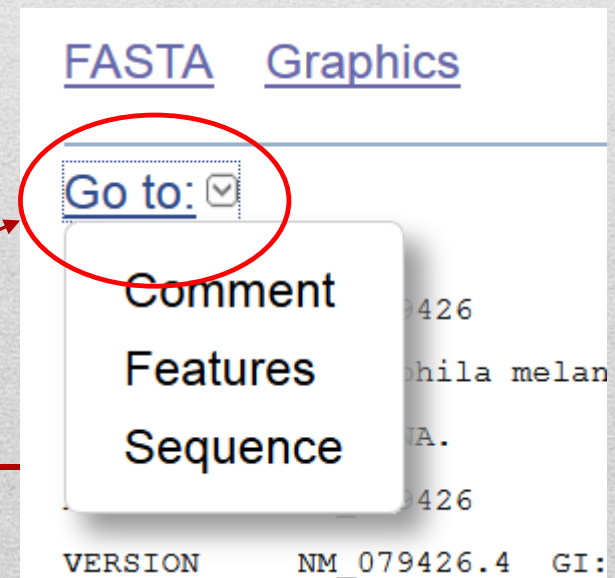
Modalità di visualizzazione ENTRY :

La entries sono simili a schede informative.

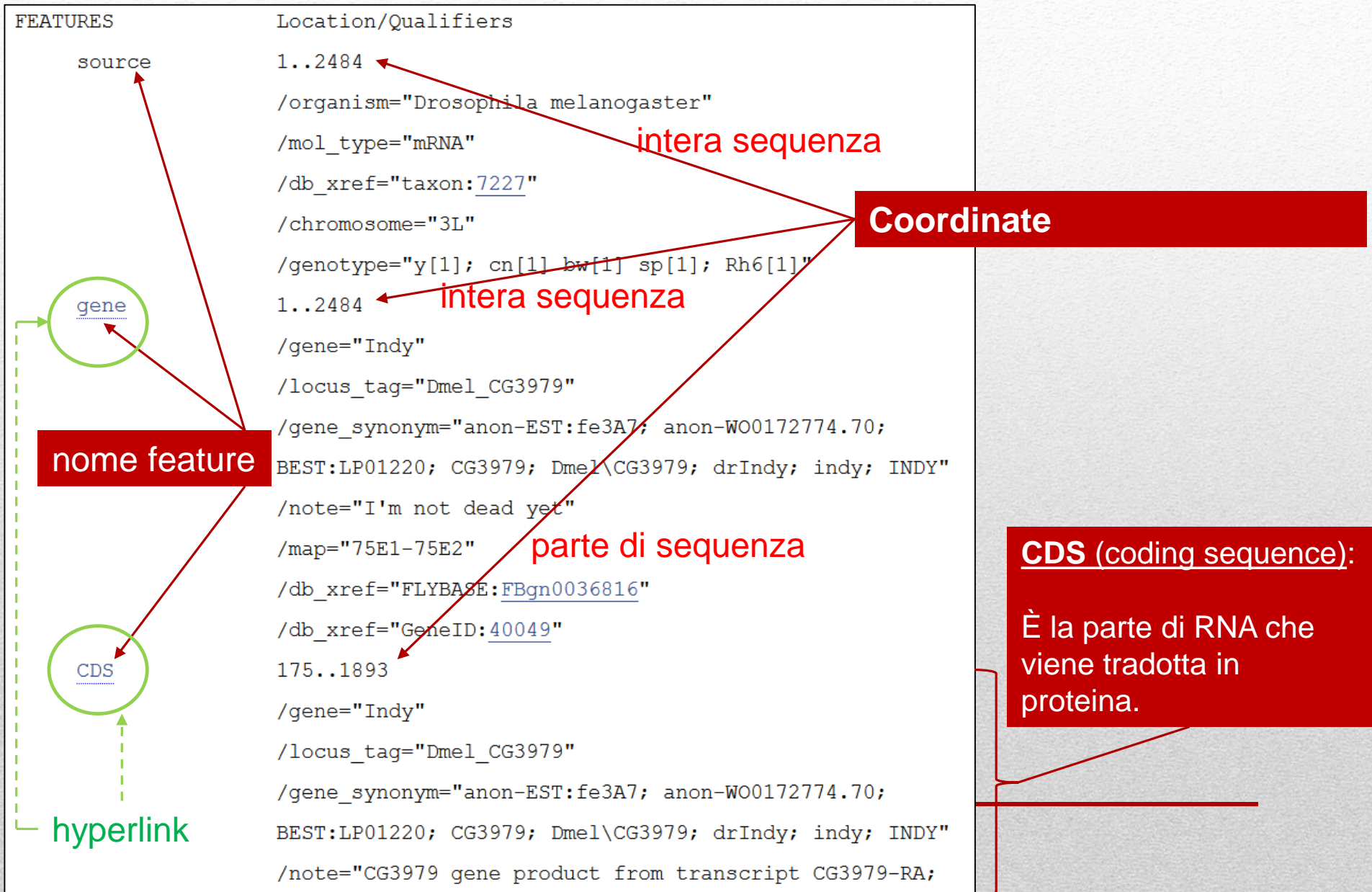
Le entries sono composte da **1 sezione** principale (descritta nella slide precedente) e **3 sezioni aggiuntive**:

- **Comment** : contiene commenti inseriti dai curatori della banca dati
- **Features** : Caratteristiche della sequenza, informazioni generali (che valgono per l'intera sequenza, e informazioni su **tratti specifici della sequenza** (ad es. note riguardanti il tratto di sequenza che va dal nucleotide 10 al nucleotide 20).
- **Sequence** : La sequenza vera e propria

Potete scegliere di andare in una sezione specifica della entry utilizzando il menù **Go to** posto immediatamente all'inizio della entry stessa.



Sezione ENTRY : Feature table




```

1  attcagtcgc gcatttcacc gtttccgaat cggacgaacc gggcgtgctt gctctcctgc
61  tgctttcagc atcggagctc cgataaggat ataactacaa cctaagagag aatccaagcc
121 tcttcctgcc gctagtttcc aaaaatctac acgccaccgc cactggacat caaaatggaa
181 attgaaattg gcgaacaacc ccagcctccg gtgaagtgtc ccaacttctt cgctaaccac
241 tggaaaggat tgggtgtgtt cctgggtgcc ctgctatgtc tgctgttat gctgctaaac
301  gaaggcgccc aatttcggtg catgtacctc cttttggtaa tggccatatt ttgggttacg
361  gaagccttgc ctctctatgt gacgtccatg ataccgatg tggccttccc aataatgggt
421  ataatgagct cggatcagac ttgccgcttg tacttcaagg atacgctggg gatgttcatg
481  ggcggcatta tggtcgcctt ggctgtggag tactgtaatc tacacaaacg tcttgccctg
541  agggtaatcc agatcgtggg ctgcagtccc cgcagattac actttggcct catcatgggt
601  acaatgtttt tgagcatgtg gatttcgaac gccgcctgta ctgccatgat gtgtccgatt
661  atccaagccc tgctggagga gctgcaggct caggggtgtc gcaaaatcaa ccatgagcct
721  caataccaaa tcgttggagg caacaagaaa aacaacgagg atgagccacc ataccaccac
781  aagatcacctc tgtgctacta tctgggcatt gcctacgcct cctcgtggg tggctgtgga
841  accatcatcg gaactgccac caatcttacc ttcaagggca tctacgaggc tctgttcaag
901  aactcccacc aacagatgga cttcccacc ttcatgttct actcgtgccc atccatgttg
961  gtctacacct tgctgacatt cgtgttccct caatggcact tcatgggtct gtggcgtccc
1021 aagagcaagg aggcacagga agtccagagg ggaacgaggg gcgccgatgt cgcacaaaag
1081 gttatcgatc agcgtacaaa ggatctgggt cccatgtcca ttcacgagat ccaagtgatg
1141 attctgttca tttttatggt tgtgatgtac ttcaccgcga agcccggcat ctttttggga
1201 tgggcccatt tgctgaattc caaggacatt cgtaactcct tgcccactat ttttgcgtc
1261 gtcatgtgct tcatgctgcc cgcacaattat gctttcctac gctactgcac cagacgcggt
1321 ggtccagctc ccacgggtcc cactccatcg ctgatcacct ggaagtccat ccagaccaag
1381 gtgccatggg gtctggtgtt cctgcttggc ggtggcttcg ctttggccga aggcagcaag
1441 cagagcggca tggccaagct gattggcaat gctctgattg gattgaaggt tctgcccaac
1501 tctgtctctc tactggtggt catcctggtg gctgtgttcc tgaccgcctt agctccaat
1561 tgggcgattg ccaacattat tattcccgtt ctggccgaga tctcccctggc cattgagatc
1621 catcctctgt acctgactct gcccgctggc ttggcctgca gtatggcctt ccacctgccg
1681 gttagtactc cgcacaacgc ttgggttgct ggctatgcca acattaggac gaaggacatg
1741 gccattgctg gaatcggctc gaccatcatt accatcatca cctgtttgt tttctgccaa
1801  acctggggcc tggctgteta tccgaacctt aactcgttcc ccgaatgggc tcagatttat
1861  gccgcggcag cactgggaaa caagacgcac tagatagtta gtaattagtg taaataacta
1921  acataccctg cacagcgata aagttgagga aaatttaggg aattttaaac gaaaagtgcc
1981  tttgctgaca gcgaaaaaat tgaaaaatat ttaactatgt atacttgcac ttcagagttg
2041  cgaaaagtgt tgatacaaaa gcattaccta ctgtttagaa aaatgtgtta aaaaaaaac
2101  gtatcgcaat atactgttaa tcaggaattg aacacctggg ctacgcactc agctaaaat
2161  ttaaatacaa attaatgtta cttaattggt gcatttagca taaaaatgga aaagatttg
2221  aaaagttaga acagtttgtt caatggcagc cctggcctgc taatatttta aataactaga
2281  ctgagagaaac ttacatattc atacatgttt ttcaacttgt aaaaaatttt aaatgaacaa
2341  ctactcaat acttcattgc gaacccaaat gaacacacaa atagcggtag gctaagctta
2401  aatgatactg tgtacatttt cagatgattt atgttttata tagtttgtaa aaaaatttaa
2461  ataataaaaa gctcaaacga caat

```

Sezione ENTRY : Sequence

Ogni riga contiene 60 caratteri (in questo caso nucleotidi)...

Divisi in gruppi di 10 caratteri (per facilitare conteggi)

Ogni riga inizia con il numero del primo carattere (nucleotide) della riga stessa

Modalità di visualizzazione alternative :

Oltre a visualizzare le informazioni sulla sequenza in modalità ENTRY (a volte detta modalità GenBank) è possibile visualizzare le informazioni in modo diverso.

E' possibile selezionare la modalità di visualizzazione grazie al menù **Display Settings** posto immediatamente all'inizio della entry.

- Summary: solo informazioni principali
- FASTA / FASTA(text): solo sequenza
- GenBank (full) informazioni estese

Dopo aver selezionato la modalità di display premete il pulsante **Apply**. Provate Fasta(text).

Display Settings: GenBank

Format

- Summary
- GenBank
- GenBank (full)
- FASTA
- FASTA (text)
- Graphics
- ASN.1
- Revision History
- Accession List
- GI List

Apply

VERSION NM_079426.4 GI:442633232
KEYWORDS RefSeq.


```
>gi|442633232|ref|NM_079426.4| Drosophila melanogaster I'm not dead yet  
transcript variant A, mRNA
```

```
ATTAGTCGCGCATTTCACCGTTTCCGAATCGGACGAACCGGGCGTGCTTGCTCTCCTGCTGCTTTTCGAG  
ATCGAGTCCCGATAAAGGATATAACTACAACCTAAAGAGGAATCCAAGCCTCCTCCTGCCGCTAGTTTCG  
AAAATCTACACGCCACCGCCACTGGACATCAAATGGAAATTGAAATTGGCGAACAACCCAGCCTCCG  
BTGAGTGCTCCAACCTTCTTCGCTAACCACTGGAAGGGATTGGTTGTGTTTCTGGTGCCGCTGCTATGTC  
TGCCTGTTATGCTGCTAAACGAAGGCGCCGAATTTCCGTTGCATGTACCTCCTTTTGGTAATGGCCATATT  
TTGGTTACGGAAGCCTTGCCCTCTCTATGTGACGTCCATGATACCGATTGTGGCCCTTCCAATAATGGGT  
ATAATGAGCTCGGATCAGACTTGCCGCTTGTACTTCAAGGATACGCTGGTGATGTTTCATGGGCGGCATTA  
TGGCGCCCTGGCTGTGGAGTACTGTAATCTACACAAACGCTTGCCTTGAGGGTAATCCAGATCGTGGG  
CTGAGTCCCCGAGATTACACTTTGGCCTCATCATGGTTACAATGTTTTGAGCATGTGGATTTCGAAC  
ECCGCTGTACTGCCATGATGTGTCCGATATCCAAAGCCGTGCTGGAGGAGCTGCAGGCTCAGGCTGTCT  
GCAAAATGAGCTCGGATCAGACTTGCCGCTTGTACTTCAAGGATACGCTGGTGATGTTTCATGGGCGGCATTA  
ATAACCGGATGAGCTCGGATCAGACTTGCCGCTTGTACTTCAAGGATACGCTGGTGATGTTTCATGGGCGGCATTA  
ACCATGAGCTCGGATCAGACTTGCCGCTTGTACTTCAAGGATACGCTGGTGATGTTTCATGGGCGGCATTA  
AACAGTGGAGTGGAGTACTGTAATCTACACAAACGCTTGCCTTGAGGGTAATCCAGATCGTGGG  
CGTGTGAGTGGAGTACTGTAATCTACACAAACGCTTGCCTTGAGGGTAATCCAGATCGTGGG  
EGACGATGGAGTGGAGTACTGTAATCTACACAAACGCTTGCCTTGAGGGTAATCCAGATCGTGGG  
TTCACCGATGGAGTGGAGTACTGTAATCTACACAAACGCTTGCCTTGAGGGTAATCCAGATCGTGGG  
CTTTTTGAGTGGAGTGGAGTACTGTAATCTACACAAACGCTTGCCTTGAGGGTAATCCAGATCGTGGG  
ETCATGAGTGGAGTGGAGTACTGTAATCTACACAAACGCTTGCCTTGAGGGTAATCCAGATCGTGGG  
CCACGGATGGAGTGGAGTACTGTAATCTACACAAACGCTTGCCTTGAGGGTAATCCAGATCGTGGG  
CCTGCTTGGCGGTGGCTTCGCTTTGGCCGAAGCCAGCAAGCAGAGCGGCATGGCCAAGCTGATTGGCAAT  
GCTCTGATTGGATTGAAGGTTCTGCCAACTCTGCTCTTACTGGTGGTCATCCTGCTGGCTGTGTTCC  
TGACCGCCTTCAGCTCCAATGTGGCGATTGGCCACATTATTATTCCCGTTCTGGCCGAGATGTCCCTGGC  
CATTGAGATCCATCCTCTGTACCTGATCCTGCGCTGGCTTGGCCTGCAGTATGGCCTTCCACCTGCCG  
BTAGTACTCCGCCAACGCTTTGGTTGTGCTGGCTATGCCAACATTAGGACGAAGGACATGGCCATTGCTG  
GAATCGGTCCGACCATCATTACCATCATCACCCTGTTTTGTTTTCTGCCAAACCTGGGGCCTGGTCTTA  
TCCGAACCTTAACTCGTTCCCCGAATGGGCTCAGATTTATGCCGCGGCAGCACTGGGAAACAAGACGCAC  
TAGATAGTTAGTAATTAGTGTAATAAATAACTAACATACCCGTCACAGCGATAAAGTTGAGGAAAAATTTAGGG  
AATTTTAAACGAAAAGTGCCTTTGCTGACAGCGAAAAATGTGAAAAATATTTAACTATGTATACTTGCAT  
TTCAGAGTTGCGAAAAGTTTTGATACAAAAGCATTACCTACTGTTTAGAAAAATGTGTTAAAAAAAAC  
STATCGCAATATACTGTTAATCAGGAATTGAACACCTGGTCTACGCACTCAGCTAAATATTTAAATACAA  
ATTAATGTTACTTAAATGTTGCATTTAGCATAAAAAATGGAAAAGATTGGAAAAGTTAGAACAGTTTGT  
CAATGGCAGCCCTGGCCTGCTAATATTTTAAATAACTAGACTGAGAGAACCTTACATATTCATACATGTTT  
TTCACCTTGTAATAAATTTTTAAATGAACAACCTCACTCAATACTTCATGCGAACCAAAATGAACACACAA  
ATAGCGGTAGGCTAAGCTTAAATGATACTGTGTACATTTTCAGATGATTTATGTTTTATATAGTTTGTA  
AAAATATTAATAATAAAAAAGCTCAAACGACAAT
```

FORMATO FASTA:

- **Prima riga:** simbolo > seguito da informazioni sulla sequenza
- **Dalla seconda riga in poi:**
Sequenza

Display mode: **FASTA** **(text)**

Il formato di visualizzazione FASTA (text) mostra solo la sequenza (senza spazi vuoti e senza numeri). Inoltre **non contiene caratteri invisibili** (a differenza delle altre modalità di visualizzazione della sequenza).

In questo formato la sequenza può essere usata come input per programmi che effettuano analisi su di essa (ad esempio composizione: frequenza caratteri A,C,G e T in questo caso)

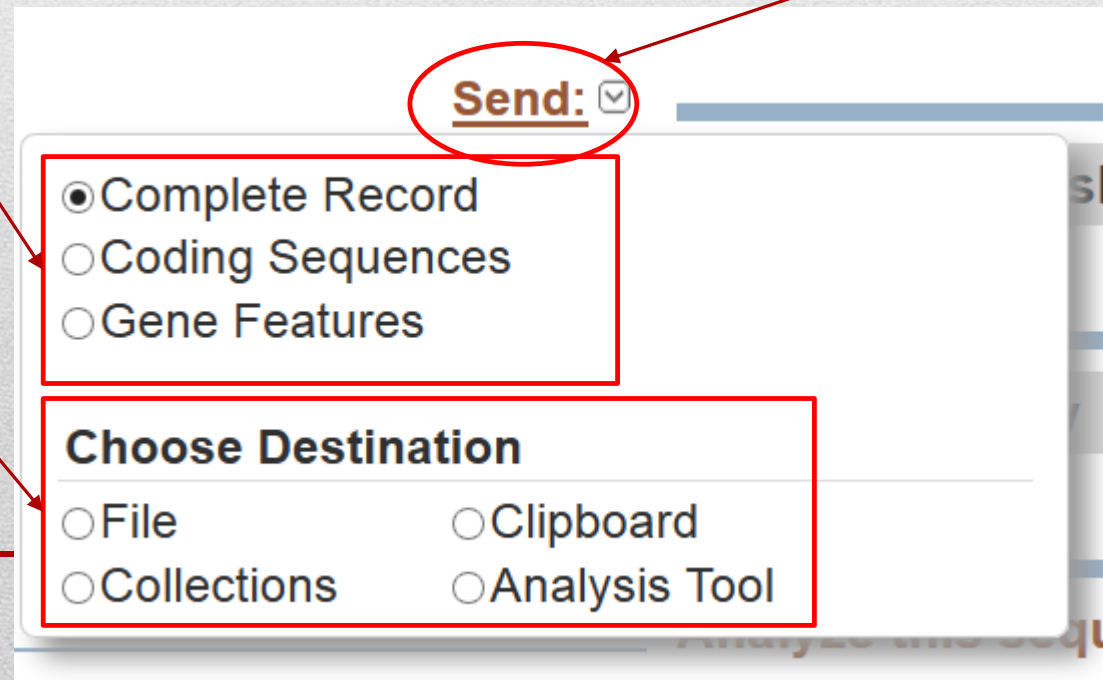
Una volta che abbiamo reperito le informazioni che cerchiamo come le utilizziamo?

Abbiamo varie opzioni ... la più banale è fare una stampa ma questa via presenta diverse limitazioni e serve solo quando abbiamo necessità, ad esempio, di aggiungere delle note manuali. Esiste un modo migliore per **esportare e/o salvare in un file** solo le informazioni di cui abbiamo bisogno.

E' possibile esportare dati (operazione in **2 passaggi**) tramite menù **Send** :

1. **Selezionate** le informazioni a cui siete interessati (ad es. **Complete Record**)

2. **Selezionate** la destinazione di queste informazioni (appunti, file, strumenti di analisi, ...)

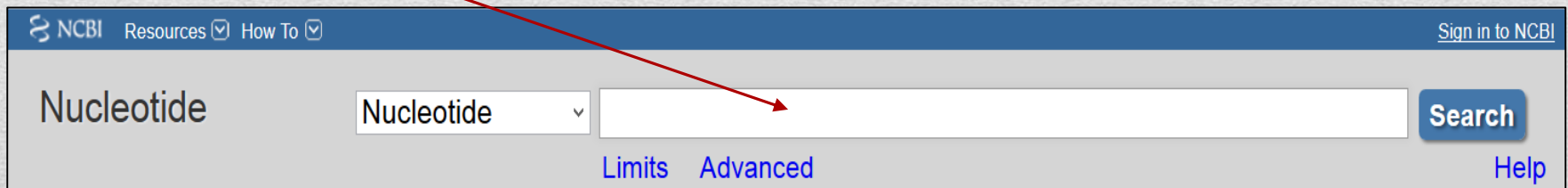


Riepilogo : ricerca per parola chiave

Quella che abbiamo visto è la modalità classica di ricerca di informazioni nelle banche dati biologiche. Essa si definisce «ricerca per parola chiave».

Nel nostro esempio la parola chiave era il **nome del gene** (INDY).

In origine non esistevano tutti gli strumenti che abbiamo visto (i menù nell'interfaccia web). L'unica modalità con cui era possibile raffinare la ricerca era quella di costruire complesse stringhe di testo da inserire nella casella di ricerca dove avevamo inserito la parola **INDY**.



The screenshot shows the top navigation bar of the NCBI website. On the left, there is a logo and the text 'NCBI Resources' with a dropdown arrow, and 'How To' with a dropdown arrow. On the right, there is a link 'Sign in to NCBI'. Below the navigation bar, there is a search interface. On the left, the word 'Nucleotide' is displayed. Next to it is a dropdown menu currently showing 'Nucleotide'. To the right of the dropdown is a large white search input field. A red arrow points from the underlined text 'casella di ricerca' in the previous paragraph to this search input field. To the right of the search field is a blue 'Search' button. Below the search field, there are two links: 'Limits' and 'Advanced'. To the right of the search field, there is a 'Help' link.

Anche se in maniera non evidente ... il sistema funziona ancora così ... solo che la stringa di ricerca viene costruita dinamicamente in base alle scelte che operiamo sugli strumenti dell'interfaccia web

Costruzione **dinamica** stringa di ricerca (I)

Nucleotide [Save search](#) [Limits](#) [Advanced](#) [Help](#)

The information on this web site remains accessible; but, due to the lapse in government funding, the information may not be up to date, and the agency may not be able to respond to inquiries until appropriations are enacted. For updates regarding government operating status see [USA.gov](#).

Display Settings: Summary, 20 per page, Sorted by Default order **Send to:** **Filter your results:**

Limits Activated: Molecule: mRNA, Source database: RefSeq [Change](#) | [Remove](#)

Results: 10

- [Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant D, mRNA](#)
1. 2,581 bp linear mRNA
Accession: NM_001169994.2 GI: 442633233
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)
- [Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant A, mRNA](#)
2. 2,484 bp linear mRNA
Accession: NM_079426.4 GI: 442633232
[GenBank](#) [FASTA](#) [Graphics](#) [Related Sequences](#)
- [Drosophila melanogaster I'm not dead yet \(Indy\), transcript variant C, mRNA](#)
3. 2,600 bp linear mRNA
Accession: NM_168779.2 GI: 442633231

Filter your results:

- All (10)
- Bacteria (0)
- INSDC (GenBank) (0)
- [mRNA \(10\)](#)
- [RefSeq \(10\)](#)

[Manage Filters](#)

Top Organisms [Tree]

- Homo sapiens (4)
- Drosophila melanogaster (4)
- Mus musculus (1)
- Oryctolagus cuniculus (1)

Analyze these

Quando cerchiamo la parola chiave INDY ed impostiamo dei filtri tipo di molecola: mRNA e source database: RefSeq

Costruzione **dinamica** stringa di ricerca (II)

Mano a mano che aggiungiamo dei filtri il sistema aggiorna una stringa di testo che specifica tutti i passaggi della nostra ricerca ...

E li utilizza per aggiornare una stringa di testo che, se salvata ci permetterà di ritornare al sito e rieffettuare la ricerca incollando la stringa nella casella posta di lato al pulsante **Search** (e premendo **Search**)

Search details

```
INDY[All Fields]
AND
(biomol_mRNA[PROP]
AND
srcdb_refseq[PROP])
```

Search

[See more...](#)

INDY[All Fields] AND (biomol_mRNA[PROP] AND srcdb_refseq[PROP])

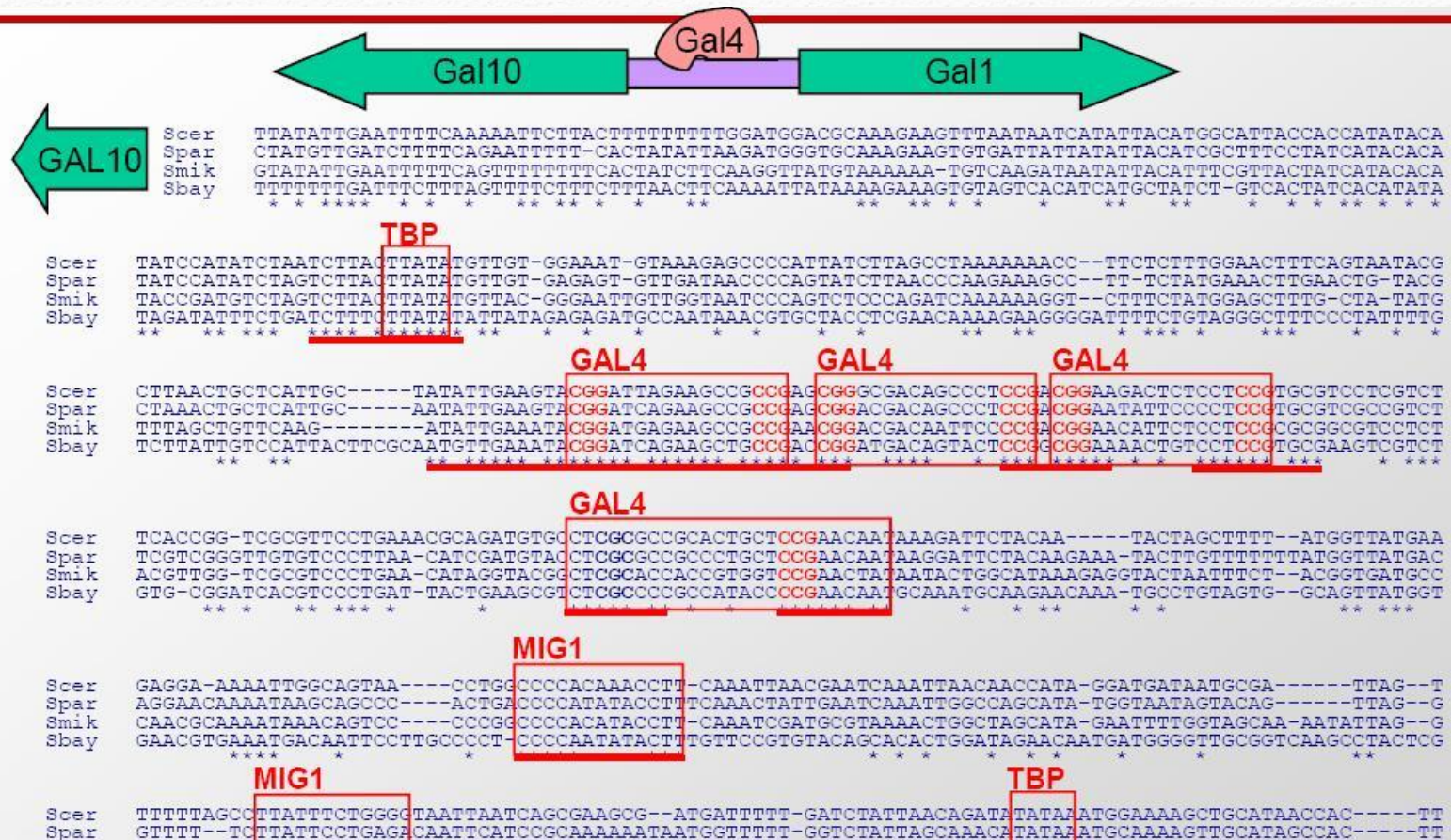
Quando cerchiamo la parola chiave INDY ed impostiamo nei limiti (Limits) tipo di molecola: mRNA e source database: RefSeq

Ricerca mediante similarità di sequenza

Abbiamo potuto trovare informazioni sul gene INDY secondo le modalità descritte nelle slide precedenti solo perchè avevamo a disposizione una parola chiave da utilizzare (**INDY**) ... ma come possiamo cercare informazioni per una molecola **di cui non sappiamo nulla?**

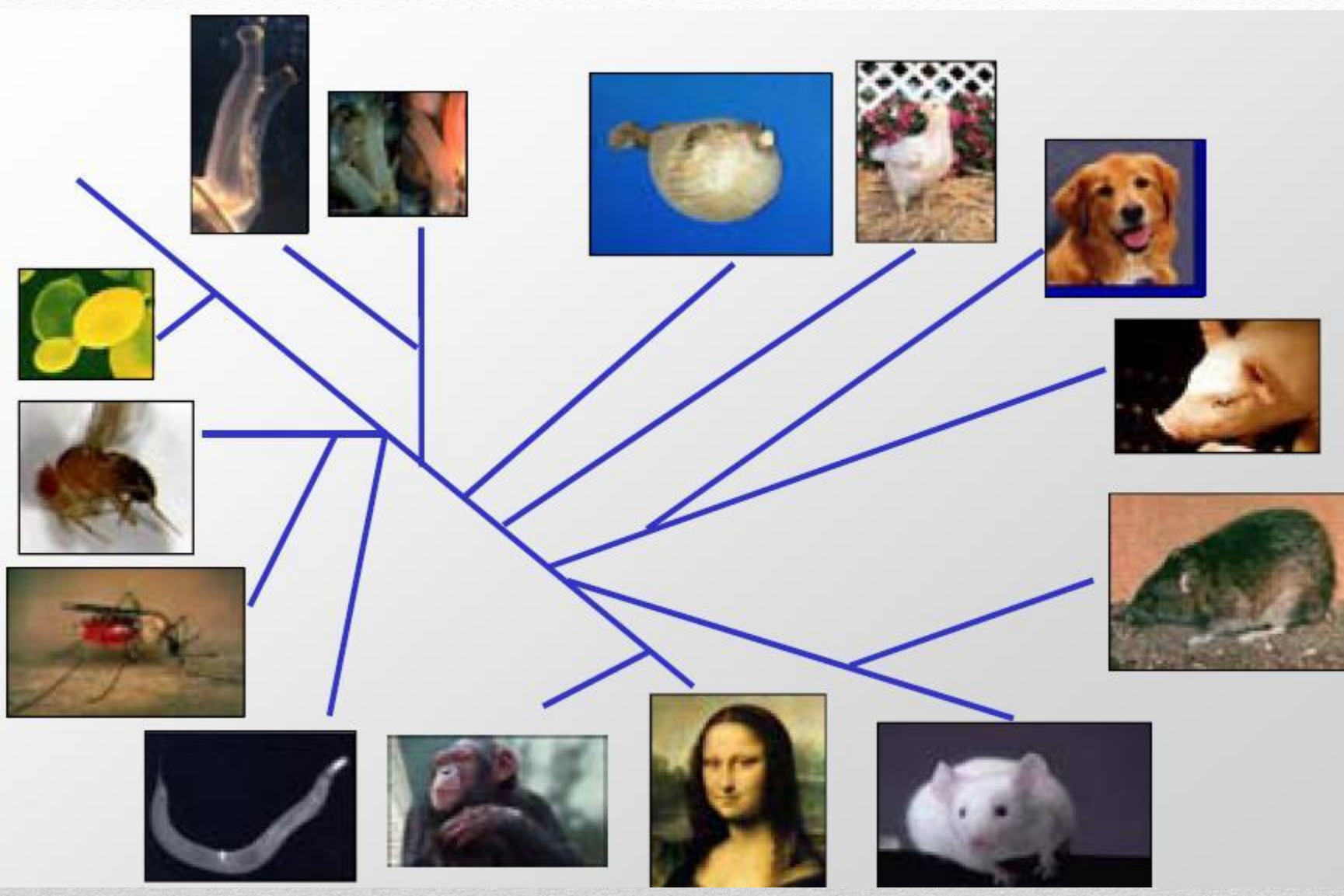
Questa situazione si verifica quando otteniamo la sequenza in laboratorio. L'unica cosa che otteniamo è, appunto, la sequenza. Nient'altro. **In questo caso una ricerca per parola chiave è impossibile da realizzare!**

Elementi conservati nel corso dell'evoluzione



Possiamo “LEGGERE” l’evoluzione per trovare elementi funzionali

Come possiamo ALLINEARE due biosequenze?



I geni/genomi cambiano nel tempo

Stato iniziale

A C G T C A T C A

mutazione

A C G T **G** A T C A

delezione

A **X** G T G **X** T C A

inserzione

A G T G T C A

T A G T G T C A

Stato finale

T A G T G T C A

Obiettivo dell'allineamento:

Stato iniziale

A	C	G	T	C	A	T	C	A
---	---	---	---	---	---	---	---	---

?



Stato finale

T	A	G	T	G	T	C	A
---	---	---	---	---	---	---	---

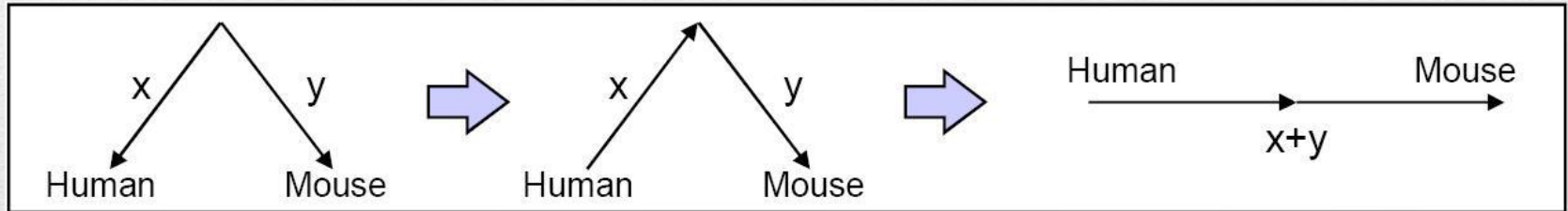
INFERENZA DELLE OPERAZIONI DI MODIFICA (editing)

Bio

CS

FORMALIZZAZIONE DEL PROBLEMA

1. Definizione di un **set di operazioni evolutive** (mutazione, inserzione e delezione)



Simmetria delle operazioni ($A \rightarrow C, C \rightarrow A$) permette reversibilità rispetto al tempo. Questa è una scelta di “disegno” della formalizzazione del problema.

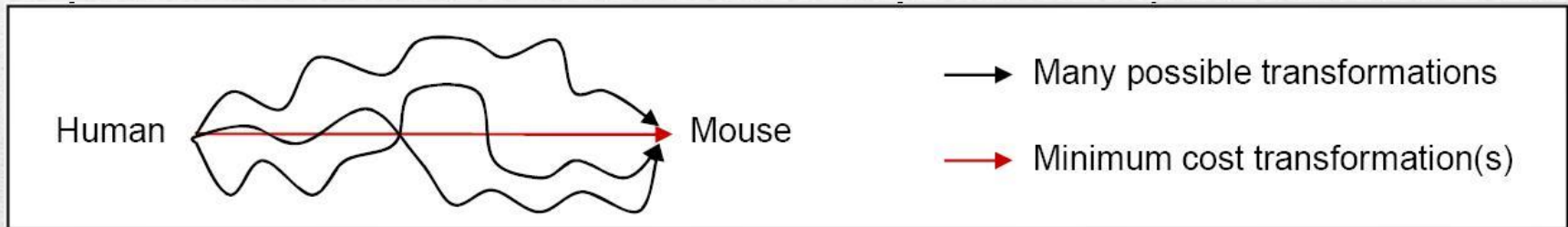
Eccezione (Bio) : dinucleotidi CpG metilati \rightarrow TpG/CpA
(perdita simmetria)

Bio

CS

FORMALIZZAZIONE DEL PROBLEMA

2. Definizione di un **criterio di ottimalità (minimo numero operazioni, minimo costo complessivo,...)**



E' **IMPOSSIBILE** inferire la serie **ESATTA** di operazioni che portano dalla sequenza A alla sequenza B!

Rasoio di Occam: “Avendo a disposizione diverse ipotesi equivalentemente competitive rispetto ad un dato problema è buona norma scegliere la più semplice e scartare le altre”.

Bio

CS

FORMALIZZAZIONE DEL PROBLEMA

3. Definizione di un **algoritmo in grado di raggiungere questa condizione di ottimalità** (o in grado di approssimarla)

Bio

CS

Predicibilità

Correttezza

Rilevanza

Casi speciali

Trade-off

Algoritmi

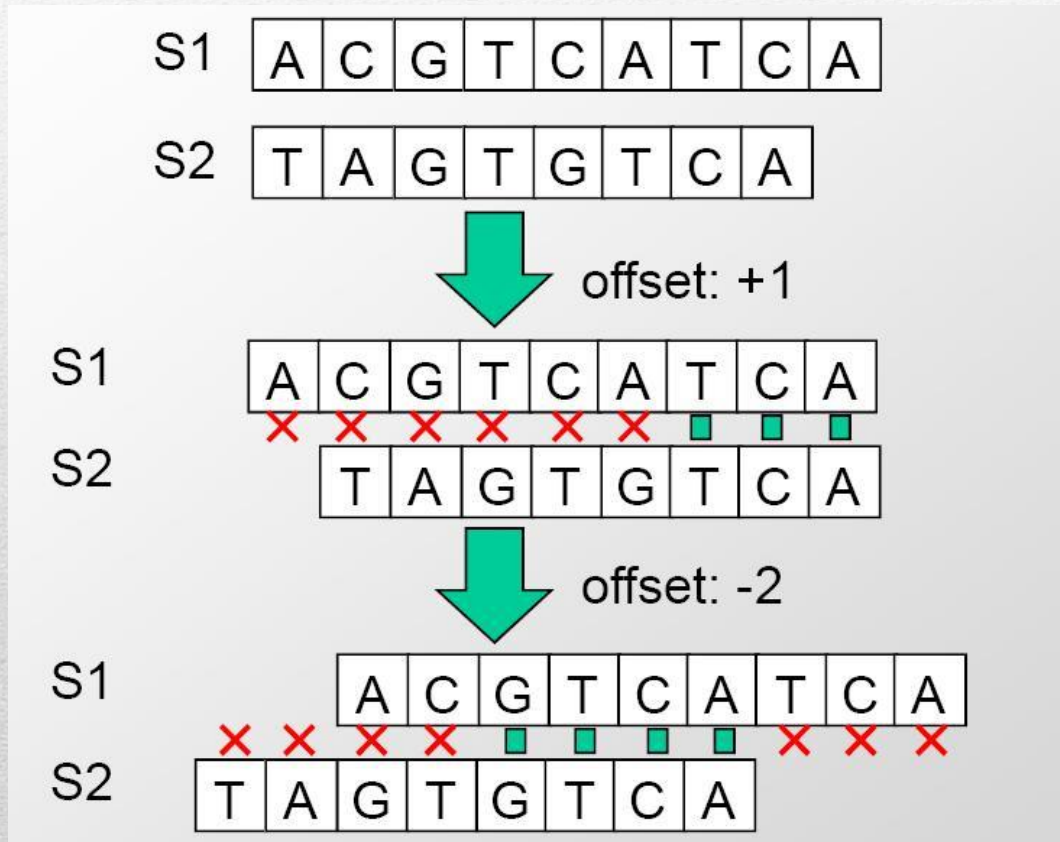
Assunzioni

Implementazione

Trattabilità

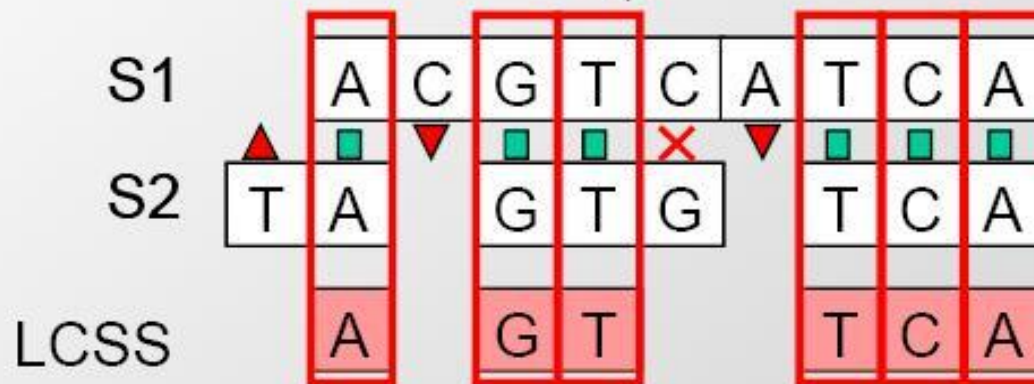
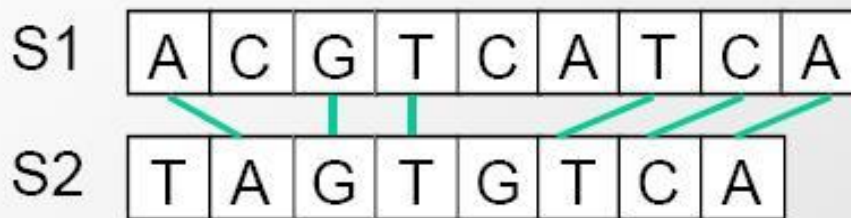
Computabilità

“Date due sequenze S1 ed S2 possibilmente originatesi da una medesima sequenza ancestrale quale è la più lunga sottosequenza in comune (LCS) ?” (gap non ammessi)



“Date due sequenze S1 ed S2 possibilmente originatesi da una medesima sequenza ancestrale quale è la più lunga sottosequenza in comune (LCS) ?”

(gap ammessi)



E' basata su **edit distance**:

- Numero di cambiamenti per passare da S1 a S2
- Funzione di scoring uniforme

- **Gap ammessi (penalità fissa):**
 - Operazioni di inserzione e delezione
 - Medesimo costo per ogni carattere inserito/deleto
- **Penalità variabili per le operazioni di editing:**
 - Transizioni (Pirimidine↔Pirimidine, Purine↔Purine)
 - Trasversioni (variazioni Pirimidine↔Purine)
 - Polimerasi confonde “relativamente spesso” A/G e C/T

Scoring function:

Match(x,x) = +1

Mismatch(A,G) = $-\frac{1}{2}$

Mismatch(C,T) = $-\frac{1}{2}$

Mismatch(x,y) = -1

	A	G	T	C
A	+1	$-\frac{1}{2}$	-1	-1
G	$-\frac{1}{2}$	+1	-1	-1
T	-1	-1	+1	$-\frac{1}{2}$
C	-1	-1	$-\frac{1}{2}$	+1

purine pyrimid.

Transitions:

$A \leftrightarrow G$, $C \leftrightarrow T$ common
(lower penalty)

Transversions:

All other operations

Molte variazioni:

es. Penalità variabili in accordo con il numero di gap

...

(riuscite a proporre delle varianti?)

Come possiamo costruire il “miglior” allineamento?

S1 A | C | G | T | C | A | T | C | A
S2 T | A | G | T | G | T | C | A

- **Data una funzione di scoring additiva :**
 - Costo mutazioni (AG, CT, altre)
 - Costo inserzione / delezione
 - Premio per match
 - **Serve algoritmo per provare il miglior allineamento:**
 - Enumerazione di tutti i possibili allineamenti?
 - Come la realizzereste?
 - Quanti sono i possibili allineamenti di due sequenze?
-

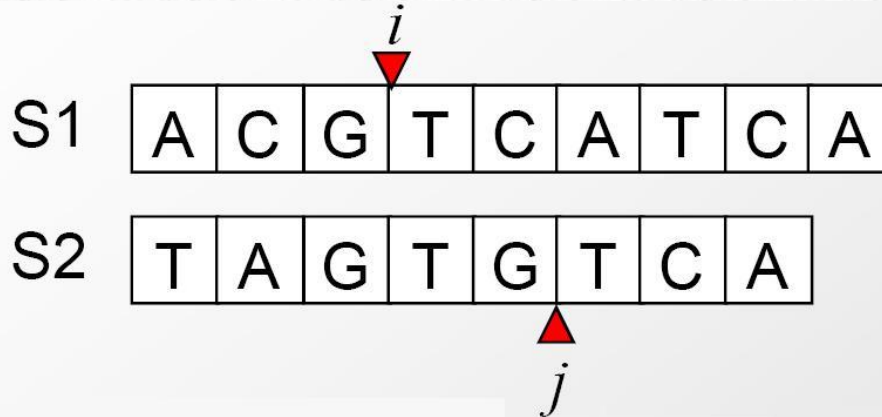
CS**Modi di allineare due sequenze di lunghezza m, n:**

$$\binom{n+m}{m} = \frac{(m+n)!}{(m!)^2} \approx \frac{2^{m+n}}{\sqrt{\pi \cdot m}}$$

- **Per due sequenze di lunghezza n:**

n	Enumerazione	Tecnica presentata oggi
10	184756	100
20	1.40E+11	400
100	9.00E+58	10000

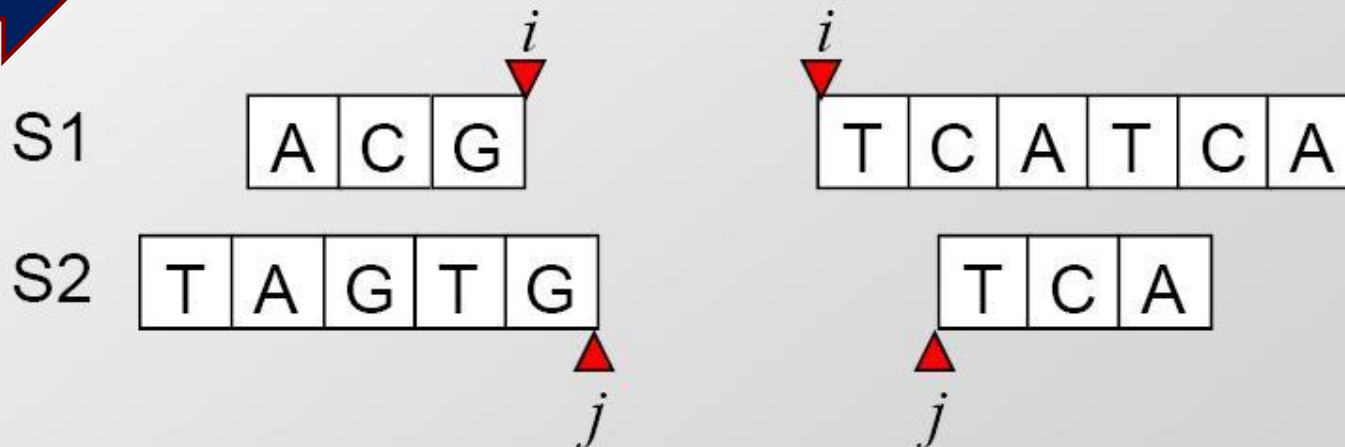
PUNTO CRUCIALE: gli score sono additivi:



• **Calcolo RICORSIVO del miglior allineamento:**

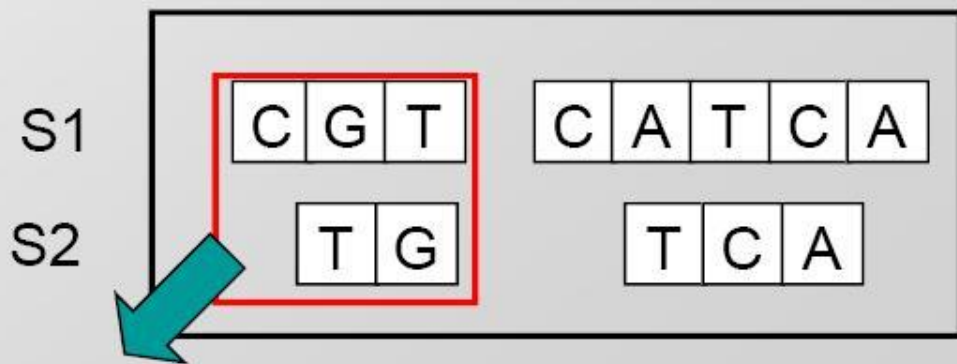
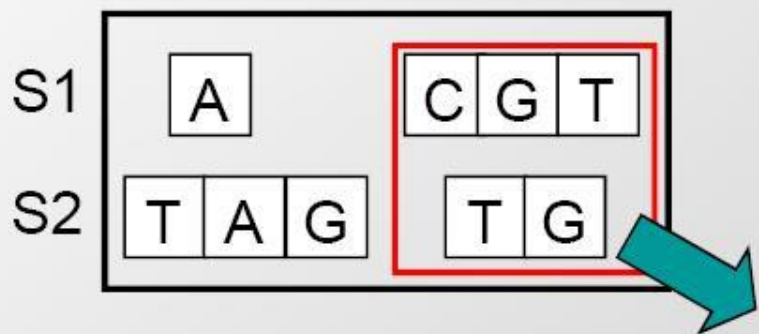
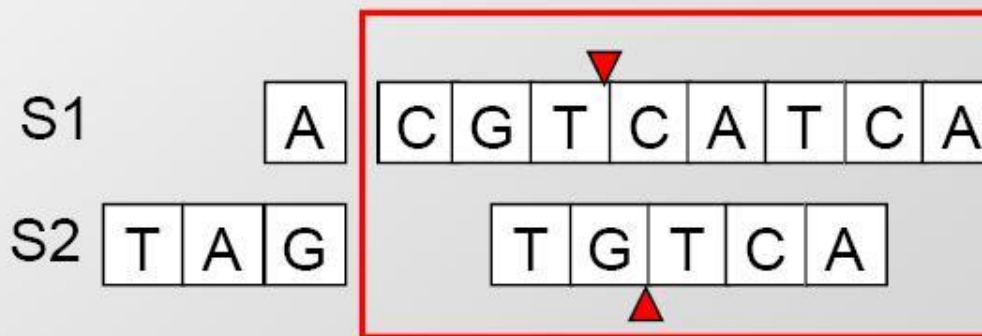
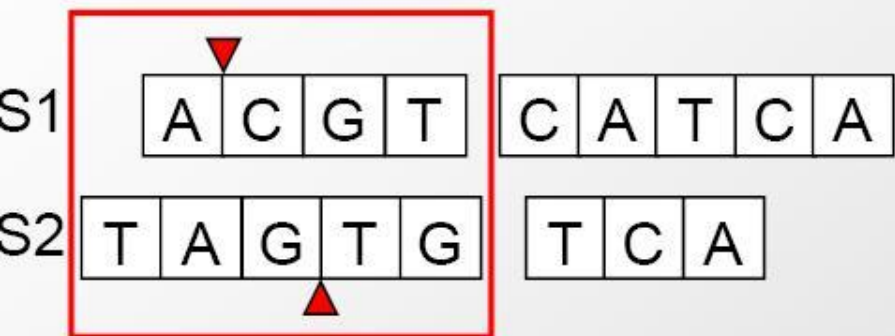
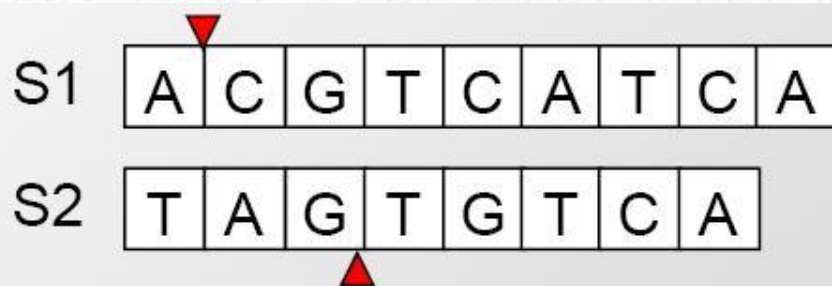
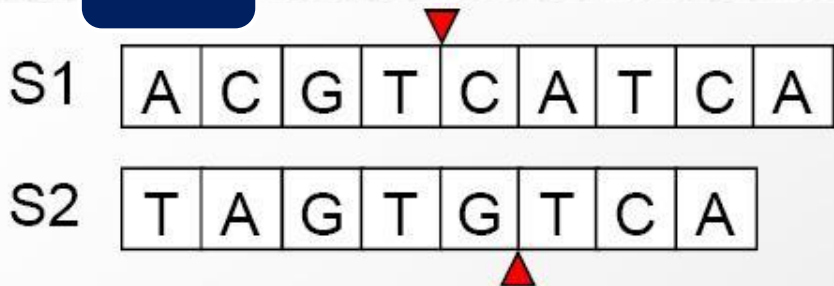
- Per una data coppia di sequenze allineate (S1,S2) il miglior allineamento è:

- Miglior allineamento di S1[1..i] e S2[1..j]
- + Miglior allineamento di S1[i..n] e S2[j..m]



CS

IDEA : riutilizzo delle operazioni svolte dal calcolatore



Sottoproblemi **IDENTICI** ! Possiamo riutilizzare più volte la stessa soluzione.

- Creazione di un DIZIONARIO (le chiavi sono le sequenze)
- Quando dobbiamo allineare due sequenze cerchiamo una soluzione corrispondente alle due sequenze
- SE ESISTE: ritorniamo il risultato
- SE NON ESISTE:
 - calcoliamo la soluzione
 - inseriamo la soluzione nel dizionario
 - restituiamo la soluzione
- Assicuriamoci di non duplicare i calcoli:
è necessario calcolare un sottoallineamento una sola volta

A diagram consisting of a large light-colored arrow pointing to the right. Inside the arrow, on the left, is a dark blue rounded square containing the white text 'CS'. To its right is a red rounded square containing the white text 'Bio'. To the right of the arrow's tip is a dark red rectangular box containing white text.

CS

Bio

**ALLINEAMENTO LOCALE DI
SEQUENZE NUCLEOTIDICHE**

Algoritmo più noto per soluzione del problema:

Smith-Waterman (1981)

Obiettivo:

“Date due sequenze relativamente lunghe, determinare la sottosequenza della prima sequenza che realizza il maggior grado di similarità con una sottosequenza della seconda sequenza”.

A diagram consisting of a large light-colored arrow pointing to the right. Inside the arrow, on the left, is a dark blue rounded square containing the white text 'CS'. To its right is a red rounded square containing the white text 'Bio'. To the right of the arrow's tip is a dark red rectangular box containing the white text 'ALLINEAMENTO LOCALE DI SEQUENZE NUCLEOTIDICHE'.

CS

Bio

ALLINEAMENTO LOCALE DI
SEQUENZE NUCLEOTIDICHE

Smith-Waterman (1981)

1) Inizializzazione

Data una coppia di sequenze $S_A = a_1, a_2, \dots, a_i, a_n$ e $S_B = b_1, b_2, \dots, b_j, b_m$ costruire una matrice $\mathbf{H} = \|\|H_{i,j}\|\|$, costituita da $n+1$ righe e $m+1$ colonne che viene inizializzata ponendo:

$$H_{i,0} = H_{0,j} \text{ per } 0 \leq i \leq n \text{ e } 0 \leq j \leq m$$

CS

Bio

ALLINEAMENTO LOCALE DI SEQUENZE NUCLEOTIDICHE

Smith-Waterman (1981)

2) Calcolo

Per ciascuna coppia di elementi a_i e b_j appartenenti alle sequenze S_A e S_B si definisce un punteggio di **similarità s** .

Per il confronto di sequenze nucleotidiche tale punteggio viene generalmente definito come segue:

$$s(a_i, b_j) = \alpha \quad \text{se } a_i = b_j \quad (\text{usualmente } > 0)$$

$$s(a_i, b_j) = \beta \quad \text{se } a_i \neq b_j \quad (\text{usualmente } \leq 0)$$

$$W_k = \gamma + \delta(k-1) \quad \text{dove } \gamma = \text{costo fisso apertura gap}$$

$$\delta = \text{penalità aggiuntiva per estensione gap}$$

CS

Bio

**ALLINEAMENTO LOCALE DI
SEQUENZE NUCLEOTIDICHE****Smith-Waterman (1981)**

2) Calcolo (continua)

I valori della matrice H_{ij} vengono determinati secondo l'equazione:

$$H_{i,j} = \max(H_{i-1,j-1} + s(a_i, b_j), H_{i-1,j} - \text{pgap}, H_{i,j-1} - \text{pgap}, 0)$$

NB: in questo esempio invece di usare $W_k = \gamma + \delta(k-1)$ come penalità dinamica per i gap di diversa lunghezza abbiamo utilizzato una penalità *costante* (**pgap**) indipendente dall'estensione dei gap.

CS

Bio

ALLINEAMENTO LOCALE DI SEQUENZE NUCLEOTIDICHE

Smith-Waterman (1981)

2) Calcolo (continua)

I valori della matrice H_{ij} vengono determinati secondo l'equazione:

$$H_{i,j} = \max(H_{i-1,j-1} + s(a_i, b_j), H_{i-1,j} - \text{pgap}, H_{i,j-1} - \text{pgap}, 0)$$

Confronto tra
nucleotidi
(match/mismatch)
score \neq

delezione
in sequenza B
di lunghezza 1

delezione
in sequenza A
di lunghezza 1



CS

Bio

ALLINEAMENTO LOCALE DI SEQUENZE NUCLEOTIDICHE

Smith-Waterman (1981)

2) Calcolo (continua)

I valori della matrice H_{ij} vengono determinati secondo l'equazione:

$$H_{i,j} = \max(H_{i-1,j-1} + s(a_i, b_j), H_{i-1,j} - \text{pgap}, H_{i,j-1} - \text{pgap}, 0)$$

INTERPRETAZIONE: Questa formula **ricorsiva** considera *tutte le possibili terminazioni* di segmenti allineati alle cui estremità si trovano a_i e b_j .

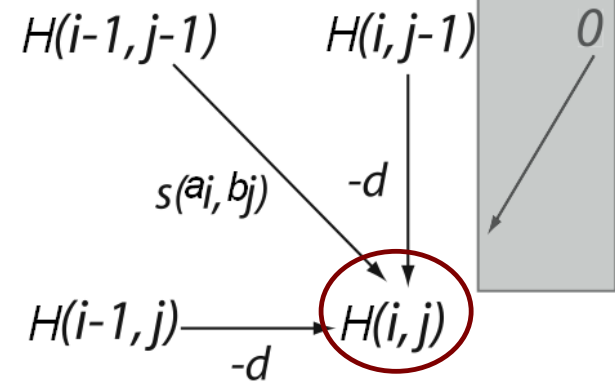
CS

Bio

ALLINEAMENTO LOCALE DI SEQUENZE NUCLEOTIDICHE

2) Calcolo

		Sequence A													
		C	A	G	C	C	U	C	G	C	U	U	A	G	
Sequence B	A	0,0	1,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	1,0	0,0	
	A	0,0	1,0	0,7	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	1,0	0,7	
	U	0,0	0,0	0,8	0,3	0,0	0,0	0,0	0,0	0,0	1,0	1,0	0,0	0,7	
	G	0,0	0,0	1,0	0,3	0,0	0,0	0,7	1,0	0,0	0,0	0,7	0,7	1,0	
	C	1,0	0,0	0,0	2,0	1,3	0,3	1,0	0,3	2,0	0,7	0,3	0,3	0,3	
	C	1,0	0,7	0,0	1,0	3,0	1,7	?							
	A														
	U														
	U														
	G														
	A														
	C														
	G														
	G														



CS

Bio

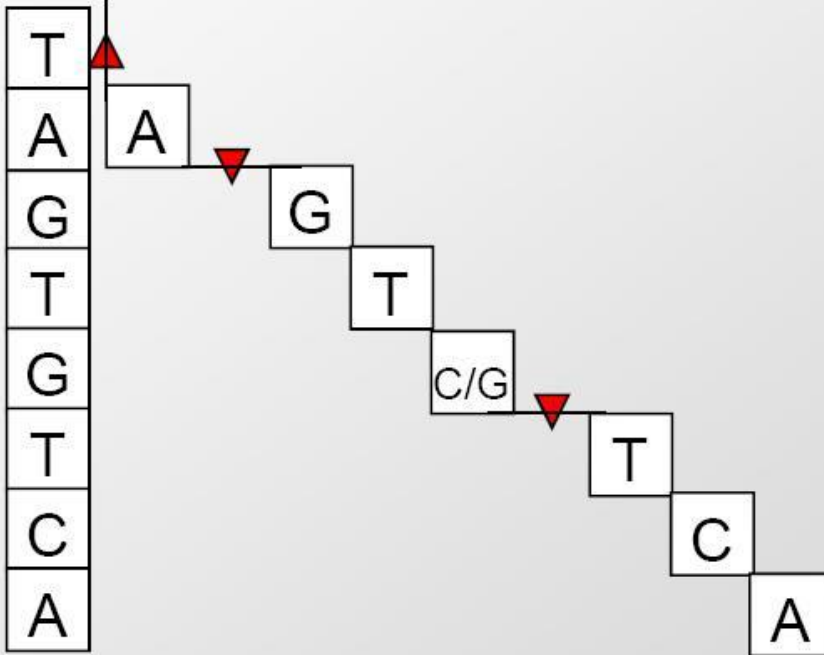
DUALITA' :

“miglior allineamento” ↔ miglior percorso attraverso la matrice !

S1

A C G T C A T C A

S2



NUOVO OBIETTIVO:

Trovare il miglior
percorso attraverso la
matrice !

CS

Bio

ALLINEAMENTO LOCALE DI SEQUENZE NUCLEOTIDICHE

3) Backtracking

-Trovare valore massimo

- Ripercorrere la matrice fino a quando non raggiungo l'angolo in alto a sinistra o trovo 0 sulla diagonale

		Sequence A												
		C	A	G	C	C	U	C	G	C	U	U	A	G
Sequence B	A	0,0	1,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	1,0	0,0
	A	0,0	1,0	0,7	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	1,0	0,7
	U	0,0	0,0	0,8	0,3	0,0	0,0	0,0	0,0	0,0	1,0	1,0	0,0	0,7
	G	0,0	0,0	1,0	0,3	0,0	0,0	0,7	1,0	0,0	0,0	0,7	0,7	1,0
	C	1,0	0,0	0,0	2,0	1,3	0,3	1,0	0,3	2,0	0,7	0,3	0,3	0,3
	C	1,0	0,7	0,0	1,0	3,0	1,7	1,3	1,0	1,3	1,7	0,3	0,0	0,0
	A	0,0	2,0	0,7	0,3	1,7	2,7	1,3	1,0	0,7	1,0	1,3	1,3	0,0
	U	0,0	0,7	1,7	0,3	1,3	2,7	2,3	1,0	0,7	1,7	2,0	1,0	1,0
	U	0,0	0,3	0,3	1,3	1,0	2,3	2,3	2,0	0,7	1,7	2,7	1,7	1,0
	G	0,0	0,0	1,3	0,0	1,0	1,0	2,0	3,3	2,0	1,7	1,3	2,3	2,7
	A	0,0	1,0	0,0	1,0	0,3	0,7	0,7	2,0	3,0	1,7	1,3	2,3	2,0
	C	1,0	0,0	0,7	1,0	2,0	0,7	1,7	1,7	3,0	2,7	1,3	1,0	2,0
	G	0,0	0,7	1,0	0,3	0,7	1,7	0,3	2,7	1,7	2,7	2,3	1,0	2,0
	G	0,0	0,0	1,7	0,7	0,3	0,3	1,3	1,3	2,3	1,3	2,3	2,0	2,0

CS

Bio

ALLINEAMENTO LOCALE DI SEQUENZE NUCLEOTIDICHE

3) Backtracking

-Trovare valore massimo

- Ripercorrere la matrice fino a quando non raggiungo l'angolo in alto a sinistra o trovo 0 sulla diagonale

		Sequence A													
		C	A	G	C	C	U	C	G	C	U	U	A	G	
Sequence B	A	0,0	1,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,0	1,0	0,0	
	A	0,0	1,0	0,7	0,0	0,0	0,0								
	U	0,0	0,0	0,8	0,3	0,0	0,0								
	G	0,0	0,0	1,0	0,3	0,0	0,0								
	C	1,0	0,0	0,0	2,0	1,3	0,3								
	C	1,0	0,7	0,0	1,0	3,0	1,7	1,3	1,0	1,3	1,7	0,3	0,0	0,0	
	A	0,0	2,0	0,7	0,3	1,7	2,7	1,3	1,0	0,7	1,0	1,3	1,3	0,0	
	U	0,0	0,7	1,7	0,3	1,3	2,7	2,3	1,0	0,7	1,7	2,0	1,0	1,0	
	U	0,0	0,3	0,3	1,3	1,0	2,3	2,3	2,0	0,7	1,7	2,7	1,7	1,0	
	G	0,0	0,0	1,3	0,0	1,0	1,0	2,0	3,3	2,0	1,7	1,3	2,3	2,7	
	A	0,0	1,0	0,0	1,0	0,3	0,7	0,7	2,0	3,0	1,7	1,3	2,3	2,0	
	C	1,0	0,0	0,7	1,0	2,0	0,7	1,7	1,7	3,0	2,7	1,3	1,0	2,0	
	G	0,0	0,7	1,0	0,3	0,7	1,7	0,3	2,7	1,7	2,7	2,3	1,0	2,0	
	G	0,0	0,0	1,7	0,7	0,3	0,3	1,3	1,3	2,3	1,3	2,3	2,0	2,0	

G C C A U U G
| | | | . |
G C C - U C G

Ricerca mediante similarità di sequenza

Supponiamo di aver ottenuto, in laboratorio, questa sequenza di cui non sappiamo nulla:

>Sequenza_sconosciuta

```
CTCGCAGGCTCCAGGGGCGGGGCGTGGCCGGGGCGCAGCGACGGGCGGGAGGTCCGGCCGGGCGCGCGC  
GCCCCGCCACACGCACGCCGGGCGTGCCAGTTTATAAAGGGAGAGAGCAAGCAGCGAGTCTTGAAGCTC  
TGTTTGGTGCTTTGGATCCATTTCCATCGGTCCTTACAGCCGCTCGTCAGACTCCAGCAGCCAAGATGGT  
GAAGCAGATCGAGAGCAAGACTGCTTTTCAGGAAGCCTTGGACGCTGCAGGTGATAAACTTGTAGTAGTT  
GACTTCTCAGCCACGTGGTGTGGGCCTTGCAAAATGATCAAGCCTTTCTTTCATGATGTTGCTTCAGAGT  
GTGAAGTCAAATGCATGCCAACATTCCAGTTTTTTAAGAAGGGACAAAAGGTGGGTGAATTTTCTGGAGC  
CAATAAGGAAAAGCTTGAAGCCACCATTAATGAATTAGTCTAATCATGTTTTTCTGAAAATATAACCAGCC  
ATTGGCTATTTAAAACCTTGAATTTTTTTAATTTACAAAAATATAAAATATGAAGACATAAACCCAGTTG  
CCATCTGCGTGACAATAAAACATTAATGCTAACACTTTTTAAAACCGTCTCATGTCTGAATAGCTTTCAA  
AATAAATGTGAAATGGTCATTTAATGTATTTTCTATATTCTCAATCACTTTTTAGTAACTTGTAGGCC  
ACTGATTATTTTAAAGATTTTAAAAATTATTATTGCTACCTTAATGTATTGCTACAAAAATCTCTTGTGG  
GGCAATGCAGGTAATAAAGTAGTATGTTGTTATTTGTAAAAAAAAAAAAAAAAAAAA
```

E' una semplice sequenza **FASTA (text)** ... potete scaricarla (sottoforma di file di testo) dalla sezione dedicata al materiale didattico presente sul sito del laboratorio.

Aprire il file, selezionare tutto e copiare negli appunti di Windows

Ricerca mediante similarità di sequenza

```
>seq_sconosciuta  
ATTCGATCTAGCGATCTA  
CTAATTCGAGGCGATCT  
TCAGCGACTAGCTA  
CGACTACGATCAC
```

Ricerca in una collezione di sequenze di sequenze simili a quella che si usa per effettuare l'interrogazione (**QUERY**)

```
>seq_1  
ATTCGGATCTAGGCTATC  
TAGCGATCGACTGACTAG  
CTAGCTAGCATCGATCAC  
>seq_2  
ATTCGAGCGATCTTTTTA  
TTATATCGGATTTCGATCG  
ATCGATCGACTAAAAAA  
>seq_3  
ATTCGGATCTAGGCTATC  
TAGCGATCGACTGACTAG  
CTAGCTAGCATCGATCAC
```

Punti da ricordare:

- Il risultato è una **lista di sequenze** ordinate dalla più simile alla nostra sequenza QUERY alla meno simile.
- Al posto di una parola chiave utilizziamo una **sequenza**.

Ricerca tramite similarità di sequenza

Apriamo il web browser e colleghiamoci al sito dello strumento di ricerca per similarità di sequenza **BLAST**, **B**asic **L**ocal **A**lignment **S**earch **T**ool (NCBI):

<http://blast.ncbi.nlm.nih.gov/Blast.cgi>

.. nella sezione **Basic Blast**, seguite il link [nucleotide blast](#)

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [Clear](#)

```
>Sequenza_sconosciuta
CTCGCAGGCTCCAGGGGCGGGGCGTGGCCGGGGCGCAGCGACGGGCGGGAGGTCCGGCCGGGCGCGCGC
GCCCCCGCCACACGCACGCCGGGCGTGCCAGTTTATAAAGGGAGAGAGCAAGCAGCGAGTCTTGAAGCTC
TGTTTGGTGCCTTGGATCCATTTCCATCGGTCTTACAGCCGCTCGTCAGACTCCAGCAGCCAAGATGGT
GAAGCAGATCGAGAGCAAGACTGCTTTTCAGGAAGCCTTGGACGCTGCAGGTGATAAACTTGTAGTAGTT
GACTTCTCAGCCACGTGGTGTGGGCCTTGCAAAATGATCAAGCCTTTCTTTCATGATGTTGCTTCAGAGT
GTGAAGTCAAATGCATGCCAACATTCCAGTTTTTTAAGAAGGGACAAAAGGTGGGTGAATTTTCTGGAGC
```

Query subrange [Clear](#)

From

To

Or, upload file

Nessun file selezionato. [Clear](#)

Job Title

Sequenza_sconosciuta

Enter a descriptive title for your BLAST search [Clear](#)

Align two or more sequences [Clear](#)

Choose Search Set

Database

Human genomic + transcript Mouse genomic + transcript Others (nr etc.):

Reference RNA sequences (refseq_rna) [Clear](#)

Organism

Optional

Enter organism name or id--completions will be suggested

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown

Exclude

Optional

Models (XM/XP) Uncultured/environmental sample sequences

Entrez Query

Optional

Enter an Entrez query to limit search [Clear](#)

Program Selection

Optimize for

Highly similar sequences (megablast)

More dissimilar sequences (discontiguous megablast)

Somewhat similar sequences (blastn)

Choose a BLAST algorithm [Clear](#)

Search database Reference RNA sequences

Show results in a new window

Incollate qui la sequenza sconosciuta

Come collezione di sequenze all'interno della quale effettuare la ricerca scegliete RefSeq

Scegliamo il tipo di ricerca che ritorna solo sequenza altamente simili (per velocizzare l'analisi)

Scegliamo di visualizzare i risultati in una nuova finestra e premiamo BLAST

Sequenza_sconosciuta

RID [5DTJ5PTT016](#) (Expires on 10-12 01:39 am)

Query ID |d|10343
Description Sequenza_sconosciuta
Molecule type nucleic acid
Query Length 826

Database Name refseq_ma
Description NCBI Transcript Reference Sequences
Program BLASTN 2.2.28+ [►Citation](#)

Other reports: [►Search Summary](#) [\[Taxonomy reports\]](#) [\[Distance tree of results\]](#)

Otteniamo un output
composto da varie parti

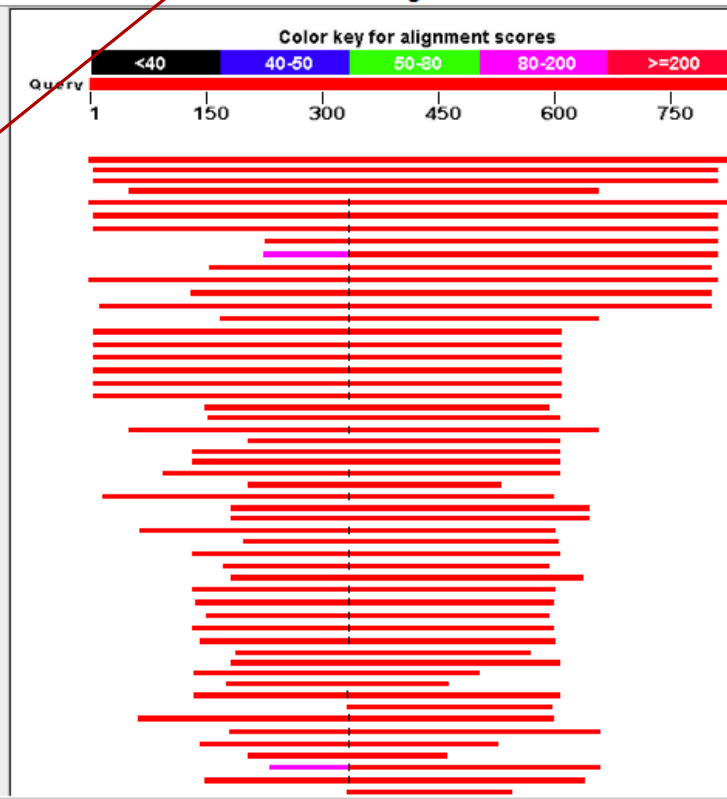
Graphic Summary



Sommario (in
forma grafica)

Distribution of 154 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments



Questa è la **sequenza QUERY** (in italiano sequenza sonda)

Abbiamo ottenuto
154 corrispondenze

Queste sono le
sequenze restituite
dalla ricerca...
La PRIMA (sequenza
più simile alla query) ha
la **stessa lunghezza**
della query

Descrizione sequenze estratte (lista risultati)

Sequences producing significant alignments:

Select: [All](#) [None](#) Selected:0

[Alignments](#) [Download](#) [GenBank](#) [Graphics](#) [Distance tree of results](#)

	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	Homo sapiens thioredoxin (TXN), transcript variant 2, mRNA	1526	1526	100%	0.0	100%	NM_001244938.1
<input type="checkbox"/>	PREDICTED: Gorilla gorilla gorilla thioredoxin, transcript variant 3 (TXN), mRNA	1443	1443	97%	0.0	99%	XM_004048428.1
<input type="checkbox"/>	PREDICTED: Pan troglodytes thioredoxin (TXN), mRNA	1437	1437	97%	0.0	99%	XM_003951491.1
<input type="checkbox"/>	PREDICTED: Macaca fascicularis thioredoxin (TXN), transcript variant X2, mRNA	950	950	73%	0.0	95%	XM_005581101.1
<input type="checkbox"/>	Homo sapiens thioredoxin (TXN), transcript variant 1, mRNA	911	1529	100%	0.0	99%	NM_003329.3
<input type="checkbox"/>	PREDICTED: Gorilla gorilla gorilla thioredoxin, transcript variant 1 (TXN), mRNA	856	1446	97%	0.0	99%	XM_004048426.1
<input type="checkbox"/>	PREDICTED: Pan troglodytes thioredoxin, transcript variant 1 (TXN), mRNA	850	1440	97%	0.0	99%	XM_001142154.2
<input type="checkbox"/>	PREDICTED: Nomascus leucogenys thioredoxin, transcript variant 2 (TXN), mRNA	828	1028	70%	0.0	98%	XM_004087002.1
<input type="checkbox"/>	PREDICTED: Nomascus leucogenys thioredoxin, transcript variant 1 (TXN), mRNA	828	1027	70%	0.0	98%	XM_003260469.2

La prima sequenza della lista copre l'intera lunghezza della sequenza Query ed è anche identica alla sequenza Query.

Abbiamo identificato la sequenza sconosciuta!

Descrizione sequenze estratte (lista risultati)

Sequences producing significant alignments:

Select: [All](#) [None](#) Selected:0

[Alignments](#) [Download](#) [GenBank](#) [Graphics](#) [Distance tree of results](#)

	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	Homo sapiens thioredoxin (TXN), transcript variant 2, mRNA	1526	1526	100%	0.0	100%	NM_001244938.1
<input type="checkbox"/>	PREDICTED: Gorilla gorilla gorilla thioredoxin, transcript variant 3 (TXN), mRNA	1443	1443	97%	0.0	99%	XM_004048428.1
<input type="checkbox"/>	PREDICTED: Gorilla gorilla gorilla thioredoxin, transcript variant 2, mRNA	1437	1437	97%	0.0	99%	XM_003951491.1
<input type="checkbox"/>	PREDICTED: Gorilla gorilla gorilla thioredoxin, transcript variant 1 (TXN), mRNA	856	1446	97%	0.0	99%	XM_004048426.1
<input type="checkbox"/>	PREDICTED: Pan troglodytes thioredoxin, transcript variant 1 (TXN), mRNA	850	1440	97%	0.0	99%	XM_001142154.2
<input type="checkbox"/>	PREDICTED: Nomascus leucogenys thioredoxin, transcript variant 2 (TXN), mRNA	828	1028	70%	0.0	98%	XM_004087002.1
<input type="checkbox"/>	PREDICTED: Nomascus leucogenys thioredoxin, transcript variant 1 (TXN), mRNA	828	1027	70%	0.0	98%	XM_003260469.2

Link ai risultati del confronto tra la query e la prima sequenza ottenuta

Link alla ENTRY della prima sequenza ottenuta

Ora abbiamo a disposizione due collegamenti che ci permettono di ottenere ulteriori informazioni. **Seguiamo il link che riporta l'accession (codice identificativo) della sequenza.** In questo esempio è : **NM_001244938.1**

Nucleotide

Nucleotide

Search

[Limits](#) [Advanced](#)[Help](#)

The information on this web site remains accessible; but, due to the lapse in government funding, the information may not be up to date, and the agency may not be able to respond to inquiries until appropriations are enacted. For updates regarding government operating status see USA.gov.

[Display Settings:](#) GenBank[Send:](#)

Homo sapiens thioredoxin (TXN), transcript variant 2, mRNA

NCBI Reference Sequence: NM_001244938.1

[FASTA](#) [Graphics](#)[Go to:](#)

LOCUS	NM_001244938	826 bp	mRNA	linear	PRI 05-OCT-2013
DEFINITION	Homo sapiens thioredoxin (TXN), transcript variant 2, mRNA.				
ACCESSION	NM_001244938				
VERSION	NM_001244938.1 GI:349732255				
KEYWORDS	RefSeq.				
SOURCE	Homo sapiens (human)				

Change region shown

Customize view

Analyze this sequence[Run BLAST](#)[Pick Primers](#)[Highlight Sequence Features](#)[Find in this Sequence](#)

In questo modo abbiamo raggiunto la **ENTRY** del gene TXN (tioredossina) umano. De qui in poi vale tutto quello che abbiamo già visto nella sezione dedicata alla ricerca per parola chiave.

Ricerca per **similarità di sequenza**

Domanda 1 :

- a) A che banca dati appartiene la entry che abbiamo appena estratto (TXN human) ?
- b) Pensate che ci sia una relazione tra la risposta della domanda 1.a e le scelte che avete fatto **PRIMA** di effettuare la ricerca BLAST (suggerimento: riguardate la slide n. 53)? Se si quale?

Esercizio 1 :

Scoprite tutto quello che potete sulla sequenza **Sequenza_sconosciuta_2** presente nello stesso file da cui avete copiato la sequenza, in formato FASTA (text), di Sequenza_Sconosciuta.

- Che tipo di molecola è (DNA o mRNA)?
 - Come si chiama il gene da cui deriva?
 - A quale organismo appartiene questa sequenza?
-

Riepilogo

Le banche dati biologiche sono collezioni di informazioni riguardanti molecole presenti nei viventi. Esse hanno dimensioni considerevoli e quindi vengono rese disponibili al pubblico **unitamente a strumenti specializzati** per svolgere ricerche al loro interno.

Esistono principalmente due tipi di ricerca all'interno di una banca dati biologica:

- Ricerca per **parola chiave** : permette di estrarre una o più sequenze fornendo una serie di parole chiave opportunamente combinate. La stringa di interrogazione può essere costruita dinamicamente grazie ad una serie di filtri progressivi. Non può essere utilizzata in assenza di informazioni sull'obiettivo della nostra ricerca.
- Ricerca per **similarità di sequenza** : al posto delle parole chiave si utilizza una sequenza sonda che serve per trovare le sequenze più simili ad essa in banca dati.

Indipendentemente dalla modalità di ricerca adottata è possibile raggiungere delle schede (entries) che contengono molti collegamenti ad altre informazioni sulla sequenza presenti in altre banche dati