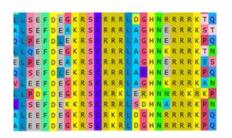
BLAST

- ricerche per similarità in banca dati
- Definizione dell'algoritmo euristico
- Come funziona il processo di allineamento
- Suite di programmi blast: BlastN, BlastP, BlastX



Implementazione del BLAST

- Utilizzo di euristiche per accelerare la ricerca di allineamenti significativi
- Descrizione delle euristiche più comuni (ad esempio, algoritmi seed-and-extend)
- Esempio di applicazione del BLAST

Concetti chiave dell'algoritmo BLAST (Basic Local Alignment Search Tool)

1. Database e Sequenza di Interrogazione:

BLAST confronta una sequenza (query) con un database di sequenze. La query è tipicamente una sequenza di interesse che un ricercatore desidera confrontare con sequenze note per trovare somiglianze.

2. Identificazione di Parole:

BLAST inizia identificando brevi parole di lunghezza fissa (k-mers) nella sequenza query. Queste parole hanno lunghezza 'k' (solitamente 3 per le proteine e 11 per i nucleotidi).

2. Creazione del dizionario:

Dopo l'identificazione delle brevi parole (k-mer) nella sequenza query, BLAST procede con la creazione di un dizionario di parole affini. Questo passaggio coinvolge l'uso di una matrice di punteggio, come BLOSUM62 per le proteine, per identificare e includere nella ricerca quelle parole del database che hanno un grado di somiglianza al di sopra di un punteggio soglia, con i k-mer trovati nella sequenza query.

Questo processo aiuta a identificare possibili corrispondenze significative che non sono esatte, ma sufficientemente simili da meritare ulteriori indagini.

I passaggi 2 e 3 sono importanti per l'efficienza di BLAST, poiché evita di confrontare l'intera sequenza di interrogazione con tutte le sequenze nel database.

4. Estensione:

Una volta identificato un seed, BLAST tenta di estendere questa corrispondenza in entrambe le direzioni per creare un allineamento, fermandosi quando il punteggio dell'allineamento scende sotto una certa soglia. Questo passaggio è cruciale poiché aiuta a identificare somiglianze biologiche significative oltre la corrispondenza iniziale del seed.

5. Punteggio e Valutazione:

Ogni allineamento viene valutato in base al numero di corrispondenze, disallineamenti e spaziature (gaps). Il punteggio viene calcolato usando una matrice di sostituzione che assegna punteggi per la sostituzione degli amminoacidi (nelle sequenze proteiche) o delle sostituzioni dei nucleotidi (nelle sequenze di DNA/RNA).

L'algoritmo calcola anche parametri statistici come il valore E, che aiuta a stimare il numero di corrispondenze che ci si può aspettare di trovare per caso quando si cerca in un database di una determinata dimensione.

6. Filtraggio e Rapporto:

Infine, BLAST filtra i risultati, mantenendo solo quegli allineamenti che soddisfano criteri specifici come punteggio minimo

Cors & Bialinfort matisand 13e corples at a Tiri ana featign per manient rivers of a strict and a surfice of the corp.

L'E-value esprime il numero di volte che ci si aspetta di trovare un allineamento con un determinato punteggio (o superiore) per puro caso, quando si effettua una ricerca contro un database di dimensioni specifiche.

Un E-value basso indica che l'allineamento trovato è altamente significativo, cioè è improbabile che sia avvenuto per caso.

Al contrario, un E-value alto suggerisce che l'allineamento potrebbe non essere significativo e potrebbe essere risultato dal caso.

Ad esempio un E-value di 1 indica che ci si aspetta di trovare una corrispondenza con un punteggio simile per caso una volta in un database delle dimensioni date. Un E-value di 0,01 indica che ci si aspetta una corrispondenza casuale con quel punteggio o superiore solo una volta su 100 ricerche.

6. Filtraggio e Rapporto:

nel filtraggio e nel report dei risultati di BLAST, si mantengono gli allineamenti che hanno un E-value basso, il che indica che sono meno probabili essere il risultato di casualità e quindi più probabili essere biologicamente significativi. Gli allineamenti con un E-value alto, al contrario, sono più probabili essere coincidenze casuali e quindi meno probabili rappresentare una corrispondenza biologicamente rilevante.

6. Filtraggio e Rapporto:

nel filtraggio e nel report dei risultati di BLAST, si mantengono gli allineamenti che hanno un E-value basso, il che indica che sono meno probabili essere il risultato di casualità e quindi più probabili essere biologicamente significativi. Gli allineamenti con un E-value alto, al contrario, sono più probabili essere coincidenze casuali e quindi meno probabili rappresentare una corrispondenza biologicamente rilevante.

Il BLAST è quindi progettato per bilanciare velocità e sensibilità. Ciò viene realizzato concentrandosi sui segmenti piccoli più promettenti (seeds) e estendendoli per trovare allineamenti significativi, piuttosto che eseguire ricerche esaustive, che sarebbero computazionalmente intensive.

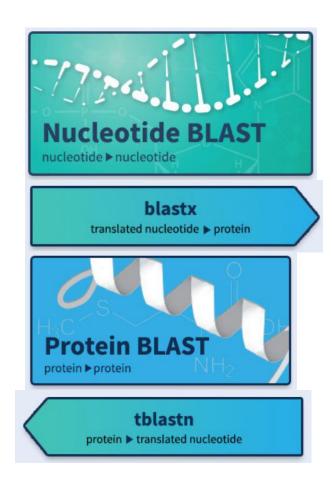
• Suite di programmi blast

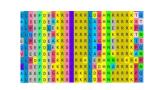
BLASTN programs search nucleotide databases using a nucleotide query.

BLASTP programs search protein databases using a protein query.

BLASTX search protein databases using a translated nucleotide query.

TBLASTN search translated nucleotide databases using a protein query.





Gli algoritmi seed-and-extend sono tra le euristiche più comuni utilizzate nei programmi di allineamento di sequenze, tra cui anche Blast.

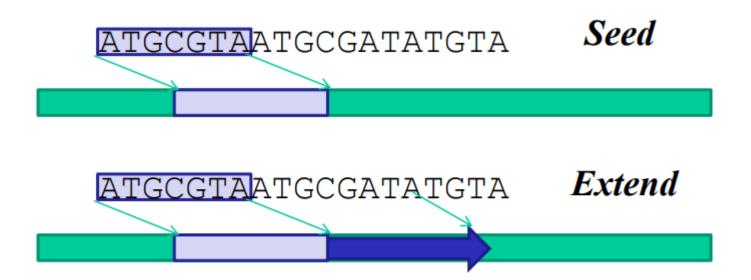
L'idea alla base dell'algoritmo seed-and-extend è di trovare queste sottosequenze conservate, chiamate "semi-omologhe" o "semi-corrispondenze", e utilizzarle per guidare la ricerca degli allineamenti completi.

In particolare, l'algoritmo ricerca le semi-corrispondenze tra la query e il database di riferimento, utilizzando una funzione di scoring per assegnare un punteggio a ciascuna semi-corrispondenza in base alla sua somiglianza con la sequenza di query.

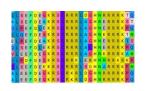
Una volta identificata una semi-corrispondenza, l'algoritmo utilizza la tecnica dell'estensione, ovvero cerca di estendere la semi-corrispondenza per ottenere un allineamento completo. In genere, l'estensione viene eseguita utilizzando un'altra funzione di scoring che tiene conto anche delle lacune (gap) presenti nelle sequenze.

Algoritmi seed-and-extend



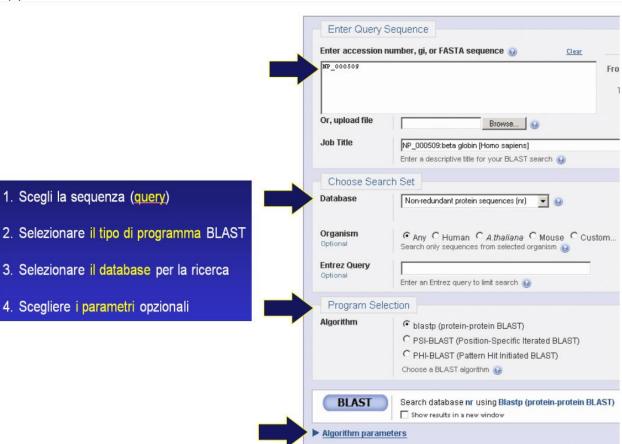


Esempio di applicazione del BLAST



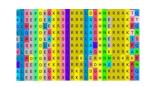
Le quattro fasi di una ricerca BLAST

- 1. Scegli la sequenza (query)
- 2. Selezionare il tipo di programma BLAST
- 3. Selezionare il database per la ricerca
- 4. Scegliere i parametri opzionali Quindi fare clic su "BLAST"





Esempio di applicazione del BLAST



Passo 1: Scelta della sequenza La sequenza può essere inserita in formato FASTA o come accession number (RefSeq)

Esempio di formato FASTA per una query BLAST



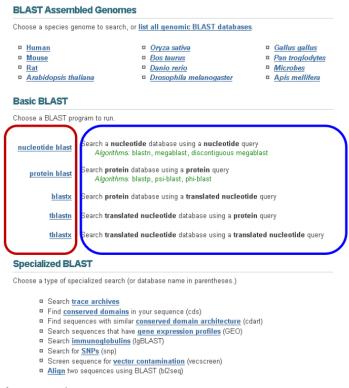
Esempio di applicazione del BLAST

Passo 2:

Scegli il programma BLAST

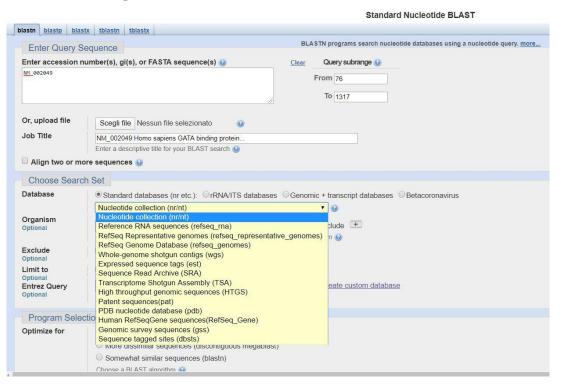
- blastn (nucleotide BLAST)
- blastp (protein BLAST)
- blastx (BLAST tradotto n ▶ P)
- tblastn (BLAST tradotto p N...)
- tblastx (BLAST tradotto n ▶ P... P► N)

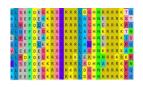
Passo 2: Scegli il programma



Esempio di applicazione del BLAST

Passo 3: scegliere il database





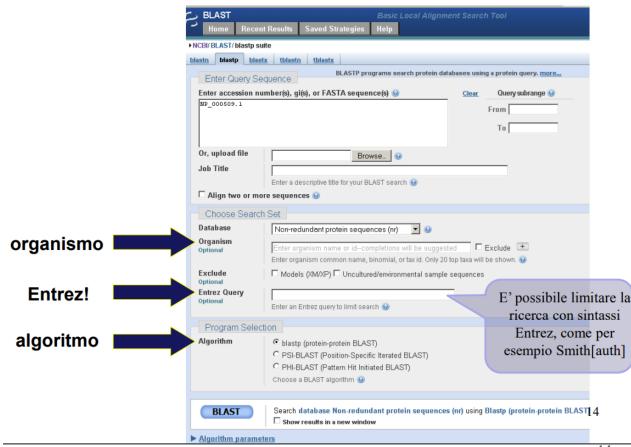
- nr = non ridondante (database più generale, ritorna una sequenza e diversi riferimenti in database in cui la stessa è presente)
- refseq = solo sequenze validate
- -refseq dbest = database di EST
- -dbsts = database di sequenze localizzate gss = genome sequence surveys (BAC, Yac, ecc)
- -Altri... pdb, genomi completi, solo sequenze oggetto di brevetti (pat) ecc. ecc.

Esempio di applicazione del BLAST

Fase 4: parametri opzionali

Si può ...

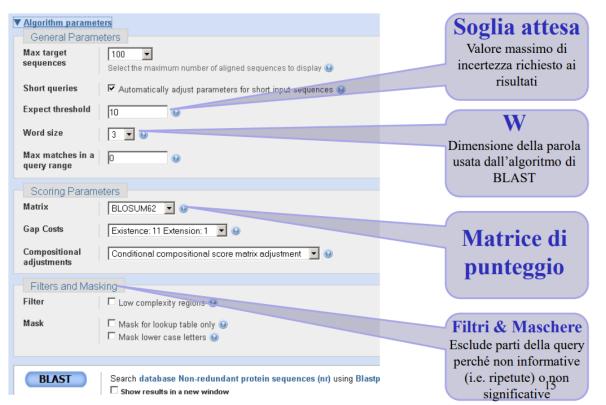
- · scegliere l'organismo di ricerca
- · attivare filtri
- · modificare la matrice di
- punteggio
- · cambiare il valore minimo di affidabilità dei risultati
- modificare la dimensione della parola W

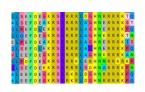


1

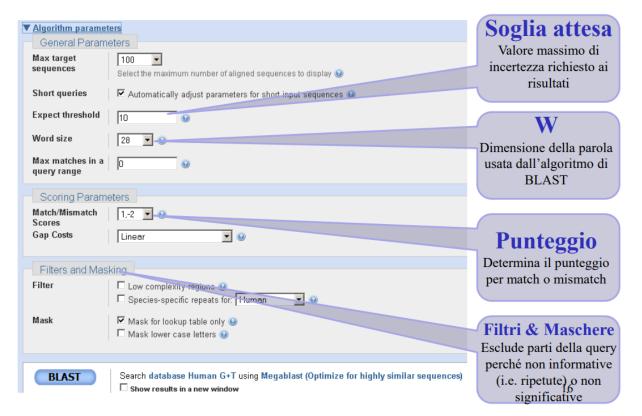
Esempio di applicazione del BLAST

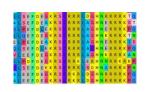
Fase 4a: parametri opzionali di blastp





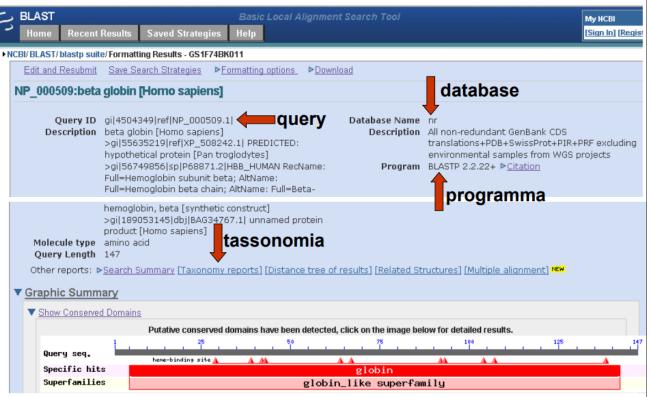


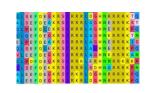




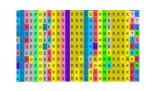
Esempio di applicazione del BLAST

BLAST output: parte superiore

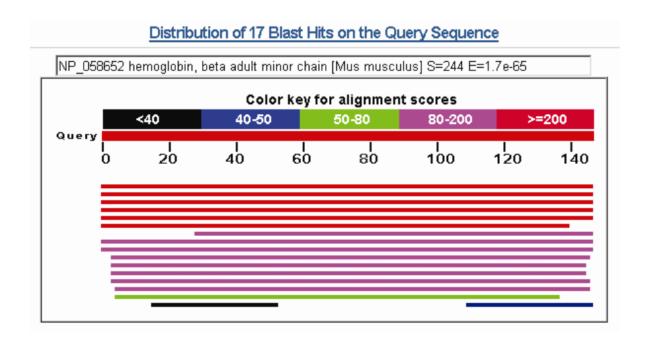




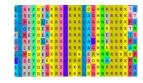
Esempio di applicazione del BLAST



BLAST output: output grafico



Esempio di applicazione del BLAST

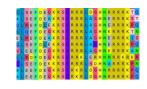


I risultati del Blast conterranno una classifica delle sequenze del database ordinate secondo valori crescenti di **E-value** rispetto alla sequenza query. Il numero di sequenze riportate sarà uguale al numero massimo N impostato nella pagina di input del programma se ci sono almeno N sequenze nella banca dati il cui allenamento con la sequenza query a E-value minore di 1.

Per ognuna delle sequenze della banca dati inclusi nella classifica vengono riportate le seguenti informazioni:

- 1. descrizione della sequenza è relativo numero di accesso della banca dati
- 2. punteggio dell'allenamento con la query
- 3. query coverage, ovvero, trattandosi di allineamenti locali, quale percentuale di residui della query è stata inclusa nell' allineamento locale con la sequenza della banca dati
- 4. E-value associato al punteggio dell'allineamento
- 5. percentuale di identità tra la sequenza e la query

Esempio di applicazione del BLAST

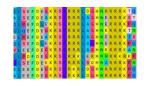


Accession	Description	Max score	Total score	Query coverage	E value	Max ident	Links
NP_000184_I	sonic hedgehog protein preproprotein [Homo sapiens] >gi 6094283 sp Q15465	241	941	100%	0.0	100%	UGM
NP_002172.2	indian hedgehog protein precursor [Homo saplens]	559	459	94%	3e-160	59%	UGM
NP_066382_1	desert hedgehog protein preproprotein [Homo sapiens]	552	442	91%	20:153	56%	UGM
NP 284941.2	mitofusin-1 [Homo sapiens]	31.2	31.2	17%	0.044	29%	UGM
NP.002572.2	pappalysin-1 preproprotein [Homo sapiens]	21.2	31.2	9%	0.048	34%	UGM
NP 001093861,1	RNA-binding Raly-like protein isoform 1 [Homo saplens]	29.6	29.6	7%	0.096	37%	UGM
NP_001093862.1	RNA-binding Raly-like protein isoform 2 [Homo saplens] >ref[NP_001093863.1]	29.6	29.6	7%	0.10	37%	UGM
NP 031393.2	RNA-binding protein Raly isoform 2 [Homo sapiens]	29.6	29.6	7%	0.11	38%	UGM
NP_057951.1	RNA-binding protein Raly isoform 1 [Homo sapiens]	29.6	29.6	7%	0.11	38%	UGM
NP 001013653.1	heterogeneous nuclear ribonucleoprotein C-like 1 [Homo sapiens]	27.7	27.7	2%	0.42	35%	GM
NP_001139653.1	heterogeneous nuclear ribonucleoprotein C-like [Homo sapiens]	22.2	27.7	7%	0.43	35%	UGM
NP_000278.3	peroxisome biogenesis factor 6 [Homo sapiens]	27.7	27.7	12%	0.47	38%	UGM
NP_001130033.2	heterogeneous nuclear ribonucleoprotein C-like [Homo sapiens]	27.3	27.3	7%	0.55	35%	UGM
NP 002580,2	protocadherin-7 isoform a precursor [Homo sapiens]	26.6	26.6	15%	1.2	29%	UGM
NP_115832.1	protocadherin-7 isoform b precursor [Homo sapiens]	25.6	26.6	15%	1.3	29%	GM
NP_001166994.1	protocadherin-7 isoform d precursor (Homo sapiens)	26.6	26.6	15%	1.4	29%	UGM
NP. 115833.2	protocadherin-7 isoform c precursor [Homo sapiens]	25.2	26.2	15%	1.5	29%	UGM
NP 002178.2	interleukin-12 subunit beta precursor [Homo sapiens]	25.8	25.8	20%	1.5	22%	UGM
NP_057715.2	GC-rich sequence DNA-binding factor 1 isoform 1 [Homo sapiens]	25.2	26.2	6%	1.5	50%	UGM
NP_037461.2	GC-rich sequence DNA-binding factor 1 isoform 2 [Homo sapiens]	25.8	25.8	6%	1.8	50%	UGM
NP_275859.3	immunoglobulin superfamily member 22 [Homo sapiens]	25.8	25.8	11%	1.9	29%	UGM
NP_079060.1	zinc finger and BTB domain-containing protein 3 [Homo sapiens]	25.4	25.4	6%	2.8	44%	UGM
NP_570969.2	protein FAM718 [Homo sapiens]	25.0	25.0	9%	3.4	35%	UGM
NP_005215.1	AT-rich interactive domain-containing protein 3A (Homo sapiens)	25.3	24.3	5%	6.3	46%	UGM
NP_689820.2	uncharacterized protein Clorf177 isoform 1 [Homo sapiens]	23.9	23.9	11%	7.0	36%	UGM

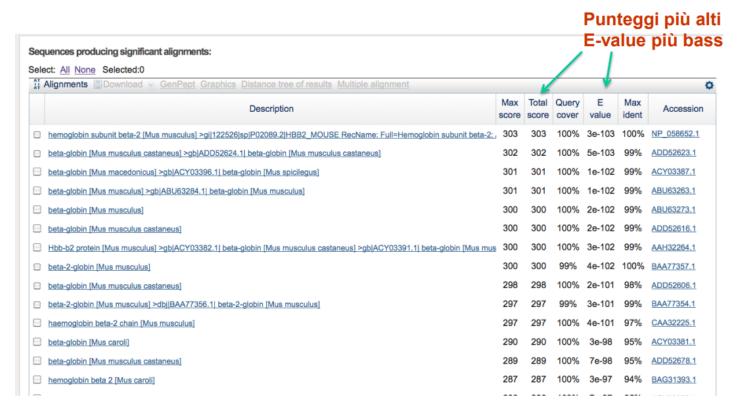
Figura 5.6

Nel cerchio con linea continua gli allineamenti che possiamo considerare significativi perché hanno un *E-value* sufficientemente basso, nel rettangolo con linea tratteggiata quelli non significativi.

Esempio di applicazione del BLAST



BLAST output: valori in dettaglio



Esempio di applicazione del BLAST

Il valore atteso E: è il numero di allineamenti (high scoring segment pairs o HSP) con punteggio maggiore o uguale a un punteggio S che dovrebbero verificarsi per caso in quella ricerca sul database.

Esso rappresenta la stima di probabilità che l'allineamento osservato non sia dovuto al caso. In altre parole si ha una stima del fatto che l'allineamento osservato sia frutto di una storia evolutiva comune tra le sequenze o se sia semplicemente dovuto al caso.

Si può pensarlo come un indice di incertezza.

Un valore E è correlato a un valore di probabilità p.

L'equazione fondamentale che descrive un valore E è:

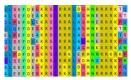
$$E = KMN e^{-\lambda S}$$

Con: M: lunghezza della query

N: lunghezza della sequenza nel database

S: score

 λ ,K: parametri che dipendono dallo scoring system (λ) e dal database usato (K)



- · Il valore di **E decresce esponenzialmente con l'aumentare S**. Valori più elevati di S corrispondono a migliori allineamenti e infatti hanno *E values* più bassi.
- Ottenere un allineamento con **E = 1** significa che esiste un altro allineamento con lo stesso score S che è risultato per caso.
- · La stessa ricerca, su un database più piccolo o più grande, anche se restituisce lo stesso allineamento deve avere un valore di E diverso (ciò dipende da K)

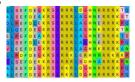
Esempio di applicazione del BLAST

Quindi...

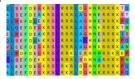
L'E-value è un valore che va da 0 al numero totale di sequenze nel database, ed è ancora una volta tanto più piccolo di 1 quanto migliore è l'allineamento.

L'E-value associato a un allineamento locale di punteggio S ottenuto da ricerca per similarità in banca dati deve essere letto come "il numero atteso di sequenze della banca dati che - per caso- producono un allineamento con la query con punteggio maggiore uguale a S".

L'E-value è inversamente proporzionale al punteggio dell'allineamento: maggiore è il punteggio, più piccolo è l'E-value



Esempio di applicazione del BLAST



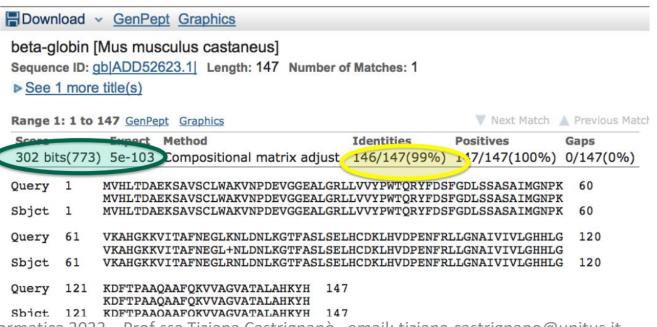
Se ritroviamo nell'output una sequenza per cui l'E-value associato al suo allineamento con la query è >= 1, significa che il risultato che stiamo osservando è quello che ci aspetteremmo per caso.

Quindi, possiamo ipotizzare che l'allineamento che osserviamo non sia indicativo di una storia evolutiva comune ma viceversa sia effetto proprio del caso.

Esaminando i risultati di un allineamento blast possiamo quindi fermarci quando nella classifica incontriamo la prima sequenza il cui E-value è > 1 (E-value>e-5) in quanto le relative sequenze, almeno dal punto di vista statistico, difficilmente avranno un legame evolutivo o funzionale con la query nonostante Blast abbia determinato degli HSP nel loro allineamento con la query.

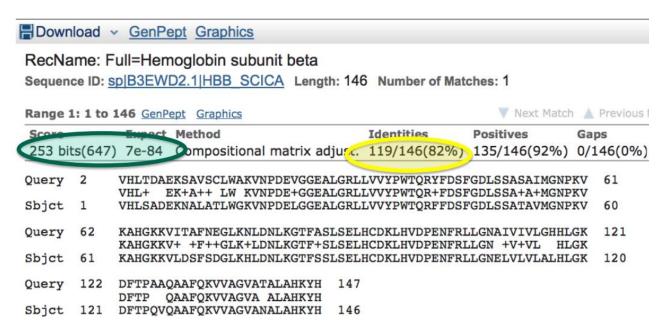
Esempio di applicazione del BLAST

BLAST output: confrontiamo gli allineamenti: seconda hit



Esempio di applicazione del BLAST

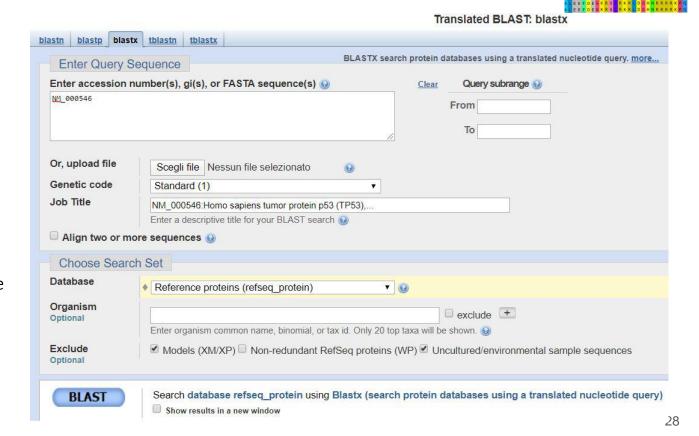
BLAST output: confrontiamo gli allineamenti: ultima hit (su 100)

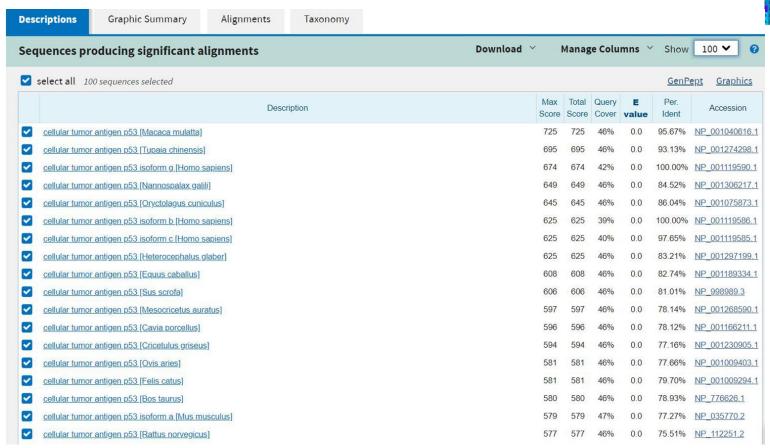


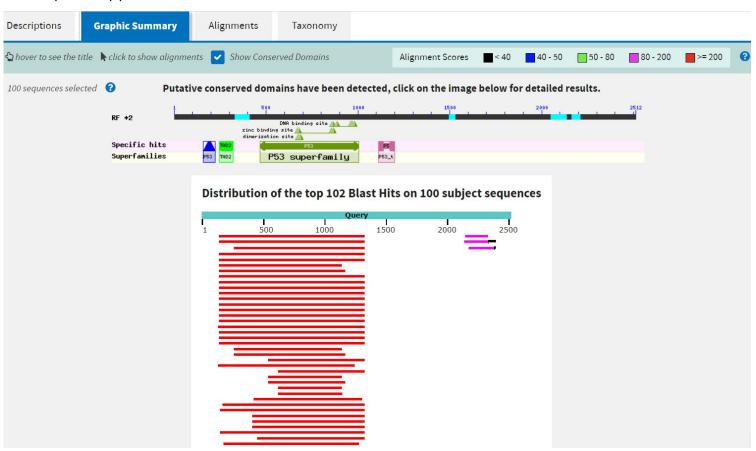
Esempio di applicazione del BLAST

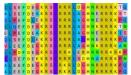
- Andate al sito nell'NCBI: http://www.ncbi.nlm.nih.g
- ☐ Selezionate la banca dati dei nucleotidi e cercate la entry NM_000546

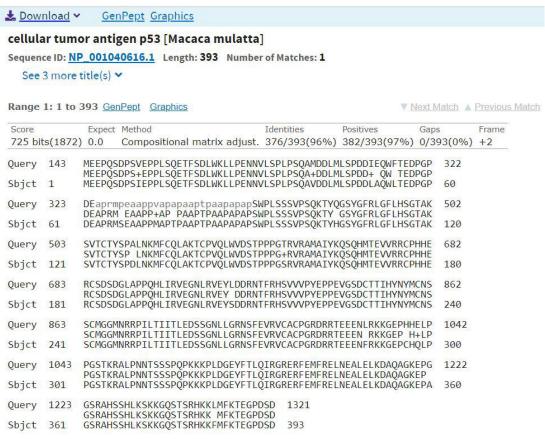
Allineiamo la sequenza NM_000546, contro il database Reference Proteins (refseq_protein), escludendo le sequenze proteiche predette (Models XM/XP) ed i campioni ambientali (Uncultured environmental sequences).

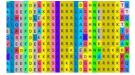












Reports	Lineage Organism	Taxonomy			
100 sequences s	elected 3				
	Organism	Blast Name	Score	Number of Hits	Description
Bilateria		animals		148	
. Deuterostomia		animals		146	
Chordata		chordates		145	
<u>Gnathostomata</u>		<u>vertebrates</u>		143	
<u>Euteleostomi</u>		vertebrates		141	
<u>Tetrapoda</u>		vertebrates		125	
	A <u>mniota</u>	vertebrates		122	
<u>Boreoeutheria</u>		placentals		120	
<u>Euarchontoglires</u>		placentals		105	
	Catarrhini	primates		<u>61</u>	
	Macaca	primates		<u>4</u>	
	Macaca mulatta	primates	725	2	Macaca mulatta hits
	Macaca fascicularis	primates	725	2	Macaca fascicularis hits
	Homo sapiens	primates	674	41	Homo sapiens hits
	Pan paniscus	primates	270	11	Pan paniscus hits
	Gorilla gorilla gorilla	primates	269	2	Gorilla gorilla gorilla hits
	Pan troglodytes	primates	268	3	Pan troglodytes hits
	Tupaia chinensis	placentals	695	1	Tupaia chinensis hits
	Nannospalax galili	rodents	649	2	Nannospalax galili hits
	Oryctolagus cuniculus	rabbits & hares	645	1	Oryctolagus cuniculus hits
	Heterocephalus glaber	rodents	625	1	Heterocephalus glaber hits
	Mesocricetus auratus	rodents	597	1	Mesocricetus auratus hits
	Cavia porcellus	rodents	596	4	Cavia porcellus hits

Esempio di applicazione del BLAST



https://www.youtube.com/watch?v=hTTvq8KJthA

https://www.youtube.com/watch?v=HXEpBnUbAMo

https://www.youtube.com/watch?v=RzC-V67z5LA&t=64s

https://www.youtube.com/watch?v=JKD5laNtwSc&t=76s

E per finire la lista dei tutorial di ncbi:

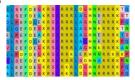
https://www.youtube.com/playlist?list=PLH-TjWpFfWrtjzMCIvUe-YbrIIeFQIKMq

Esempio di applicazione del BLAST



Rispondete alle seguenti domande relative alla entry NM_000546:

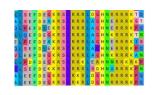
- -Di che tipo di molecola si tratta e quanto è lunga?
- -Quando è stata aggiornata l'ultima volta la entry?
- -Cosa dice la sua definizione?
- -Da che organismo proviene?
- -Quale gene produce questa molecola? Questo gene ha nomi alternativi?
- E' un trascritto codificante? Se sì, dove inizia la 5' UTR? E la 3' UTR?
- Quanto sono lunghe? E la
- CDS invece quanto è lunga? E la proteina codificata?
- Riesci a trovare un modo per recuperare la sequenza di
- NM_000546 in formato FASTA?
- Su che cromosoma si trova il gene che produce questo trascritto?



- -Cerca di capire se esistono malattie umane note associate a mutazioni di questo gene.
- -Procedi alla entry della proteina prodotta da questo trascritto.
- -Qual è identificativo (Accession, o ID) della proteina? La sua definizione?
- -Esistono dei siti noti di fosforilazione annotati su questa proteina? In che posizione?
- -Ottieni la sequenza in formato FASTA della proteina.
- -Vai alla entry del gene che codifica questa proteina.
- -Su che filamento è annotato? Ci sono altri geni nelle vicinanze?
- -Cerca di capire se il gene produce altri trascritti oltre a quello che codifica per questa proteina.

ALLINEAMENTI TRA SEQUENZE

Limiti e sfide di Blast:

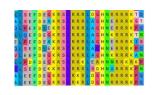


Pur essendo uno degli strumenti di allineamento di sequenze più utilizzati in bioinformatica ci sono alcune limitazioni e sfide nell'uso di Blast che possono influenzare i risultati dell'allineamento.

- 1. Sequenze altamente **simili**: BLAST può avere difficoltà a distinguere tra sequenze altamente simili, poiché le sequenze condividono molte regioni identiche. In tali casi, BLAST potrebbe non essere in grado di distinguere tra le due sequenze o di identificare le regioni di divergenza.
- 2. Lunghezza della sequenza: la lunghezza della sequenza può influenzare la precisione di BLAST. Sequenze troppo corte possono non fornire informazioni sufficienti per l'allineamento corretto, mentre sequenze troppo lunghe possono richiedere troppo tempo per l'allineamento.

ALLINEAMENTI TRA SEQUENZE

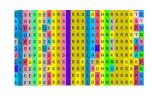
Limiti e sfide di Blast:



- 3. **Sovrapposizione** di domini proteici: BLAST può avere difficoltà a rilevare le sovrapposizioni tra i domini proteici all'interno di una singola sequenza. In tali casi, BLAST può non essere in grado di identificare la corretta struttura proteica o di riconoscere le relazioni di omologia tra sequenze.
- 4. **Bias** di composizione: la composizione nucleotidica può influenzare la precisione di BLAST. Sequenze con una forte composizione di un particolare nucleotide possono produrre falsi positivi o falsi negativi nell'allineamento.

ALLINEAMENTI TRA SEQUENZE

Limiti e sfide di Blast:



- 5. Selezione del **database** di riferimento: la scelta del database di riferimento può influenzare la capacità di BLAST di identificare relazioni di omologia tra le sequenze. Un database di riferimento troppo ampio o troppo piccolo può compromettere la precisione dell'allineamento.
- 6. Scalabilità: BLAST può essere limitato nella sua capacità di gestire grandi quantità di sequenze, soprattutto quando si tratta di sequenze di dimensioni diverse.
- Tuttavia, nonostante queste sfide, BLAST rimane uno strumento utile per l'analisi di sequenze e l'identificazione di relazioni di omologia tra le sequenze.