# Lezione 8

## Ricerche in banche dati (databases) attraverso l'uso di BLAST

# BLAST: Basic Local Alignment Search Tool

Basic Local Alignment Search Tool.  Altschul et al. 1990,1994,1997

- Sviluppato per rendere ancora più veloci le ricerche nelle banche dati rispetto a FASTA, senza perdere in sensibilità e selettività

- Metodo euristico per allineamenti locali

- Pensato specificamente per ricerche in database

- Basato sulle stesse assunzioni di FASTA: un buon allineamento contiene corti frammenti di match esatti

# BLAST: Basic Local Alignment Search Tool

Basic Local Alignment Search Tool. Altschul et al. 1990,1994,1997

- Input:
  - Query sequence Q (la vostra sequenza!)
  - Database of sequences DB
  - Minimal score S

- Output:
  - Sequenze presenti nel DB (Seq), per le quali Q e Seq abbiano uno score > S

# BLAST Fundamentals

- BLAST tells you about non-chance similarities between biological sequences.
- If similarities are not due chance then they must be due to something else!
  - Homology
  - Simple identification
- All BLAST searches begin with a sequence
  - protein or nucleotide
  - experimentally determined or one from database

NCBI Webinars

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

# What BLAST tells you

Here's my sequence.

1. What is it related to? (What does it do?)
   – Homology
   – Function
2. Is it already in the database? (Identification)
   – find the matching sequence in the database
   – organism of origin
3. Where is it located or how is it organized?
   – in a genome
   – other annotation problems
     - comparing sequences
     - looking for frame shifts

# BLAST Word Matching

```
MEAAVKEEISVEDEAVDKNI

MEA
 EAA
  AAV
   AVK
    VKE
     KEE
      EEI
       EIS
        ISV
         ...
```
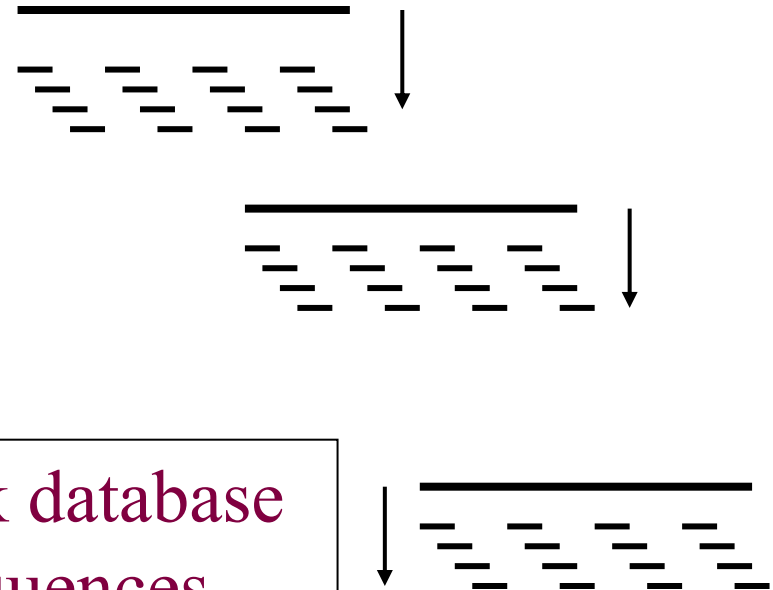
Break query into words:

Break database sequences into words:

Alignment starts with initial word of 11

**ACACTGAGTGA**
| | | | | | | | | | |
**ACACTGAGTGA**

Extension to the left has no mismatches, no penalty points
Extension to the right has mismatches and penalty points

GCACCTTTGCC**ACACTGAGTGA**GCTGCTCTATG
| | | | | | | | | | | | | | | | | | | | | | |   | | | |   | |   |
GCACCTTTGCC**ACACTGAGTGA**CCTGCACTGTA

Extension to the left has no penalty points and can continue to grow
Extension to the right accumulates too many mismatch penalty points; extension in this direction stops

CAACCTCAAGGGCACCTTTGCC**ACACTGAGTGA**GCTGCTCTATGGTCCTTTGGGG
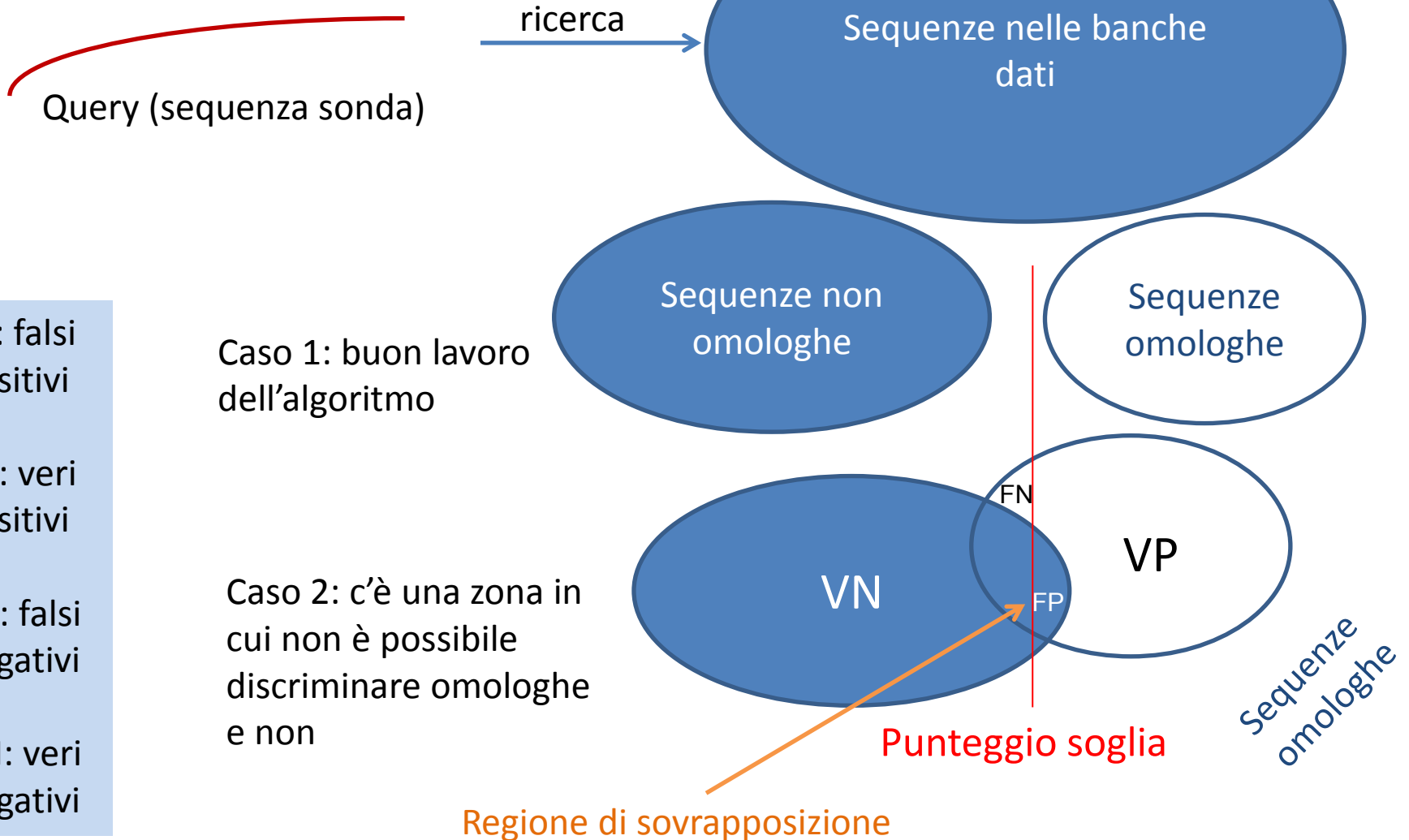| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |   | | | |   | |   |        | | | |
CAACCTCAAGGGCACCTTTGCC**ACACTGAGTGA**CCTGCACTGTAAAGTTTTGCAT

If left side cannot grow any more, the final alignment looks like this:

CAACCTCAAGGGCACCTTTGCC**ACACTGAGTGA**GCTGCTCTATG
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |   | | | |   | |   |
CAACCTCAAGGGCACCTTTGCC**ACACTGAGTGA**CCTGCACTGTA

**Figure 3.1 Simple extension example for BLASTN.** Starting with an initial match of "words," BLAST extends the alignment between query and hit, keeping track of penalty points against, and increasing significance for, extending the alignment.

# Ricerche in database

L'algoritmo deve identificare le sequenze omologhe e non omologhe separate da un valore soglia

Query (sequenza sonda)

ricerca →

Sequenze nelle banche dati

Sequenze non omologhe

Sequenze omologhe

Caso 1: buon lavoro dell'algoritmo

FP: falsi positivi

VP: veri positivi

FN: falsi negativi

VN: veri negativi

Caso 2: c'è una zona in cui non è possibile discriminare omologhe e non

VN

FN

VP

FP

Sequenze omologhe

Punteggio soglia

Regione di sovrapposizione

# BLAST Statistics

```
Score = 18.5 bits (36),  Expect = 47992
Identities = 5/5 (100%), Positives = 5/5 (100%), Gaps = 0/5 (0%)

Query   1     ELVIS   5
              ELVIS
Sbjct   8     ELVIS   12
```

- Number of chance alignments = 48 thousand!
- Indistinguishable from chance

**The most important statistic: Expect value (e-value)**
Expected number of random alignments with a particular score or better

```
Score = 89.7 bits (204),  Expect = 7e-18
Identities = 50/103 (49%), Positives = 54/103 (52%), Gaps = 18/103 (17%)

Query   1     MKLLAATVL---LLTICSLEGALVR
              MK L    VL    LL +CSLEGA V
Sbjct   1     MKVL---VLAMVLLCVCSLEGAVVM

Query   54    SPELQAEAKSYFEKSKEQLTPLIKKAGTELVNFLSYFVELGTQ   96
              E    +AK Y E    EQ  P  K    TE         F +L TQ
Sbjct   5
```

- Number of chance alignments = $7 \times 10^{-18}$
- Not due to chance

- The e-value depends directly on the size of the search space (database)
- Search the smallest database likely to contain the sequence of interest

Attesa (**E**xpectation) di trovare PER CASO uno Score come quello osservato

# Scoring: Nucleotide

Number of Chance Alignments = $2 \times 10^{-73}$

```
Score =   288 bits (318),  Expect = 2e-73
Identities = 262/325 (81%), Gaps = 8/325 (2%)
Strand=Plus/Plus

Query  1923   TCAGCCTACCATGAGAATAAGAGAAAGA-AAATGAAGATCAAAAGCTTATTCATCTGTTT   1981
              ||||  |||| ||||||||||||||||| ||||||||| |  | ||||| |||| |||
Sbjct  33774  TCAGACTACCCTGAGAATAAGAGAAAGAGAAATGAAGACCTAGA-CTTATCCATCTCTTT   33832

Query  1982   TTCTTTTTCGTTGGTGTAAAGCCAACACCCTGTCTAAAAAACATAAATTTCTTTAATCAT   2041
              ||  |||| |    ||||||||||||||     ||||   ACAAATTTCTTTAAATAT
Sbjct        TGC                                                          33892
```

Match=+2

Mismatch=-3

```
Query  2042   TTTGCCTCTTTTC                        AGAATCTAATAGAGTGGT   2100
              |||||||||||||                        |||||| || |
Sbjct  33893  TTTGCCTCTTTTCTCTGTGCTACAATTAATAAAAAAATGAAAAGAATCTAATTTAATTGT   33952

Query  2101   ACAGCACTGTTA-T                           GGTTCTGTGG   2159
              | |||||||||                              ||||| ||
Sbjct  33953  CTATGACTGTTATT                           GGTTCTATGA   34012
```

Gap

$-(5 + 4(2)) = -13$

```
Query  2160   AAGTTCCAGTGTTC                       TGTGGGCTA   2219
              || |||||||||                         | ||| ||
Sbjct  34013  AAATTCCACTATTCTCTCTTTCCCTATTTCAATGGAGGACTTCTAGTTCCTTCTGGATTA   34072

Query  2220   AT----TAAATAAATCATTAATACT   2240
              ||    ||||  || |||||||||
Sbjct  34073  ATTGCATAAAAGAAACATTAATACT   34097
```

# Scoring: Protein

Number of Chance Alignments = 4 X $10^{-50}$

Score =  176 bits (447),  Expect = 4e-50, Method: Compositional matrix adjust.
Identities = 98/232 (42%), Positives = 139/232 (60%), Gaps = 14/232 (6%)

```
Query  30   MAKVLTLELYKKLRDKETPSGFTVDDVIQTGV--DNPGHPFIMTVGCVAGDEESYEVFKE  87
            + K LT +L+++ +D+     GF+     I +G    N G       VG  AG  +SY  F
Sbjct  26   LQKCLTKDLWEQCKDRRDKYGFSFKQAIFSGSKWTNSG------VGVYAGSHDSYYAFAP  79

Query  88   LFDPIISDRHGGYKPTDKHKTDLNHENLKGG---DDLDPNYVLSSRVRTGRSIKYTLPP  144
```

| | | | | |
|K| |D| DKH |Q| D D  + S+R+R |D| 137 |
|FMD| | | SDKHIS | | PADED-KMINSTRIRVA | | |
|K +5| |HCS E +2| RALNSI F −3 | SMTEKEQQQLIDDHFLR E +2 |204|
| | | | AL | | +M++ E++QLI DHFLR | | |

```
Sbjct  138  AVTRKERKEIEHLVTSALGEFTGELKGKYKS              196

Query  205  SGMARDWPDARGIWHNDNKSFLVWVNEED
            +G+ RDWP+ARGI+HND K+FLVWVNEED
Sbjct  197  AGLERDWPEARGIFHNDAKTFLVWVNEED
```

Gap
$-(11 + 4(1)) = -14$

NCBI Webinars

Scores from **BLOSUM62**, a position independent matrix
 – Same substitution gets the same score at all positions
 – All positions equally likely to change

# E value: significatività statistica

Non si interpretano come p values dove

$$p < 0.05$$

sono generalmente considerati significativi

**Regola generale**

E values $< 10^{-6}$ sono molto probabilmente significativi.

$10^{-6} <$ E values $< 10^{-3}$ meritano una seconda occhiata.

E values $< 10^{-3}$ andrebbero scartati (ci aspettiamo di trovare 0.001 sequenze non correlate alla nostra-falsi positivi- che ottengono un punteggio superiore a quell'S).

# BLAST Programs

## BLAST has five programs

Differ in the types of sequences they align and at what level

| Program | Query Seq. Type | Database Seq. Type | Alignment Level |
|---------|-----------------|---------------------|-----------------|
| blastn | nucleotide | nucleotide | nucleotide |
| blastp | protein | protein | protein |
| blastx | nucleotide | protein | protein |
| tblastn | protein | nucleotide | protein |
| tblastx | nucleotide | nucleotide | protein |

**Six-frame translation**

# BLAST Homepage

blast.ncbi.nlm.nih.gov

## BLAST Assembled Genomes

Find Genomic BLAST pages:

[Enter organism name or id--completions will be suggested] **GO**

- Human
- Mouse
- Rat
- Cow
- Pig
- Dog
- Rabbit
- Chimp
- Guinea pig
- Fruit fly
- Honey bee
- Chicken
- Zebrafish
- Clawed frog
- Arabidopsis
- Rice
- Yeast
- Microbes

## Basic BLAST

Choose a BLAST program to run.

| | |
|---|---|
| nucleotide blast | Search a **nucleotide** database using a **nucleotide** query<br>*Algorithms: blastn, megablast, discontiguous megablast* |
| protein blast | Search **protein** database using a **protein** query<br>*Algorithms: blastp, psi-blast, phi-blast, delta-blast* |
| blastx | Search **protein** database using a **translated nucleotide** query |
| tblastn | Search **translated nucleotide** database using a **protein** query |
| tblastx | Search **translated nucleotide** database using a **translated nucleotide** query |

## Specialized BLAST

Choose a type of specialized search (or database name in parentheses.)

- Make specific primers with **Primer-BLAST**
- Cluster multiple sequences together with their database neighbors using **MOLE-BLAST**
- Find **conserved domains** in your sequence (cds)
- Find sequences with similar **conserved domain architecture** (cdart)
- Search sequences that have **gene expression profiles** (GEO)
- Search **immunoglobulins and T cell receptor sequences** (IgBLAST)
- Screen sequence for **vector contamination** (vecscreen)
- **Align** two (or more) sequences using BLAST (bl2seq)
- Search **protein** or **nucleotide** targets in PubChem BioAssay
- Search **SRA by experiment**
- Constraint Based Protein **Multiple Alignment Tool**
- Needleman-Wunsch **Global Sequence Alignment Tool**
- Search **RefSeqGene**
- Search **trace archives**
- Search bacterial and fungal rRNA sequences with **Targeted Loci BLAST**

# Nucleotide Databases

## Choose Search Set

**Database**

Genomic plus Transcript
    Human genomic plus transcript (Human G+T)
    Mouse genomic plus transcript (Mouse G+T)

⦿ Others (nr etc.):

**Organism**
Optional

Other Databases
✓ **Nucleotide collection (nr/nt)**
    Reference RNA sequences (refseq_rna)
    Reference genomic sequences (refseq_genomic)

lude  +
ill be shown

**Exclude**
Optional

nces

**Limit to**
Optional

**Entrez Query**
Optional

    NCBI Genomes (chromosome)
    Expressed sequence tags (est)
    Genomic survey sequences (gss)
    High throughput genomic sequences (HTGS)
    Patent sequences(pat)
    Protein Data Bank (pdb)
    Human ALU repeat elements (alu_repeats)
    Sequence tagged sites (dbsts)
    Whole-genome shotgun contigs (wgs)
    Transcriptome Shotgun Assembly (TSA)
    16S ribosomal RNA sequences (Bacteria and Archaea)
    Sequence Read Archive (SRA)

eate custom database

Services
megablast
blastn
tblastn
tblastx

# Nucleotide Databases

- Default database (nr/nt) is <u>not</u> comprehensive
  - Traditional GenBank and RefSeq RNA
  - Useful subsets: RefSeq RNA, 16S rRNA reference sequences

- What is <u>not</u> in nr/nt? the majority of nucleotide data
  - Bulk sequences (EST, GSS, HTGS, STS)
  - RefSeq Genomic Sequences (Chromosome, RefSeq Genomic)
  - US, European and Asian Patents (pat)
  - Whole Genome Shotgun Contigs (WGS) (Second Largest)
  - Transcriptome Shotgun Assemblies (TSA)
  - Next-Gen Reads (SRA) (Largest set of data)

NCBI Webinars

Ricordiamo che l'efficienza della ricerca aumenta se limitiamo il database che interroghiamo

# Limiting Databases

Search the smallest database likely to contain the sequence of interest.

**Choose Search Set**

**Database**    Non–redundant protein sequences (nr)

**Organism**
Optional    bacteria (taxid:2)    ☐ Exclude    +

Enterobacteriales (taxid:91347)    ☑ Exclude

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

Organism limit

Exclude    ☐ Models (XM/XP) ☑ Uncultured/environmental sample sequences
Optional

**Entrez Query**    25000:30000[Molecular Weight]
Optional
Enter an Entrez query to limit search

Exclude predicted and uncultured

Limit with Entrez query

NCBI/ BLAST Home

BLAST finds regions of similarity between biological sequences.   more...

**New** Try SmartBLAST for an improved protein-protein search

## BLAST Assembled Genomes

Find Genomic BLAST pages:

[Enter organism name or id--completions will be suggested]   GO

- Human
- Mouse
- Rat
- Cow
- Pig
- Dog

- Rabbit
- Chimp
- Guinea pig
- Fruit fly
- Honey bee
- Chicken

- Zebrafish
- Clawed frog
- *Arabidopsis*
- Rice
- Yeast
- Microbes

# Genome Databases

Shortcuts to popular organisms

- Comprehensive search for genomic data
- Finds the best set (most assembled) of genomic sequences

ng a **protein** query
ast, phi-blast, delta-blast

ng a **translated nucleotide** query

e database using a **protein** query

e database using a **translated nucleotide** query

name in parentheses )

NCBI Webinars

# Algorithm Parameters: General



General Parameters

| | |
|---|---|
| **Max target sequences** | 100 ⬍ |
| | Select the n ... umber of aligned sequences to display 🔵 |
| **Short queries** | ☑ Automa...ust parameters for short input sequences 🔵 |
| **Expect threshold** | 10 🔵 |

Dropdown options:
```
  10
  50
✓ 100
  250
  500
  1000
  5000
  10000
  20000
```

- Increase Max target sequences
- Decrease Expect threshold

Set to more stringent value:
- 1e-6
- 0.001

**Let Expect threshold govern output not Max target sequences**

Threshold = soglia (vedi diapositiva 5)

**BLAST**   Search database Reference proteins (refseq_protein) using Blastp (protein-protein BLAST)
☐ Show results in a new window

⊟ Algorithm parameters
Note: Parameter values that differ from the default are highlighted in yellow and marked with ♦ sign

### General Parameters

**Max target sequences**
`100` ▼
Select the maximum number of aligned sequences to display ⓦ

**Short queries**
☑ Automatically adjust parameters for short input sequences ⓦ

**Expect threshold**
`10` ⓦ

**Word size**
`3` ▼ ⓦ

**Max matches in a query range**
`0` ⓦ

Verranno presentate tutte le hits (sequenze trovate) sotto questa soglia di E values (cioè con E < 10)
https://www.youtube.com/watch?v=nO0wJgZRZJs

### Scoring Parameters

**Matrix**
`BLOSUM62` ▼ ⓦ

**Gap Costs**
`Existence: 11 Extension: 1` ▼ ⓦ

**Compositional adjustments**
`Conditional compositional score matrix adjustment` ▼ ⓦ

### Filters and Masking

**Filter**
☐ Low complexity regions ⓦ

**Mask**
☐ Mask for lookup table only ⓦ
☐ Mask lower case letters ⓦ

**BLAST**   Search database Reference proteins (refseq_protein) using Blastp (protein-protein BLAST)
☐ Show results in a new window

Ricordiamo che l'E risponde alla domanda: quante sequenze mi aspetto che abbiano **per caso** uno score maggiore o uguale a quello che ho osservato (falsi positivi!)

Questo filtro è importante: permette di effettuare ricerche escludendo regioni con molte ripetizioni come omopolimeri

Dopo aver deciso se cerchiamo nucleotidi contro nucleotidi, proteine contro proteine etc, possiamo anche decidere in che specifico db cercare, ad esempio **Refseq**

# NCBI RefSeq Database

- *Goal:* Provide a single reference sequence for each molecule of the central dogma (DNA, mRNA, and protein)

- Distinguishing features

  - Non-redundancy
  - Updates to reflect the current knowledge of sequence data and biology
  - Includes biological attributes of the gene, gene transcript, or protein
  - Encompasses a wide taxonomic range, with primary focus on mammalian and human species
  - Ongoing updates and curation (both automated and manual review), with review status indicated on each record

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# RefSeq Accession Number Prefixes

*From curation of GenBank entries:*

| | |
|---|---|
| **NT_** | Genomic contigs |
| **NM_** | mRNAs |
| **NP_** | Proteins |
| **NR_** | Non-coding transcripts |

*From genome annotation:*

| | |
|---|---|
| **XM_** | Model mRNA |
| **XP_** | Model proteins |

Complete list of molecule types in Chapter 18 of the NCBI Handbook
*http://ncbi.nlm.nih.gov/books/NBK21091*

E' possibile limitare la ricerca ad uno specifico gruppo tassonomico o ad uno specifico organismo

E' possibile definire specifici parametri per la ricerca

**BLAST** — Search database Reference proteins (refseq_protein) using Blastp (protein-protein BLAST)
☐ Show results in a new window

⊟ Algorithm parameters — Note: Parameter values that differ from the default are highlighted in yellow and marked with ♦ sign

**General Parameters**

Max target sequences — 100 ▾ — Select the maximum number of aligned sequences to display

Short queries — ☑ Automatically adjust parameters for short input sequences

Expect threshold — 10

Word size — 3 ▾

Max matches in a query range — 0

Verranno presentate tutte le hits (sequenze trovate) sotto questa soglia di E values (cioè con E < 10)
https://www.youtube.com/watch?v=nO0wJgZRZJs

**Scoring Parameters**

Matrix — BLOSUM62 ▾

Gap Costs — Existence: 11 Extension: 1 ▾

Compositional adjustments — Conditional compositional score matrix adjustment ▾

**Filters and Masking**

Filter — ☐ Low complexity regions

Mask — ☐ Mask for lookup table only
☐ Mask lower case letters

**BLAST** — Search database Reference proteins (refseq_protein) using Blastp (protein-protein BLAST)
☐ Show results in a new window

Ricordiamo che l'E risponde alla domanda: quante sequenze mi aspetto che abbiano **per caso** uno score maggiore o uguale a quello che ho osservato (falsi positivi!)

**BLAST** Search database Reference proteins (refseq_protein) using Blastp (protein-protein BLAST)
☐ Show results in a new window

⊟Algorithm parameters                     Note: Parameter values that differ from the default are highlighted in yellow and marked with ♦ sign

### General Parameters

**Max target sequences**   `100 ▼`
Select the maximum number of aligned sequences to display ⓪

**Short queries**   ☑ Automatically adjust parameters for short input sequences ⓪

**Expect threshold**   `10`  ⓪

**Word size**   `3 ▼` ⓪   Numero di residui con cui si inizia la ricerca

**Max matches in a query range**   `0`  ⓪

PQG
PEG
PRG
PKG
PNG
PDG
PHG
PMG
PSG
PQA
PQN
etc.

### Scoring Parameters

**Matrix**   `BLOSUM62 ▼` ⓪

**Gap Costs**   `Existence: 11 Extension: 1 ▼` ⓪

**Compositional adjustments**   `Conditional compositional score matrix adjustment ▼` ⓪

### Filters and Masking

**Filter**   ☐ Low complexity regions ⓪

**Mask**   ☐ Mask for lookup table only ⓪
☐ Mask lower case letters ⓪

**BLAST** Search database Reference proteins (refseq_protein) using Blastp (protein-protein BLAST)
☐ Show results in a new window

Ricordiamo che l'E risponde alla domanda: quante sequenze mi aspetto che abbiano **per caso** uno score maggiore o uguale a quello che ho osservato (falsi positivi!)

**BLAST**

Search database Reference proteins (refseq_protein) using Blastp (protein-protein BLAST)
☐ Show results in a new window

⊟ Algorithm parameters — Note: Parameter values that differ from the default are highlighted in yellow and marked with ♦ sign

### General Parameters

| | |
|---|---|
| **Max target sequences** | 100 ▼ Select the maximum number of aligned sequences to display ② |
| **Short queries** | ☑ Automatically adjust parameters for short input sequences ② |
| **Expect threshold** | 10 ② |
| **Word size** | 3 ▼ ② |
| **Max matches in a query range** | 0 ② |

### Scoring Parameters

| | |
|---|---|
| **Matrix** | BLOSUM62 ▼ ② |
| **Gap Costs** | Existence: 11 Extension: 1 ▼ ② |
| **Compositional adjustments** | Conditional compositional score matrix adjustment ▼ ② |

### Filters and Masking

| | |
|---|---|
| **Filter** | ☐ Low complexity regions ② |
| **Mask** | ☐ Mask for lookup table only ② <br> ☐ Mask lower case letters ② |

**BLAST**

Search database Reference proteins (refseq_protein) using Blastp (protein-protein BLAST)
☐ Show results in a new window

Vedi lezioni precedenti per matrice e gap

Questa terza voce permette di controllare per la composizione AA delle sequenze analizzate

Questo filtro è importante: permette di effettuare ricerche escludendo regioni con molte ripetizioni come omopolimeri

# Esercizi con BLAST

- Proviamo ad effettuare una ricerca con le sequenze disponibili nel file

- [BLAST]

- Basic BLAST
  - blastp, creatine kinases
    - COBALT extension
- Genome BLAST
  - blastn, tomato ETR2
    - Potato genome BLAST
    - Formatting options
    - Genome context
- SRA BLAST
  - Potato RNA-Seq
- Primer BLAST
  - BRCA1 Exon Primers
- Microbial Genomes BLAST
  - Chicken Gut 16S
- MOLE-BLAST
  - Clustering Bovine Rumen 16S