



Align Two Sequences Using NCBI BLAST

Creating alignments for pairs of custom sequences

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

National Center for Biotechnology Information • National Library of Medicine • National Institutes of Health • Department of Health and Human Services

Align two (or more) sequences using BLAST

NCBI “Align Two Sequences” service has been fully integrated into the NCBI distributed computing system (*splitd*) [1]. Using “Align Two Sequences”, two groups of sequences can be directly compared. Search results are identified using the assigned unique request ids (RIDs) and are accessible for up to 36 hours. With an NCBI login [2], a search strategy can be saved for future reference. In addition, search results can be displayed or downloaded in various formats using the same RID to highlight different features. The dot-matrix graph presentation is also available if the query and the subject boxes each contains only a single input sequence. The new report format [3] provides additional functions, such as displaying the hits graphically in Sequence Viewer [4] for interactive examination.



Additional search settings

In a BLAST search form, the “Blast 2 sequences” checkbox (A) activates the “Align Two Sequences” function and displays the subject sequence input box (B) while removing the elements pertaining to database selection. The “Align Two Sequences” also adds a new set of parameters for fine tuning searches:

- blastn, megablast or discontinuous megablast algorithms as well as organism-specific repeat filters are available for nucleotide searches (C);
- the “Automatically adjust parameter for short input sequences” settings (D) under the “Algorithm parameters” section is on by default to automatically optimize this type of searches;
- organism-specific repeat filters are available for masking repeat regions in nucleotide searches (E);
- different composition-based statistics can be selected (F) to adjust the significance of the protein sequence alignment.

The screenshot shows the NCBI BLAST search interface. Annotations A-F highlight specific settings:

- A**: "Align two or more sequences" checkbox.
- B**: "Enter Subject Sequence" input box.
- C**: "Optimize for" radio buttons: "Highly similar sequences (megablast)", "More dissimilar sequences (discontinuous megablast)", "Somewhat similar sequences (blastn)".
- D**: "Automatically adjust parameters for short input sequences" checkbox under "Algorithm parameters".
- E**: "Species-specific repeats for: Human" dropdown menu under "Filters and Masking".
- F**: "Conditional compositional score matrix adjustment" option under "Compositional adjustments".

References

1. Bealer K, Coulouris G, Dondoshansky I, Madden T, Merezuk Y, Raytselis Y. A Fault-Tolerant Parallel Scheduler for BLAST. <ftp.ncbi.nlm.nih.gov/blast/documents/blast-sc2004.pdf>
2. My NCBI help manual. www.ncbi.nlm.nih.gov/books/NBK3843/
3. The New BLAST Result Page. ftp.ncbi.nlm.nih.gov/pub/factsheets/Howto_NewBLAST.pdf
4. The Graphical Sequence Viewer. ftp.ncbi.nlm.nih.gov/pub/factsheets/Factsheet_Graphical_SV.pdf

More formatting options

“Align Two Sequences” uses BLAST formatter to display search results. This makes the following possible:

- displaying the alignment in more formats (A, D), such as “Pairwise with identities”;
- downloading alignment (B) in XML, plain text, Hit table and comma delimited format for future reference;
- adding a CDS translation to the nucleotide alignment (C, E) when an accession/gi with CDS feature is used;
- getting more information from other high value NCBI databases through Related Information (F) when sequences, provided as accession or gi, are cited in those databases;
- sorting matched hits in alternative ways using the Description table header when batch subject sequences are used and/or multiple alignment segments are present in the results (not shown);
- on demand dot-matrix display (G) to provide a summary on overall similarity between the two input sequences; and
- graphical display for an alignment (H) in the context of Sequence Viewer.

Click Formatting Options to see the available options.

Check CDS Features to see the translation of ORFs (See alignment display below).

Change the result display back to traditional format

For P12235 vs P12234 alignment, dot-matrix display reveals possible internal repeats in P12234, represented by lines parallel to the main diagonal.

Download GenBank Graphics

Range 1: 225 to 547 GenBank Graphics

| Score | Expect | Identities | Gaps | Strand |
|---------------|--------|--------------|-----------|-----------|
| 560 bits(303) | 7e-164 | 317/324(98%) | 3/324(0%) | Plus/Plus |

Query: 6 GAAGAAGGTTTATNCTTCAGTCCAGGAAGCAATTCACATNCAGCTGTtagagaagcaga
 Subjct: 225T.....T.....GT
 CDS:prefoldin subuni 65 G R M F I L Q S K E A I H S Q L L E K Q

Query: 66 aaatagcagaagaaaaaataaagaactagaacagaaaaagTCCTACTGGAGCGAAG-G 124
 Subjct: 285C.....
 CDS:prefoldin subuni 85 K I A E E K I K E L E Q K K S Y L E R S

Related Information: Gene - associated gene details, UniGene - clustered expressed sequence tags, Map Viewer - aligned genomic context, GEO - microarray expression data

BLAST Results for: (5) - emb|BX648844| (2633 letters)