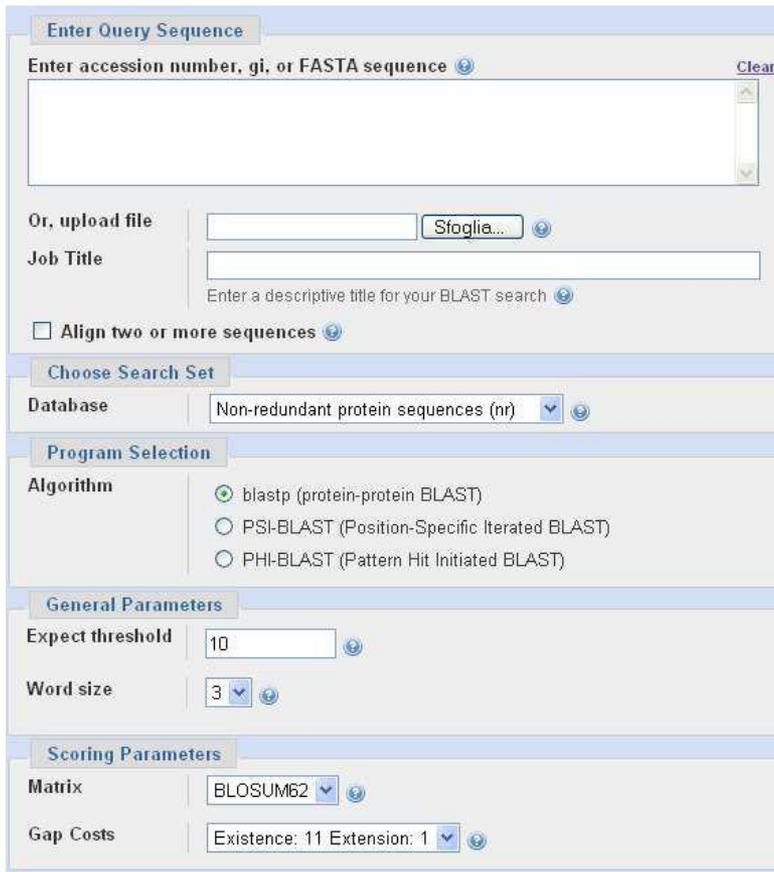


BLAST

26.04.2018

Basic Local Alignment Search Tool (BLAST)

BLAST (Altschul-1990) is an heuristic Pairwise Alignment composed by six-steps that search for local similarities.



The screenshot displays the NCBI BLAST web interface with the following sections and parameters:

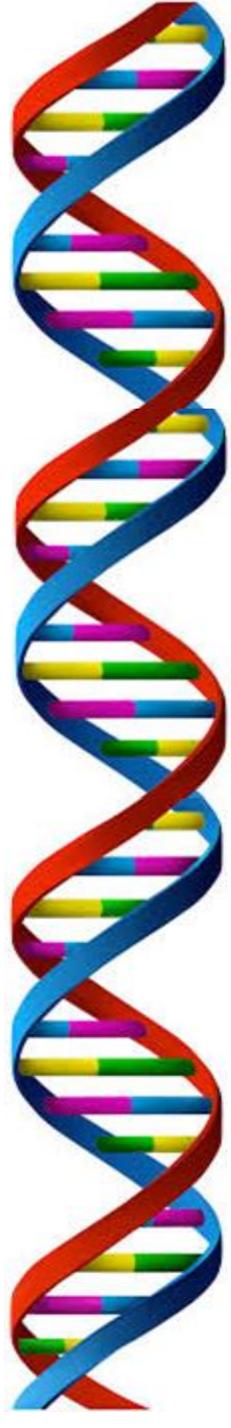
- Enter Query Sequence:** A text input field for "Enter accession number, gi, or FASTA sequence" with a "Clear" button.
- Or, upload file:** A file upload button labeled "Sfoglia...".
- Job Title:** A text input field for "Enter a descriptive title for your BLAST search".
- Align two or more sequences:** An unchecked checkbox.
- Choose Search Set:** A dropdown menu for "Database" set to "Non-redundant protein sequences (nr)".
- Program Selection:** Radio buttons for "Algorithm":
 - blastp (protein-protein BLAST)
 - PSI-BLAST (Position-Specific Iterated BLAST)
 - PHI-BLAST (Pattern Hit Initiated BLAST)
- General Parameters:**
 - "Expect threshold" set to 10.
 - "Word size" set to 3.
- Scoring Parameters:**
 - "Matrix" set to BLOSUM62.
 - "Gap Costs" set to "Existence: 11 Extension: 1".

The most used access point to BLAST is NCBI.

As for the FASTA mask allows set the search parameters and the minimum significance required.

It is also possible to choose the type of BLAST algorithm to use.

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



Basic Local Alignment Search Tool (BLAST)

Phases 1 to 3:

Also BLAST bases its heuristic on indexing, but introduces a sophistication level of the k-tuples, which are also called "word".

It starts by creating a list of words of length W and creation of "w-mers", i.e.. words of length W . W are words that, according a substitution matrix, have a score $>T$ if aligned with the input sequence.

TFDER **LSH**GVQQTFWECIKGD

VSH = 16

ISH = 14

LAH = 13

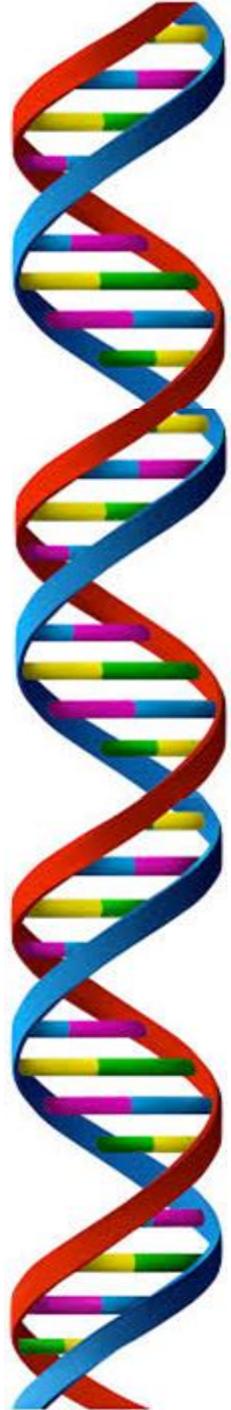
LTH = 13

LSH = 13

Word size (W) = 3

Threshold (T) = 13

For each word they are also generated all possible w-mers, so there are many more words than in FASTA.



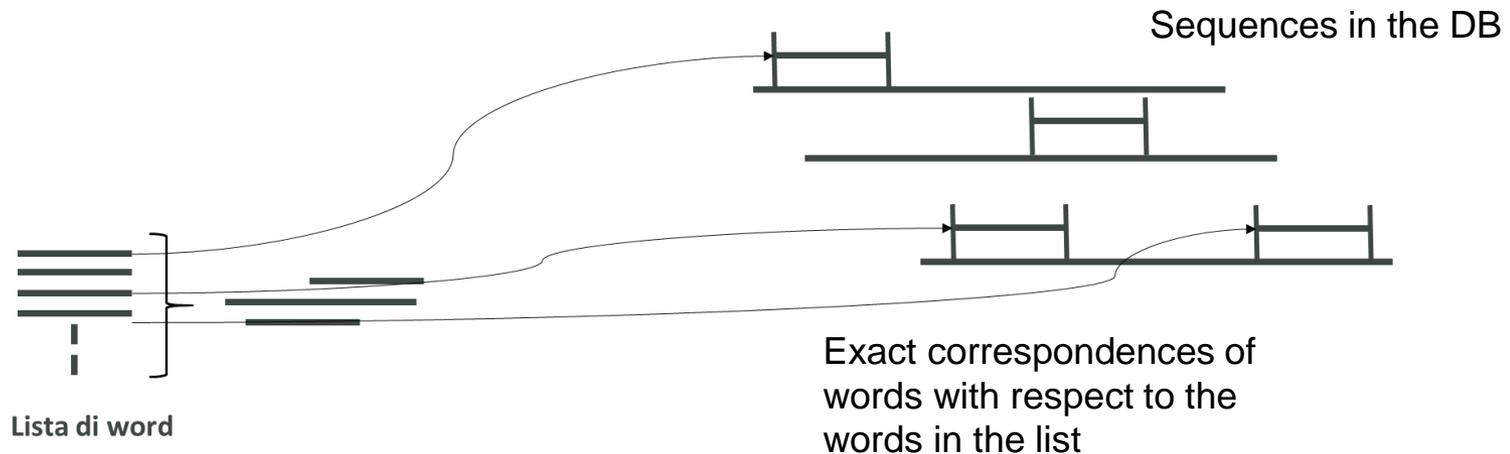
BLAST

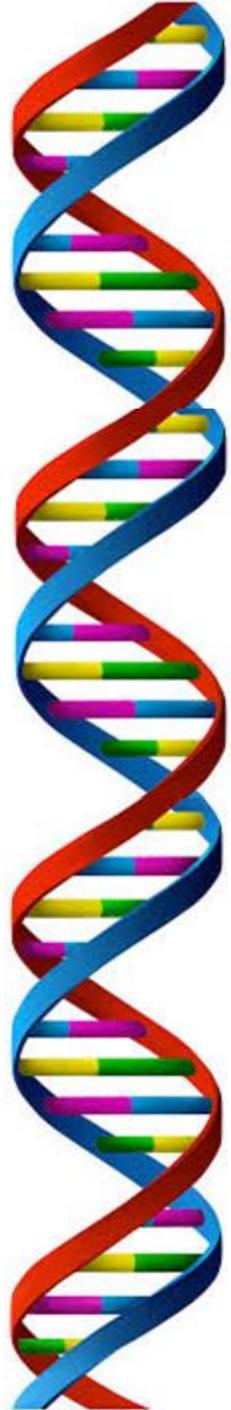
Phase 1:

It is extracted from the query sequence a list of words, each having the same length w . We obtain $L-W+1$ words.

Phase 2:

Each word in the list is compared with sequences of the same length belonging to the all sequences of the considered database.





BLAST

Phase 3:

For each alignment it is assigned a score using a substitution matrix (for example, BLOSUM 62) and is compared with a pre-established value T .

The sequences with values equal to or greater than T (with values between 13 and 15) are further analyzed, while those with values less than T are discarded.

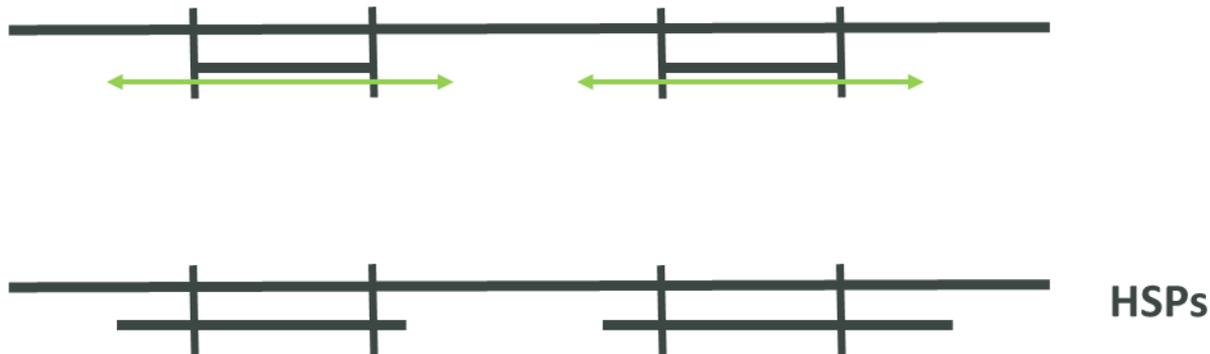
The substitution matrices assigned a positive score for each identity or for a substitution with amino acids of the same type (hydrophobic with hydrophobic, positively charged with positively charged etc ...) and negative for a substitution with amino acids of different type (eg. Basic amino acid with amino acid etc ..).

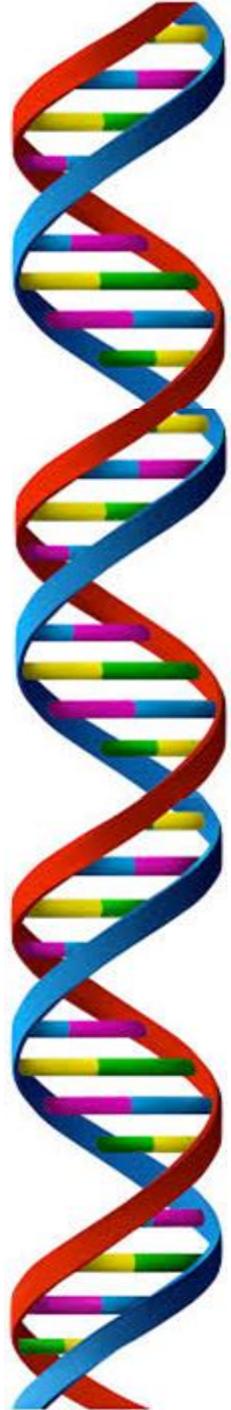
Allineamento singolo euristico: BLAST

Phase 4:

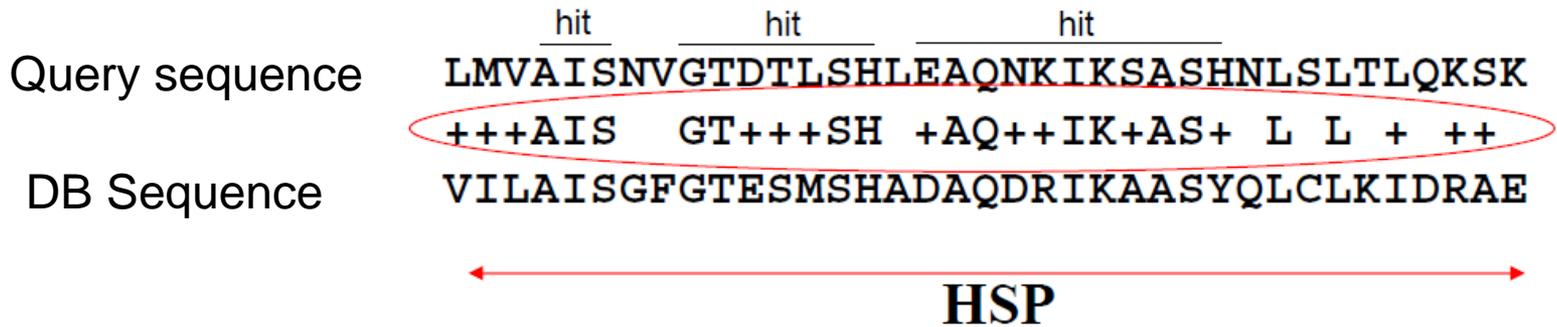
Extending the aligned zone. In this phase are taken into account, one at a time, the characters immediately adjacent (both right and left) to the **word**.

For each character added to the word you recalculate the score with the sequence in the database, and as long as the similarity score increases will be added to the adjacent characters the word (the pairing areas so found are called **HSP - high scoring segment pair**) .





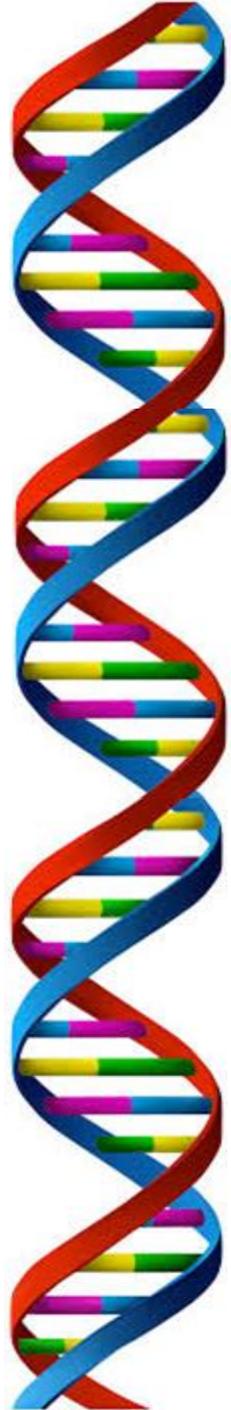
BLAST



Phase 5:

The HSP score (S) below a certain threshold (chosen empirically) are discarded. The output of BLAST phase 5 will contain the set of sequences whose HSP were not discarded with their scores.

A threshold X is used to determined the maximum tolerated score decrease.



BLAST

Phase 6:

BLAST generally operates on databases that contain millions of sequences

So not all alignments have a biological relevance, or not all of the output sequences are homologous to the query sequence / input.

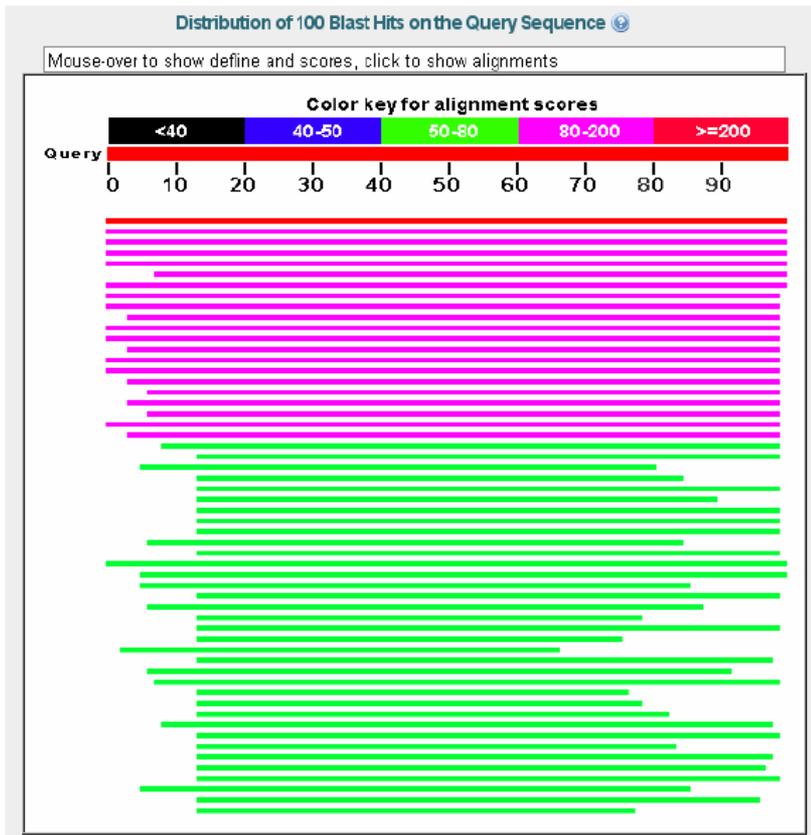
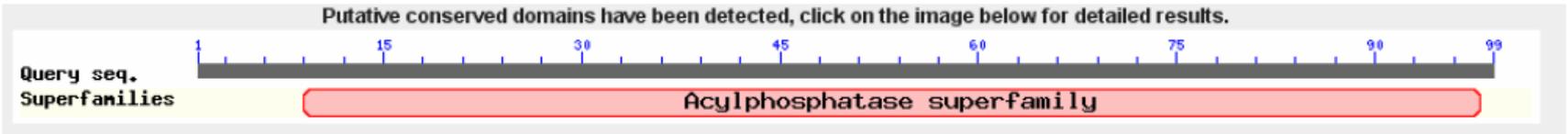
Therefore it is necessary a statistical evaluation of alignment obtained by BLAST

BLAST adds to each sequence present in the output the **E-Value**, a value that, properly interpreted by the researcher, will indicate how likely it is that the S score indicates a biological correlation between the two sequences.

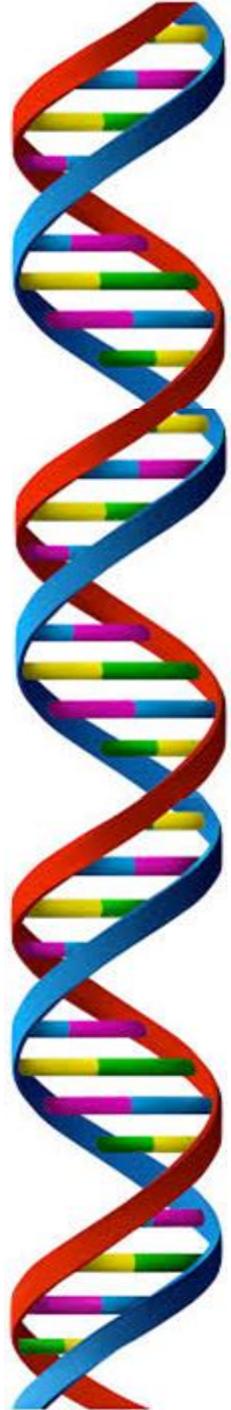
The more E-value is close to zero the more attest to a biological correlation between the sequences.

The result page of BLAST

If BLAST recognizes a wide similarity with one or more known protein families, it indicates it in a graphical way.



It is also present a summary graph that allows to quickly explore the first 100 hits, showing them aligned with the query sequence.



The result page of BLAST

BLAST includes in its output a textual summary of the sequences that were found: for each it shows a link (in this case to Entrez Protein), the name, the score (in bits) and The Expect value, according to which the results are sorted .

Sequences producing significant alignments:		Score (Bits)	E Value	
sp P07032.2 ACYP1 CHICK	RecName: Full=Acylphosphatase-1; AltN...	203	2e-52	G
sp P41500.2 ACYP1 BOVIN	RecName: Full=Acylphosphatase-1; AltN...	143	2e-34	G
sp P24540.2 ACYP1 PIG	RecName: Full=Acylphosphatase-1; AltNam...	141	9e-34	G
sp P56376.2 ACYP1 MOUSE	RecName: Full=Acylphosphatase-1; AltN...	139	5e-33	G
sp P07311.2 ACYP1 HUMAN	RecName: Full=Acylphosphatase-1; AltN...	135	7e-32	G
sp Q28FK7.1 ACYP1 XENTR	RecName: Full=Acylphosphatase-1; AltN...	134	2e-31	G
sp Q6DE05.1 ACYP1 XENLA	RecName: Full=Acylphosphatase-1; AltN...	132	3e-31	G
sp P14621.2 ACYP2 HUMAN	RecName: Full=Acylphosphatase-2; AltN...	130	1e-30	G
sp P35744.2 ACYP2 CAVPO	RecName: Full=Acylphosphatase-2; AltN...	130	2e-30	G
sp P56375.1 ACYP2 MOUSE	RecName: Full=Acylphosphatase-2; AltN...	130	2e-30	G
sp P00819.2 ACYP2 PIG	RecName: Full=Acylphosphatase-2; AltNam...	129	3e-30	G
sp P00820.2 ACYP2 RABIT	RecName: Full=Acylphosphatase-2; AltN...	129	4e-30	G
sp P35745.2 ACYP2 RAT	RecName: Full=Acylphosphatase-2; AltNam...	129	4e-30	G
sp P07033.2 ACYP2 BOVIN	RecName: Full=Acylphosphatase-2; AltN...	127	1e-29	G
sp P00818.2 ACYP2 HORSE	RecName: Full=Acylphosphatase-2; AltN...	126	3e-29	G
sp P07031.2 ACYP2 CHICK	RecName: Full=Acylphosphatase-2; AltN...	124	8e-29	G
sp P14620.2 ACYP2 ANAPL	RecName: Full=Acylphosphatase-2; AltN...	124	1e-28	G
sp P00821.2 ACYP2 MELGA	RecName: Full=Acylphosphatase-2; AltN...	122	4e-28	G
sp Q5EBE1.1 ACYP2 XENTR	RecName: Full=Acylphosphatase-2; AltN...	119	4e-27	G
sp Q9VF36.1 ACYP2 DROME	RecName: Full=Acylphosphatase-2; AltN...	97.4	1e-20	G
sp P56544.3 ACYP1 DROME	RecName: Full=Acylphosphatase-1; AltN...	83.6	2e-16	G
sp Q83AB0.2 ACYP COXBU	RecName: Full=Acylphosphatase; AltName...	69.7	3e-12	G
sp A9KH12.1 ACYP COXBN	RecName: Full=Acylphosphatase; AltName...	69.7	3e-12	G
sp A9FGA8.1 ACYP SORCS	RecName: Full=Acylphosphatase; AltName...	68.2	1e-11	G
sp A0LI66.1 ACYP SYNFM	RecName: Full=Acylphosphatase; AltName...	67.4	2e-11	G
sp P0AB66.1 ACYP EC057	RecName: Full=Acylphosphatase; AltName...	65.9	6e-11	G
sp A5VFP2.1 ACYP SPHWV	RecName: Full=Acylphosphatase; AltName...	65.5	6e-11	G
sp Q83LL9.4 ACYP SHIFL	RecName: Full=Acylphosphatase; AltName...	65.1	8e-11	G

The results page of BLAST

Under the descriptions begin the actual results, i.e., HSP: in addition to the description, the result shows the identity, similarity (Positives) and the inserted gaps, if any.

```
>\[sp|P14620.2|ACYP2 ANAPL RecName: Full=Acylphosphatase-2; AltName: Full=Acylphosphate  
phosphohydrolase 2; AltName: Full=Acylphosphatase, muscle type  
isozyme  
Length=103
```

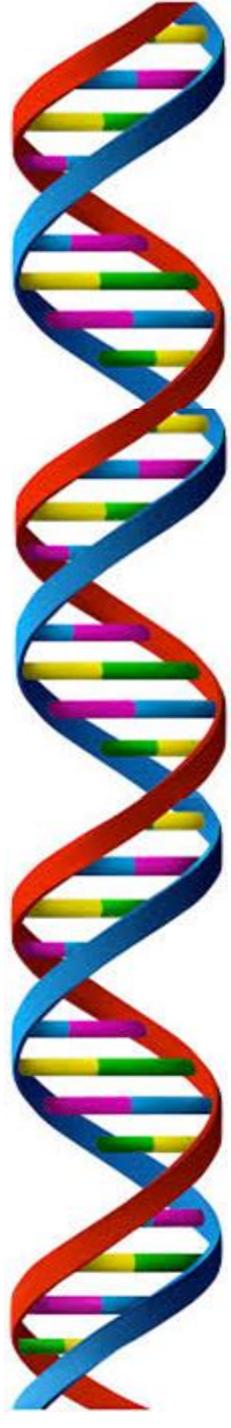
```
Score = 124 bits (311), Expect = 1e-28, Method: Compositional matrix adjust.  
Identities = 56/92 (60%), Positives = 69/92 (75%), Gaps = 0/92 (0%)
```

```
Query 7 LMSVDYEVSGRVQGVFFFRKYTQSEAKRLGLVGWVRNTSHGTVQGQAQGPAAARVRELQEWL 66  
L SVDYEV GRVQGV FR YT+ EA++LG+VGWV+NTS GTV GQ QGP +V ++ WL  
Sbjct 11 LKSVDYEVFGRVQGVCFRMYTEEEARKLGVVGVWVKNNTSQGTVTGQVQGPEDKVNAMKSWL 70  
  
Query 67 RKIGSPQSRISRAEFTNEKEIAALEHTDFQIR 98  
K+GSP SRI R F+NEKEI+ L+ + F R  
Sbjct 71 TKVGSPPSRIDRTNFSNEKEISKLDFSGFSTR 102
```

```
>\[sp|P00821.2|ACYP2 MELGA RecName: Full=Acylphosphatase-2; AltName: Full=Acylphosphate  
phosphohydrolase 2; AltName: Full=Acylphosphatase, muscle type  
isozyme; AltName: Full=Acylphosphatase isozyme TU1  
Length=103
```

```
Score = 122 bits (306), Expect = 4e-28, Method: Compositional matrix adjust.  
Identities = 56/95 (58%), Positives = 69/95 (72%), Gaps = 0/95 (0%)
```

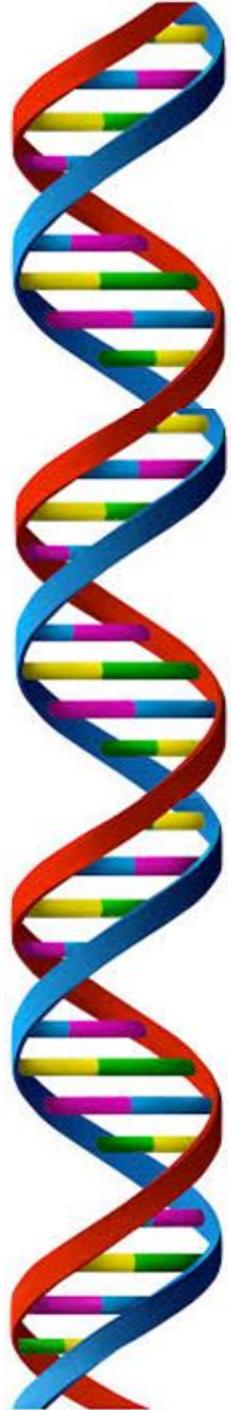
```
Query 4 SEG LMSVDYEVSGRVQGVFFFRKYTQSEAKRLGLVGWVRNTSHGTVQGQAQGPAAARVRELQ 63  
S L SVDYEV GRVQGV FR YT+ EA++LG+VGWV+NT GTV GQ QGP +V ++  
Sbjct 8 SGALKSVDYEVFGRVQGVCFRMYTEEEARKLGVVGVWVKNTRQGTVTGQVQGPEDKVNAMK 67  
  
Query 64 EWL RKIGSPQSRISRAEFTNEKEIAALEHTDFQIR 98  
WL K+GSP SRI R F+NEKEI+ L+ + F R  
Sbjct 68 SWLSKVGSPSRIDRTNFSNEKEISKLDFSGFSTR 102
```



BLAST releases

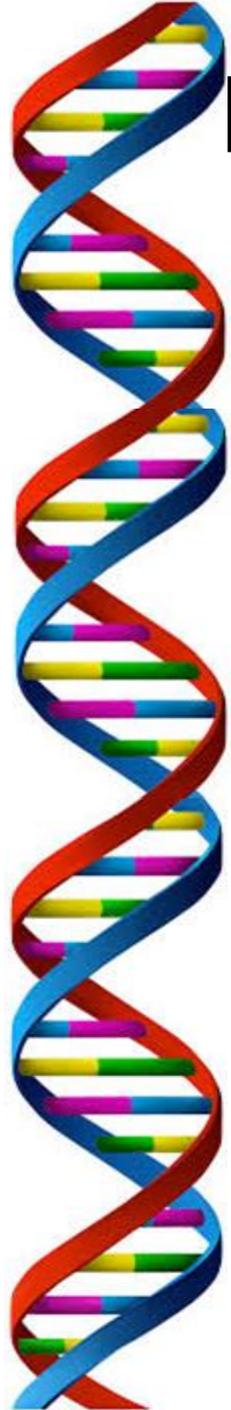
BLAST has recently been implemented with a two-hit method, which provides for the extension of the HSP regions can only happen if two independent hits occur within a number of residues (represented by threshold A), with no gaps in between. There were also implemented several specialized versions:

- **blastp**: look for similarities in protein databases from a query of amino acids.
- **blastn**: look for similarities in nucleotide databases from a query of nucleotides.
- **blastx**: look for similarities in protein databases from a query of nucleotides that is translated into all possible frame.
- **tblastn**: look for similarities in nucleotide databases from a query of amino acids, translating the entries of the database into nucleotides.
- **tblastx**: look for similarities in nucleotide databases from a query of nucleotides, translating all the entries of the database into amino acids



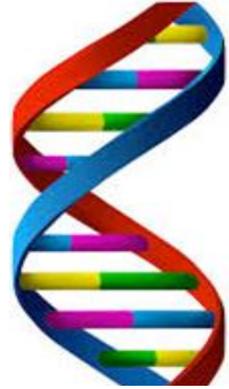
BLAST

- BLAST has 4 fundamental parameters:
 - **W**: word size, the greater the number, the smaller the number of words generated, the lower the execution time. But the sensitivity decreases considerably.
 - **T**: threshold, the lower the number, the greater the number of w-mers included in the list, the higher the running time. But it has an enhanced sensitivity.
 - **S**: score, the smaller the number, the greater the HSP length.
 - **X**: the lower the number, the greater the number of w-mers included in the list, the higher the running time. But it has an enhanced sensitivity.



Pairwise Alignment Considerations: FASTA vs BLAST

- FASTA differs from BLAST since it searches the best alignment between the entire sequence under investigation and the sequences of the database used as a reference. Instead BLAST search only local sequence homology.
- FASTA also uses a scoring matrix only during the confrontation extension while BLAST uses a scoring matrix during several phases of the alignment (that is, scanning and extension).
- FASTA examines the amino acids in pairs ($ktup = 2$) or individually taken ($ktup = 1$), BLAST uses for comparison of 3-4 amino groups (words).
- Both the algorithm that governs the early stages of FASTA and the one of BLAST do not permit the presence of gaps within each sequence segment taken into consideration.
- Unlike BLAST, however, the FASTA algorithm contemplates the possibility of insertions and deletions in the last phase alignment.



Multiple alignment vs Pairwise alignment

A

```
1: EAGFPPGVVNI PGFGPTAGAAIASHEDVDKVAFTGSTE VGH LIQVA
2: EAGFPPGVVNI VPGFGPTAGAAIASHEDVDKVAFTGSTE IGRV IQVA
3: QYMDQONLYLVVKGG-VPETTELL--KERFDHIMYTGSTAVGKIVMAA
4: NVFSPAWA-TVVEGDETI SQQLL--QEKFDHIFFTGSPRVGRLIMAA
5: EAGVPVGLVNVVQG-GAETGSLLCHHPNVAKVSFTGSVPTGKKVMEM
6: DI-FPAGVINILFGRGKTVDPLTGHPKVRMVSLTGS IATGEHI I SH
```

B

```
1: EAGFPPGVVNI VPGFGPTAGAAIASHEDVDKVAFTGSTE VGH LIQVA
2: EAGFPPGVVNI VPGFGPTAGAAIASHEDVDKVAFTGSTE IGRV IQVA
3: QYMDQONLYLVVKGG-VPETTELL--KERFDHIMYTGSTAVGKIVMAA
4: NVFSPAWA-TVVEGDETI SQQLL--QEKFDHIFFTGSPRVGRLIMAA
5: EAGVPVGLVNVVQG-GAETGSLLCHHPNVAKVSFTGSVPTGKKVMEM
6: DI-FPAGVINILFGRGKTVDPLTGHPKVRMVSLTGS IATGEHI I SH
```



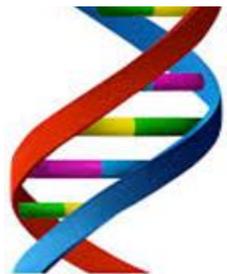


Multiple Alignment

```
TRYI_DROME : IIGGSDQLIRNAPWQVS IQISAR---HECGGVIYSKEI IITAGHCLHER-SVILMKV-----RVGA---QNHNYGG-TLVPVAAY--KVHEQFDSRFLH--- : 84
ENIK_PIG/8 : IVGGNDSREGAWPFWVALYNG---QLLOGASLVS RDWLVSAHCVYG---FNLEPSKWKAILG--LHMTSNLITS PQIVIRLIDE IVINPHYNRRRKD--- : 90
THR_BOVIN : IVEGQDAEVGLSPWQVMLFRKSPQE--LLCGASLISDRWVLTAAHCLLYPPWDKNFTVDDLLVRIK-HSRTRYERKVEKISMLDK-IYIHPRYNWKEN--- : 95
KLK1_MOUSE : IVGGFNCEKNSQPWQVAVYRFT---KYQCGGILLINVNWVLTAAHCHND---KYQV---WL GK-MNFLEDEPSAQHRLMSK-AI PHPDFNMSLLNEHT : 86
CIRA_BOVIN : IVNGEEAVPGSWPQVSLQDKTG--FHF CGGSLINENWVTAHCGVT---TSDWV---VAGEFDQGSSEK-IQK LKI AK-VFKNSKYNSLITIN--- : 85
CTRL_ANOGA : WGGEVAKNGSAPYQVSLQVPGWG--HNCGGSLINDRWVLTAAHCLVG-HAPGDLMV---LVGF---NSLKEGG-ELLKVDK-LLYHSRYNLPFRFH--- : 85
CTRL_HALRU : IVGGSNAAAGEFPWQGS LQVRS GTSWFHICGCVLYTTSKALTAHCLSN-SASSYRL-G--FGMLR-MNNVDGTEQYSSVTS-YTNHPNYNGNAAG--- : 90
```

```
TRYI_DROME : -----YDI AVLRLSTP-LTFGLSTRAINLAS---TSP--SGGITVIVTGWGH---TDNG---ALSDSLQKAQLQIIDRGECASQKFGYGAD-FVGEETI : 165
ENIK_PIG/8 : -----SDI AMMHLEFK-VNYTDYIQPICLPE---ENQVFPGRICS IAGWGK---VIYQG---SPADILQEADVPLLSNEKQQQMP-EYN---ITENMM : 171
THR_BOVIN : -----LDRDI ALLKLRP-IELSDYIHPVCLPDKQTA AKL LHAGFKGRVTGWGNRREIWTTSVAEVQPSVLQVWNLPLVERPVCKAS---TRIR---ITDNMF : 186
KLK1_MOUSE : PQEEDYSNDLM LRLKRP-ADITDVKPIDLPT---EEP--KL GSTCLASGWS---ITPVKY--EYDELQCVNLKLLPNEDCAKA---HIEK---VTDIML : 173
CIRA_BOVIN : -----NDI TLLKLS TA-ASF SQT VSAVCLPS---ASDDFAAGITCVITGWGL---TRYTNA--NTPDRLQQASLPLLSNINCKKY---WGTK---IKDAMI : 166
CTRL_ANOGA : -----NDI GLVRLEQP-VQFSELVQSVEYSE---KAVPANATVRLTGWGR---TSANG---PSP TLLQSLNVVILSNEDCNKK---GGDPGYTDVGH L : 165
CTRL_HALRU : -----YPND IAVLRLTSSMDTSSSAVGPSVWLL-----VERLCRTNMYDQR--MGKTQWRWQHPPNNLQKVDMTVLINSDCSSFWSGISGAT-VNSGHI : 175
```

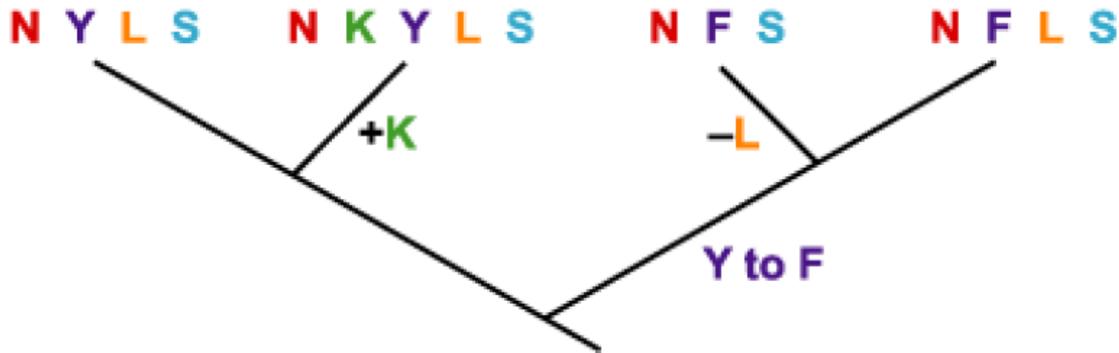
```
TRYI_DROME : CAAS---TD-ADACTGDSGGPLVASSQ-----LVGIVSWG-YRCADDNYPGVYADVAILRPWI : 218
ENIK_PIG/8 : CAGYE--EGG-IDSCQDSGGPLMCLEN--NRWLLAGVTSFG-YQCALPNRPGVYARVPKFTEWI : 230
THR_BOVIN : CAGYKPGEGKRGDACEDSGGPFVMSKSPYNNRWFYQMGIVSWG-EGCDRDGKYGFYTHVRLKKWI : 250
KLK1_MOUSE : CAGDM--DGG-KDTCAGDSGGPLICDGV-----LQGITSWGSPSGKPNVPGIYTRVLNFNWI : 229
CIRA_BOVIN : CAGA---SG-VSSCMDSGGPLVCKKN--GAWTILVGVSWG-SSTCSTSTPGVYARVIALVNW : 223
CTRL_ANOGA : CTLTK---TG-EGACDSGGPLVYEK-----LVGVVNFVGV-VECALG-YPDGFARVSYHDWV : 218
CTRL_HALRU : CIFE---SG-RSACDSGGPLVCGNT-----LITGITSWGISSCSGS-YPSVYTRVSSFYNWV : 228
```



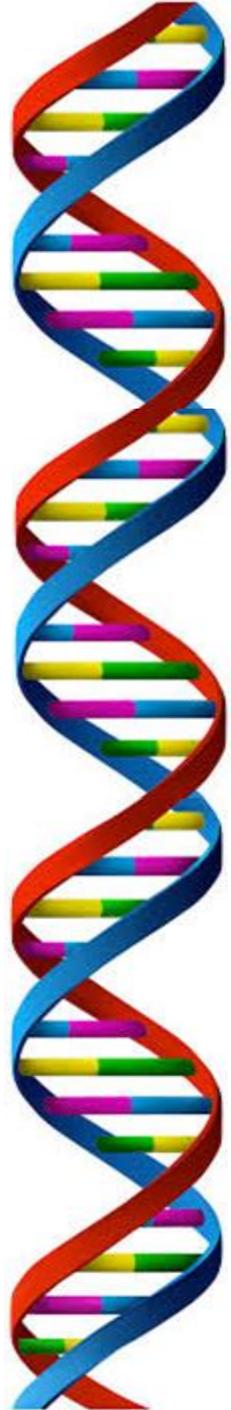
Conservation of the catalytic amino acids in some members of the trypsin family

Multiple Alignment

seqA	N	•	F	L	S
seqB	N	•	F	-	S
seqC	N	K	Y	L	S
seqD	N	•	Y	L	S



The evolutionary history is deduced from the alignment



Multiple Alignment

- Difficulty
 - objective scoring Function: weight to be assigned to the various sequences
 - Complexity of the problem: computation time
- Clustal, ClustalW and ClustalX