# 生物信息学：导论与方法
# Bioinformatics: Introduction and Methods



https://www.coursera.org/course/pkubioinfo

# 生物信息学：导论与方法
# Bioinformatics: Introduction and Methods

## 北京大学生物信息学中心 高歌、魏丽萍
## Ge Gao & Liping Wei
## Center for Bioinformatics, Peking University

# Sequence Database Search

## 北京大学生物信息学中心 高歌

### Ge Gao, Ph.D.
### Center for Bioinformatics, Peking University

# Unit 1:
# Sequence Databases

## 北京大学生物信息学中心 高歌

**Ge Gao, Ph.D.**

**Center for Bioinformatics, Peking University**

**New Best Alignment = Previous Best + Local Best**

Score of Best Previous Alignment

( *Russ Altman BMI214*)

$$
\begin{array}{ccc}
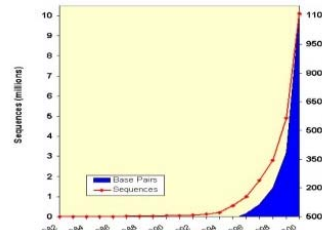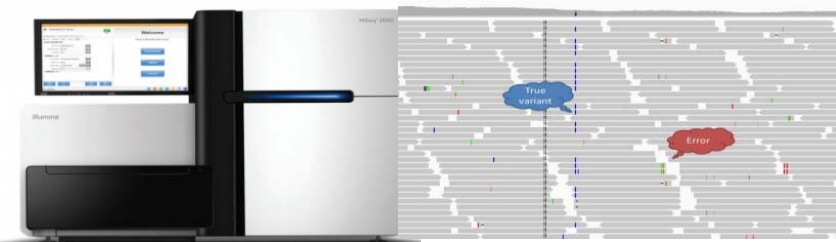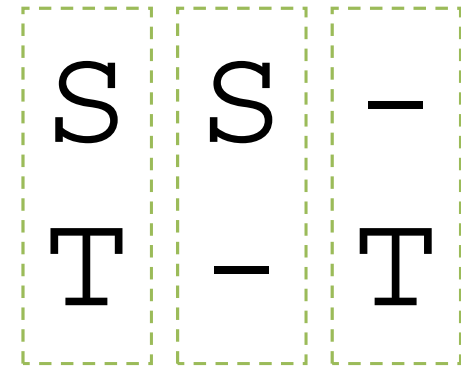\text{S} & \text{S} & - \\
\text{T} & - & \text{T}
\end{array}
$$

**Global alignment (Needleman-Wunsch)**

$$F(0,0) = 0$$

$$F(i,j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \end{cases}$$

**Local alignment (Smith-Waterman)**

$$F(0,0) = 0$$

$$F(i,j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \\ 0 \end{cases}$$

# Sequence Alignment

## "How can we determine the similarity between two sequences?"

Why is it important?

- Similar sequence ➔ Similar structure ➔ Similar function (The "*Sequence-to-Structure-to-Function Paradigm*")
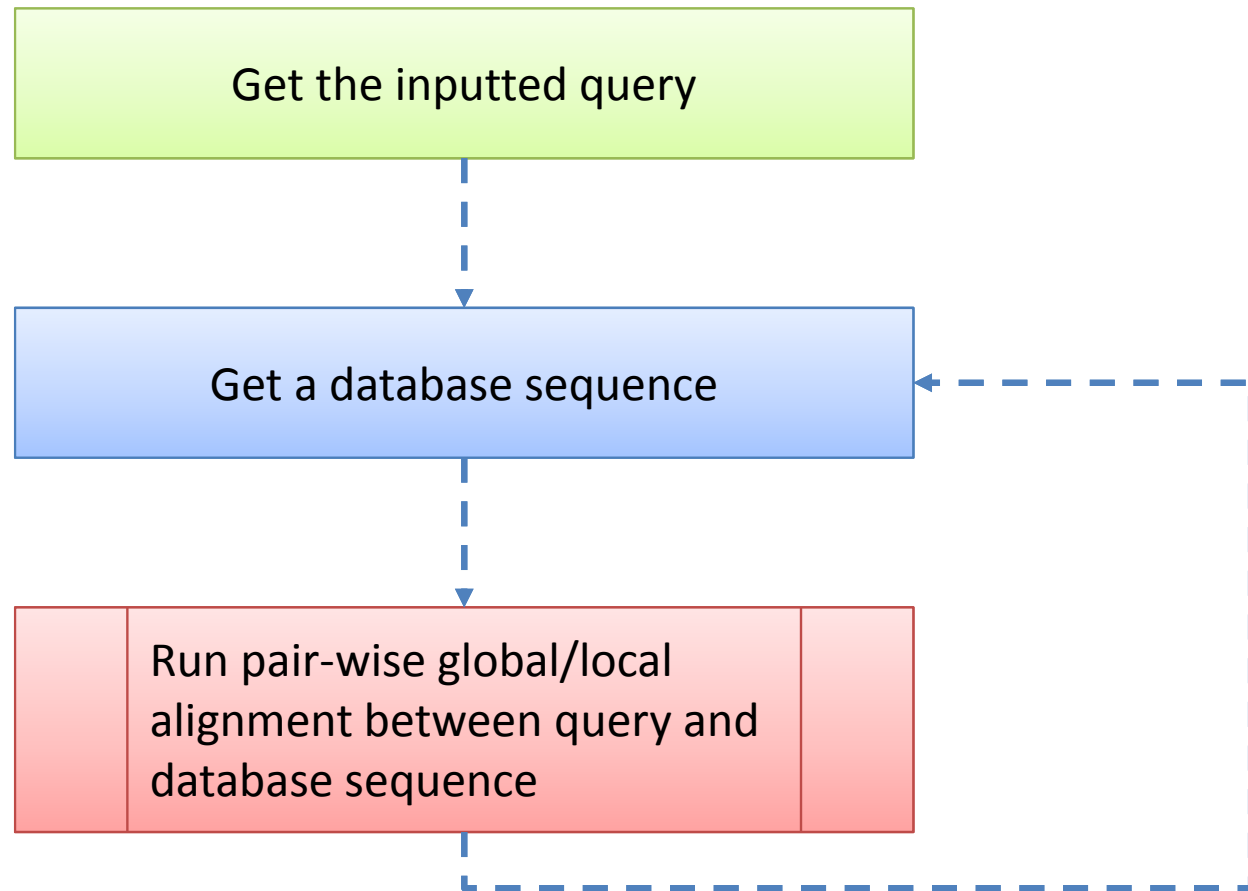- Similar sequence ➔ Common ancestor ("*Homology*")

# Sequence Database Searching

- Rather than do the alignment pair-wise, it's more often to search <span style="color:blue">sequence database</span> in a <span style="color:red">high-throughput</span> style.

- Or, identify similarities between
  - <span style="color:red">novel query sequence</span>
    whose structures and functions are usually unknown and/or uncharacterized
  - <span style="color:red">sequences in (public) databases</span>
    whose structures and functions have been elucidated and annotated.
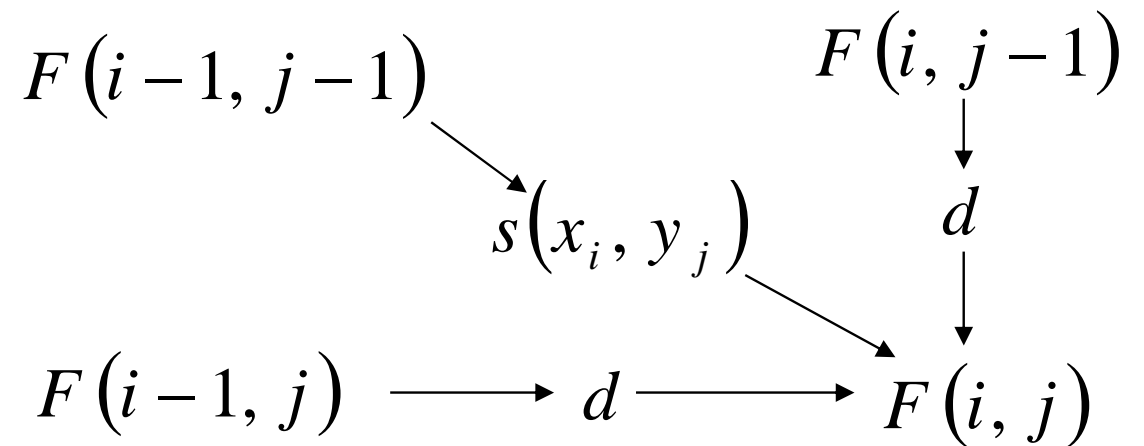
# Sequence Database Searching
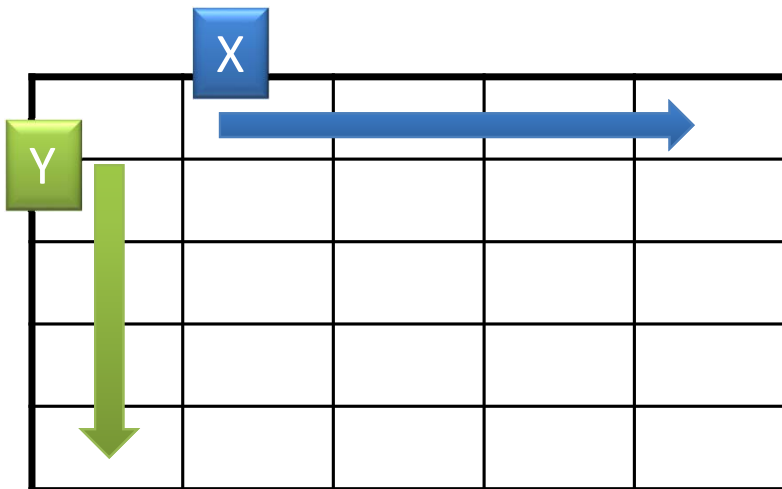
- The query sequence is compared/aligned with every sequence in the database

- Statistically significant hits are assumed to be related to the query sequence
  - Similar function/structure
  - Common evolutionary ancestor

# A (naïve) algorithm for database searching



Get the inputted query

Get a database sequence

Run pair-wise global/local alignment between query and database sequence

$$F(0,0) = 0$$

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) & \text{$x_i$ aligned to $y_j$} \\ F(i-1, j) + d & \text{$x_i$ aligned to \textit{a gap}} \\ F(i, j-1) + d & \text{$y_j$ aligned to \textit{a gap}} \end{cases}$$

There are nm entries in the matrix.
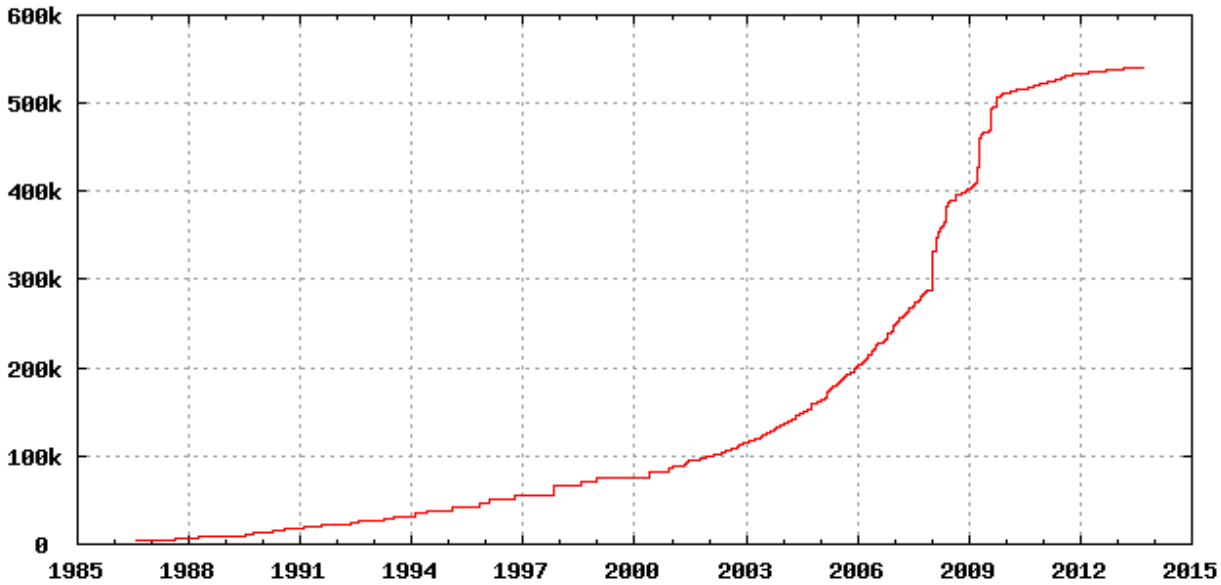
Sequence X of length m

Sequence Y of length n

Each entry requires a constant number c of operation(s).

Dynamic programming matrix

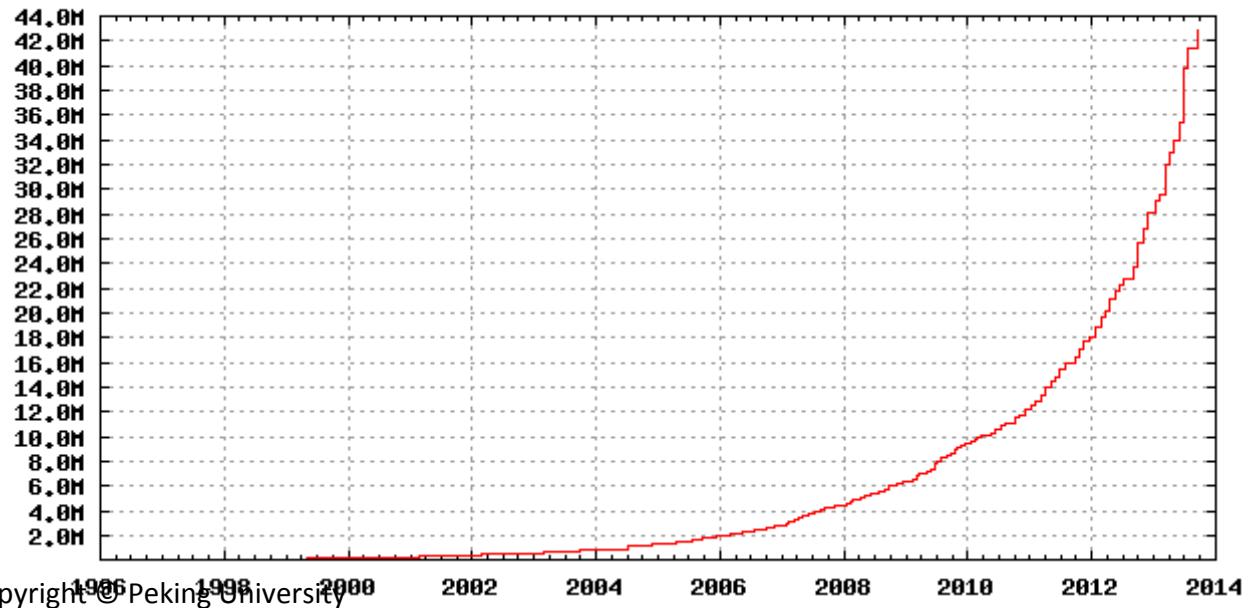c*m*n operations needed in total, for one pair-wise alignment.

- Say your query sequence (HBA_HUMAN) has 142 amino acids

- Most recent release of human-curated Swiss-Prot protein databases contains 540,958 sequences with 192,206,270 amino acids (Sept 18[th], 2013);
  - On average, the sequence length is 192,206,270/540,958 = 355.30 aa

- And assume your super-fast computer can run one operation in 1μs = (0.000001s)

- Then, you will need 7.8 hr for ONE comparison!

Number of entries in UniProtKB/Swiss-Prot

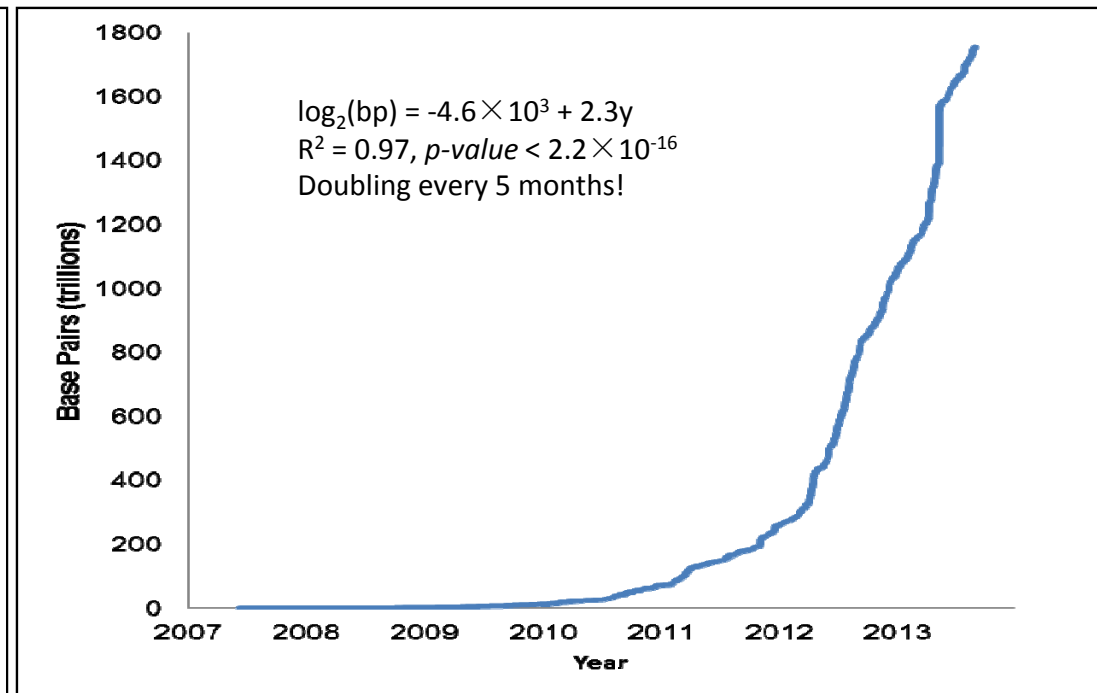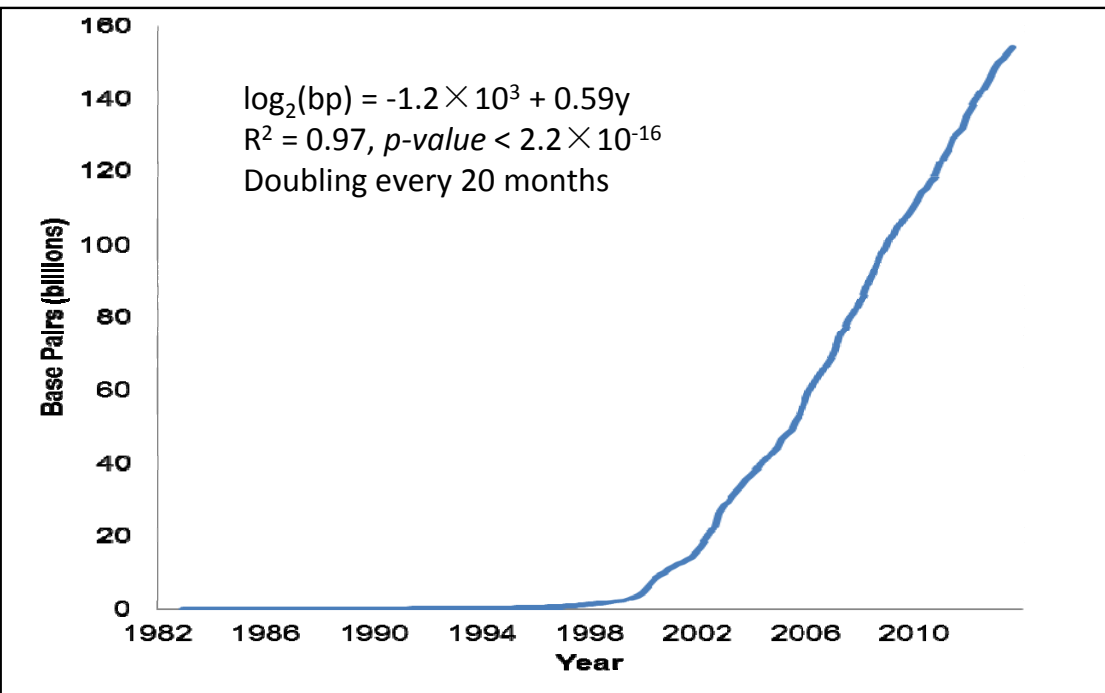Source: http://web.expasy.org/docs/relnotes/relstat.html

Source: http://www.ebi.ac.uk/uniprot/TrEMBLstats

Number of entries in UniProtKB/TrEMBL

# Genbank

# SRA

$\log_2(bp) = -1.2 \times 10^3 + 0.59y$
$R^2 = 0.97$, *p-value* $< 2.2 \times 10^{-16}$
Doubling every 20 months

$\log_2(bp) = -4.6 \times 10^3 + 2.3y$
$R^2 = 0.97$, *p-value* $< 2.2 \times 10^{-16}$
Doubling every 5 months!
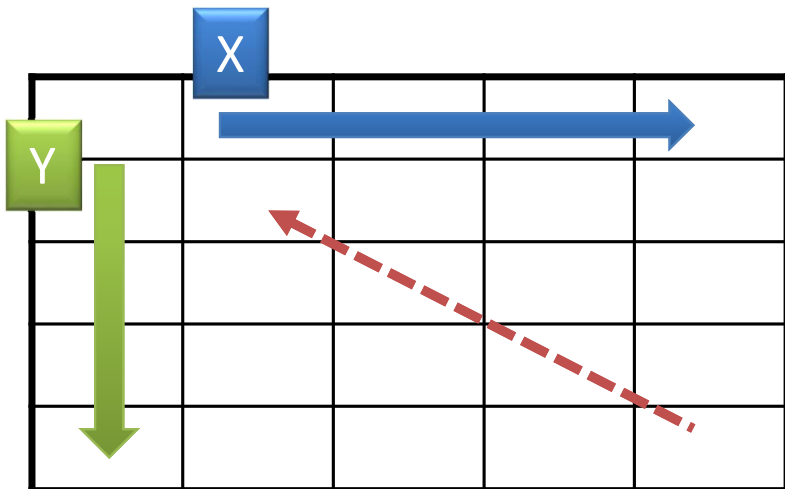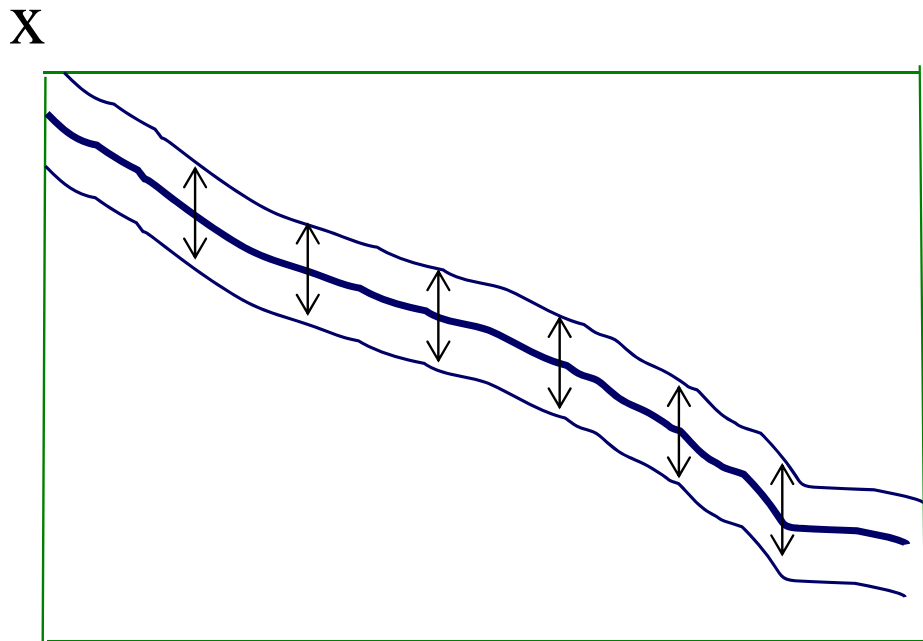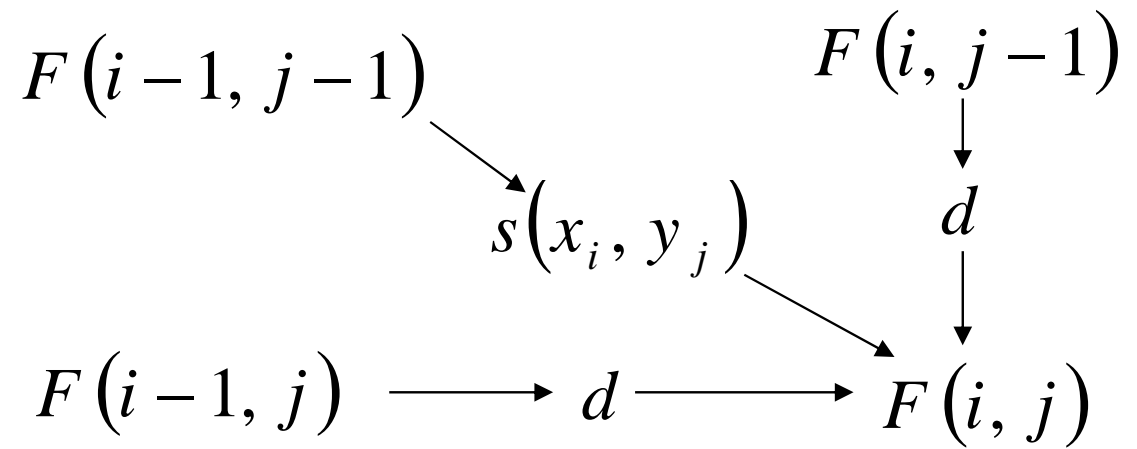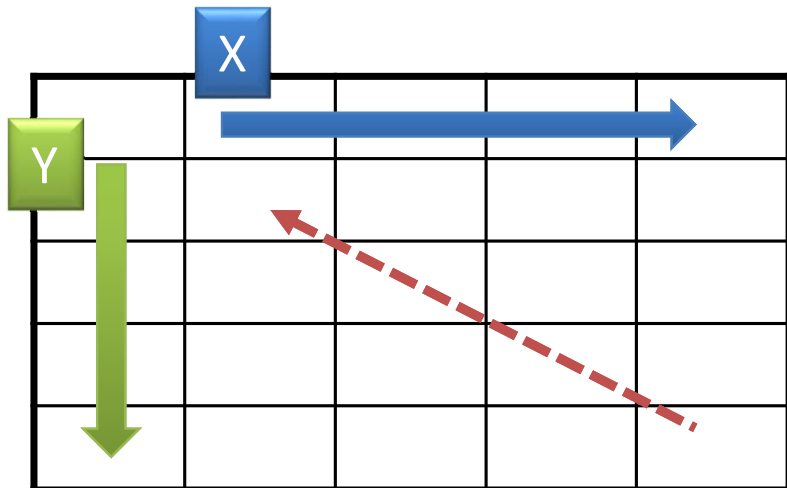
$$F(0,0) = 0$$

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) & x_i \text{ aligned to } y_j \\ F(i-1, j) + d & x_i \text{ aligned to } a \text{ gap} \\ F(i, j-1) + d & y_j \text{ aligned to } a \text{ gap} \end{cases}$$
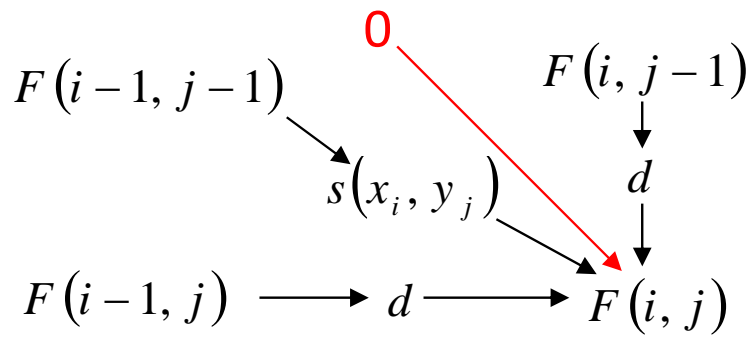


$F(i-1, j-1)$     $F(i, j-1)$

$s(x_i, y_j)$    $d$

$F(i-1, j) \longrightarrow d \longrightarrow F(i, j)$

$$F(i-1, j-1)$$

$$F(i, j-1)$$

$$s(x_i, y_j)$$

$$d$$

$$F(i-1, j) \longrightarrow d \longrightarrow F(i, j)$$

X

Y

x

y

| | | |
|---|---|---|
| HBA_HUMAN | 1 | MV-LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-D | 48 |
| | | \|\| \|:\|.:\|:.\|.\|.\|\|\|\| :.\|.\|.\|\|\|.\|:.:.:\|.\|:.:\|..\| \| |
| HBB_HUMAN | 1 | MVHLTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGD | 48 |
| HBA_HUMAN | 49 | LS-----HGSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLR | 93 |
| | | \|\| .\|:.:\|\|.\|\|\|\|\|..\|.:.:.:\|\|:\|\|...:\|\|:\|\|..\|\|. |
| HBB_HUMAN | 49 | LSTPDAVMGNPKVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLH | 98 |
| HBA_HUMAN | 94 | VDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR | 142 |
| | | \|\|\|.\|\|:\|\|.\|.:.\|:..\|\|.\|...\|\|\|\|.\|.\|:..\|:\|.\|:..\|..\|\|. |
| HBB_HUMAN | 99 | VDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH | 147 |

A    G

A    G

|     |     | A   | A   | G   |
|-----|-----|-----|-----|-----|
|     | 0   | 0   | 0   | 0   |
| A   | 0   | 2   | 2   | 0   |
| G   | 0   | 0   | 0   | 4   |
| C   | 0   | 0   | 0   | 0   |

$F(i-1, j-1)$

$s(x_i, y_j)$

$F(i, j-1)$

$d$

0

$F(i-1, j)$ ⟶ $d$ ⟶ $F(i, j)$

# BLAST: Intro

- To make the alignment effectively, a *Heuristic* algorithm BLAST (Basic Local Alignment Search Tool) is proposed by Altschul *et al* in 1990.

- BLAST finds the highest scoring locally optimal alignments between a query sequence and a database.
  - Very fast algorithm
  - Can be used to search extremely large databases
  - Sufficiently sensitive and selective for most purposes
  - Robust – the default parameters just work for most cases

# Basic Local Alignment Search Tool

**Stephen F. Altschul[1], Warren Gish[1], Webb Miller[2]
Eugene W. Myers[3] and David J. Lipman[1]**

[1]*National Center for Biotechnology Information
National Library of Medicine, National Institutes of Health
Bethesda, MD 20894, U.S.A.*

[2]*Department of Computer Science
The Pennsylvania State University, University Park, PA 16802, U.S.A.*

[3]*Department of Computer Science
University of Arizona, Tucson, AZ 85721, U.S.A.*

A new approach to rapid sequence comparison, basic local alignment search tool (BLAST),
directly approximates alignments that optimize a measure of local similarity, the maximal
segment pair (MSP) score. Recent mathematical results on the stochastic properties of MSP
scores allow an analysis of the performance of this method as well as the statistical
significance of alignments it generates. The basic algorithm is simple and robust; it can be
implemented in a number of ways and applied in a variety of contexts including straight-
forward DNA and protein sequence database searches, motif searches, gene identification
searches, and in the analysis of multiple regions of similarity in long DNA sequences. In
addition to its flexibility and tractability to mathematical analysis, BLAST is an order of
magnitude faster than existing sequence comparison tools of comparable sensitivity.

blast.ncbi.nlm.nih.gov/Blast.cgi

BLAST®　　　　　　　　　Basic Local Alignment Search Tool

| Home | Recent Results | Saved Strategies | 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.

› NCBI/ BLAST Home

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

**New** DELTA-BLAST, a more sensitive protein-protein search    Go

## BLAST Assembled RefSeq Genomes

Choose a species genome to search, or list all genomic BLAST databases.

- Human
- Mouse
- Rat
- Arabidopsis thaliana
- Oryza sativa
- Bos taurus
- Danio rerio
- Drosophila melanogaster
- Gallus gallus
- Pan troglodytes
- Microbes
- Apis mellifera

## Basic BLAST

Choose a BLAST program to run.

nucleotide blast — Search a **nucleotide** database using a **nucleotide** query
_Algorithms:_ blastn, megablast, discontiguous megablast

protein blast — Search **protein** database using a **protein** query
_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
- Search trace archives
- 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

### News

Update to SRA-BLAST

SRA-BLAST has undergone a dramatic update, both in terms of user interface and search performance.
Thu, 20 Jun 2013 11:00:00 EST

More BLAST news...

### Tip of the Day

Use Genomic BLAST to see the genomic context

If you are interested in the evolution of a particular gene or gene family it is often intetesting to examine the intro-exon structure even across species.

More tips...

BLAST is a registered trademark of the National Library of Medicine.

Copyright | Disclaimer | Privacy | Accessibility | Contact | Send feedback        NCBI | NLM | NIH | DHHS

**WELCOME**

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.
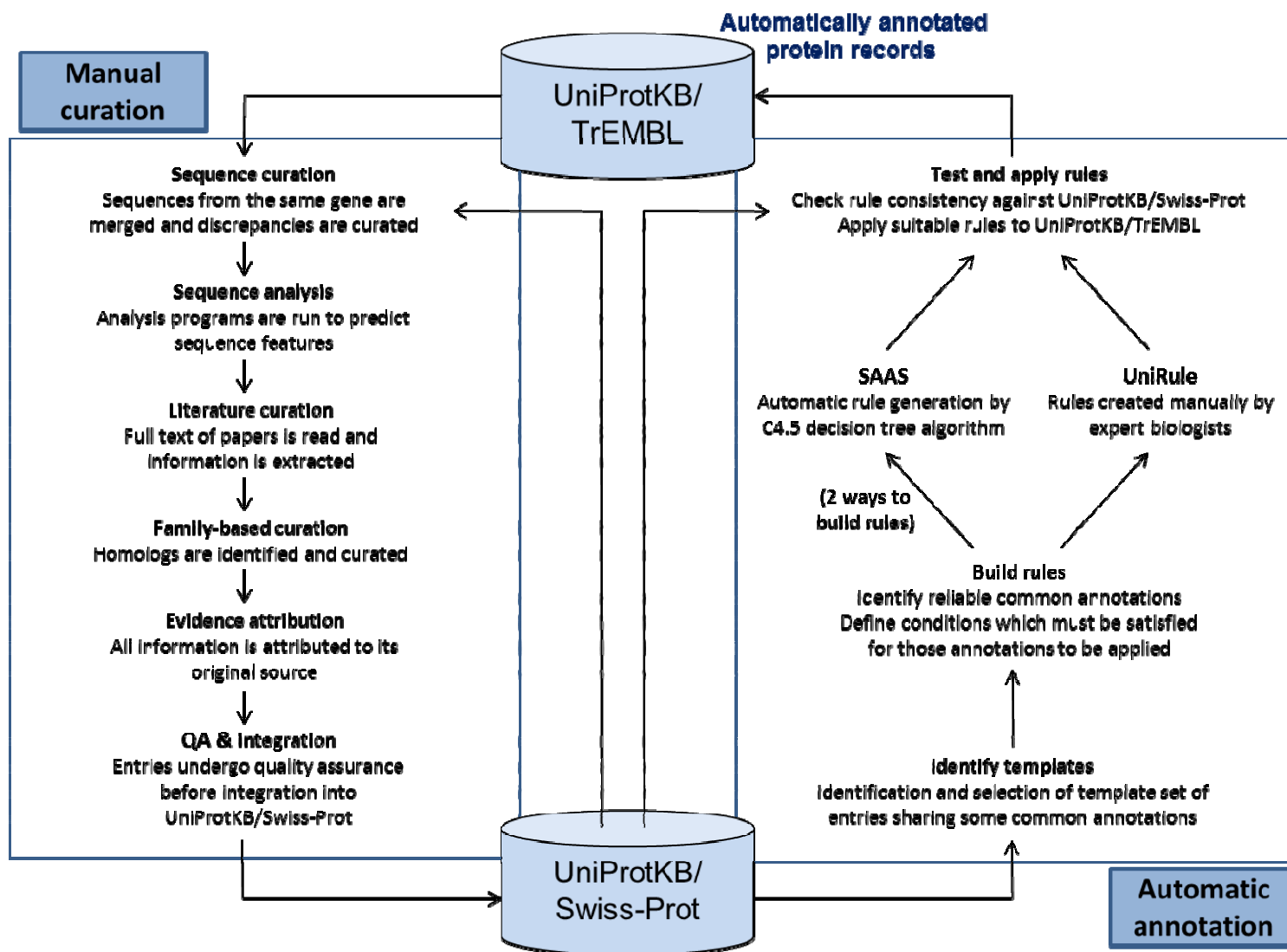
**data**
**protein sequence**

**knowledge**
**functional information**

**Manual curation**

**UniProtKB/ TrEMBL**

**Automatically annotated protein records**

UniProt

**Sequence curation**
Sequences from the same gene are merged and discrepancies are curated

**Sequence analysis**
Analysis programs are run to predict sequence features

**Literature curation**
Full text of papers is read and information is extracted

**Family-based curation**
Homologs are identified and curated

**Evidence attribution**
All information is attributed to its original source

**QA & integration**
Entries undergo quality assurance before integration into UniProtKB/Swiss-Prot

**Test and apply rules**
Check rule consistency against UniProtKB/Swiss-Prot
Apply suitable rules to UniProtKB/TrEMBL

**SAAS**
Automatic rule generation by C4.5 decision tree algorithm

**UniRule**
Rules created manually by expert biologists

(2 ways to build rules)

**Build rules**
Identify reliable common annotations
Define conditions which must be satisfied for those annotations to be applied

**Identify templates**
Identification and selection of template set of entries sharing some common annotations

**UniProtKB/ Swiss-Prot**

**Automatic annotation**

Swiss Institute of Bioinformatics

**(Modified from http://education.expasy.org/cours/Turin/UniProtKB_Turin.ppt)**

**Names and origin**

| Protein names | Recommended name:<br>**Elongation factor Tu 1**<br>Short name=EF-Tu 1<br>Alternative name(s): |
|---|---|
| Gene names | |
| Organism | |
| Taxonomic identi... | |
| Taxonomic lineage | Bacteria > Proteobacteria > Gammaproteobacteria > Enterobacteriales ><br>Enterobacteriaceae > Escherichia |

**Protein and gene names**
**Taxonomic information**

**General annotation (Comments)**

| Function | This protein promotes the GTP-dependent binding of aminoacyl-tRNA to |
|---|---|
| Subunit structure | |
| Subcellular locat... | |
| Miscellaneous | |
| | The antibiotic pulvomycin inhibits protein biosynthesis by disrupting the allosteric control mechanism of EF-Tu. HAMAP MF_00118_B |
| Sequence similarities | Belongs to the GTP-binding elongation factor family. EF-Tu/EF-1A subfamily. |

**Manual annotation**
**Function, Subcellular location,**
**Catalytic activity, Disease,**
**Tissue specificty, Pathway...**

**References**

« Hide 'large scale' references

[1] "The nucleotide sequence of the cloned tufA gene of Escherichia coli."
Yokota T., Sugisaki H., Takanami M., Kaziro Y.
Gene 12:25-31(1980) [PubMed: 7011903] [Abstract]
Cited for: NUCLEOTIDE SEQUENCE [GENOMIC DNA].

[2] "The complete gen..."
Blattner F.R., Plunke..., Riley M., Collado-Vides J., Glasner J.D., Rode C.K., Ma..., ...ick H.A., Goeden M.A., Rose D.J., Mau B., Shao Y.
Science 277:1453-1...
Cited for: NUCLEOTIDE SEQUENCE [LARGE SCALE GENOMIC DNA].
Strain: K12 / MG1655 / ATCC 47076.

[3] "Highly accurate genome sequences of Escherichia coli K-12 strains MG1655 and W3110."
Hayashi K., Morooka N., Yamamoto Y., Fujita K., Isono K., Choi S., Ohtsubo E., Baba T., Wanner B.L., Mori H., Horiuchi T.
Mol. Syst. Biol. 2:E1-E5(2006) [PubMed: 16738553] [Abstract]
Cited for: NUCLEOTIDE SEQUENCE [LARGE SCALE GENOMIC DNA].
Strain: K12 / W3110 / ATCC 27325 / DSM 5911.

**References**

**Cross-references**

**Sequence databases**

| EMBL<br>GenBank<br>DDBJ | J01690 Genomic DNA. Translation: AAA50993.1.<br>M10459 Genomic DNA. Translation: AAA24702.1. Sequence problems.<br>U00096 Genomic DNA. Translation: AAC76364.1.<br>AP009048 Genomic DNA. Translation: BAE77952.1.<br>U18997 Genomic DNA. Translation: AAA58136.1.<br>AF058450 Genomic DNA. Translation: AAC14286.1. |
|---|---|
| PIR | EFECTA. A91475. |
| RefSeq | AP_004451.1. AC_000091.1.<br>NP_417798.1. NC_000913.2. |

**3D s...**

**Cross-references**
**to over 125 databases**

| | | | | | PDBsum |
|---|---|---|---|---|---|
| | | | | 34 | [»] |
| | | | | 34 | [»] |
| | | | | 34 | [»] |
| 2HCJ | X-ray | 2.12 | A<br>B | 9-45<br>60-394 | [»] |
| 2HDN | X-ray | 2.80 | A/C/E/G/I/K<br>B/D/F/H/J/L | 9-45<br>60-394 | [»] |

| ProteinModelPortal | P0CE47. |
|---|---|
| SMR | P0CE47. Positions 9-393. |
| ModBase | Search... |

**2-D gel databases**

| 2DBase-Ecoli | P0A6N1. |
|---|---|
| ECO2DBASE | E042.0. 6TH EDITION. |

**P0CE47** (EFTU1_ECOLI) ⭐ Reviewed, UniProtKB/Swiss-Prot
Last modified November 30, 2010. Version 7. 🔶 History...

MSKE...
AFDQ...
NMITGAAQM...VPYIIVFL
NKCDMVDD...GSALKALE
GDAEWEAK...VFSISGRG
TVVTGRVE...RKLLDEGR
AGENVGVLLRGIKREEIERGQVLAKPGTIKPHTKFESEVYILSKD
EGGRHTPFFKGYRPQFYFRTTDVTGTIELPEGVEMVMPGDNIKMV
VTLIHPIAMDDGLRFAIREGGRTVGAGVVAKVLG

**One protein sequence**
**One gene**
**One species**

**Sequence annotation (Features)**

| Feature key | Position(s) | Length | Description | Graphical view |
|---|---|---|---|---|
| **Molecule processing** | | | | |
| ☐ Initiator methionine | 1 | 1 | Removed Ref.4<br>Ref.5 Ref.6 | |
| ☐ Chain | 2 - 394 | 393 | Elongation factor Tu | |

**Manual annotation**
**Post-translational modifications,**
**variants, transmembrane domains,**
**signal peptide...**

**Amino acid modifications**

| ☐ Modified residue | 2 | 1 | N-acetylserine HAMAP MF_00118_B | |
|---|---|---|---|---|
| ☐ Modified residue | 57 | 1 | N6,N6-dimethyllysine; alternate Ref.11 | |

**Alternative products:**
**protein sequences produced by**
**alternative splicing,**
**alternative promoter usage,**
**alternative initiation...**

**Ontologies**

**Keywords**

| Biological process | Antibiotic resistance<br>Protein biosynthesis |
|---|---|
| Cellular component | Cell membrane<br>Cytoplasm<br>Membrane |
| Ligand | GTP-binding<br>Nucleotide-binding |
| Molecular function | Elongation factor |
| PTM | Acetylation<br>Methylation<br>Phosphoprotein |
| Technical term | 3D-structure<br>Complete proteome<br>Direct protein sequencing |

**Gene Ontology (GO)**

| Biological process | response to antibiotic<br>Inferred from electronic annotation. Source: UniProtKB-KW |
|---|---|
| Cellular component | cytoplasm<br>Inferred from direct assay. Source: UniProtKB<br>plasma membrane<br>Inferred from electronic annotation. Source: UniProtKB-SubCell |
| Molecular function | GTP binding<br>Inferred from electronic annotation. Source: UniProtKB-KW<br>GTPase activity<br>Inferred from electronic annotation. Source: InterPro<br>protein binding<br>Inferred from physical interaction. Source: IntAct<br>translation elongation factor activity<br>Inferred from electronic annotation. Source: UniProtKB-KW |

**Manual annotation**
**Keywords**
**and**
**Gene Ontology**

**UniProtKB/Swiss-Prot**
**www.uniprot.org**
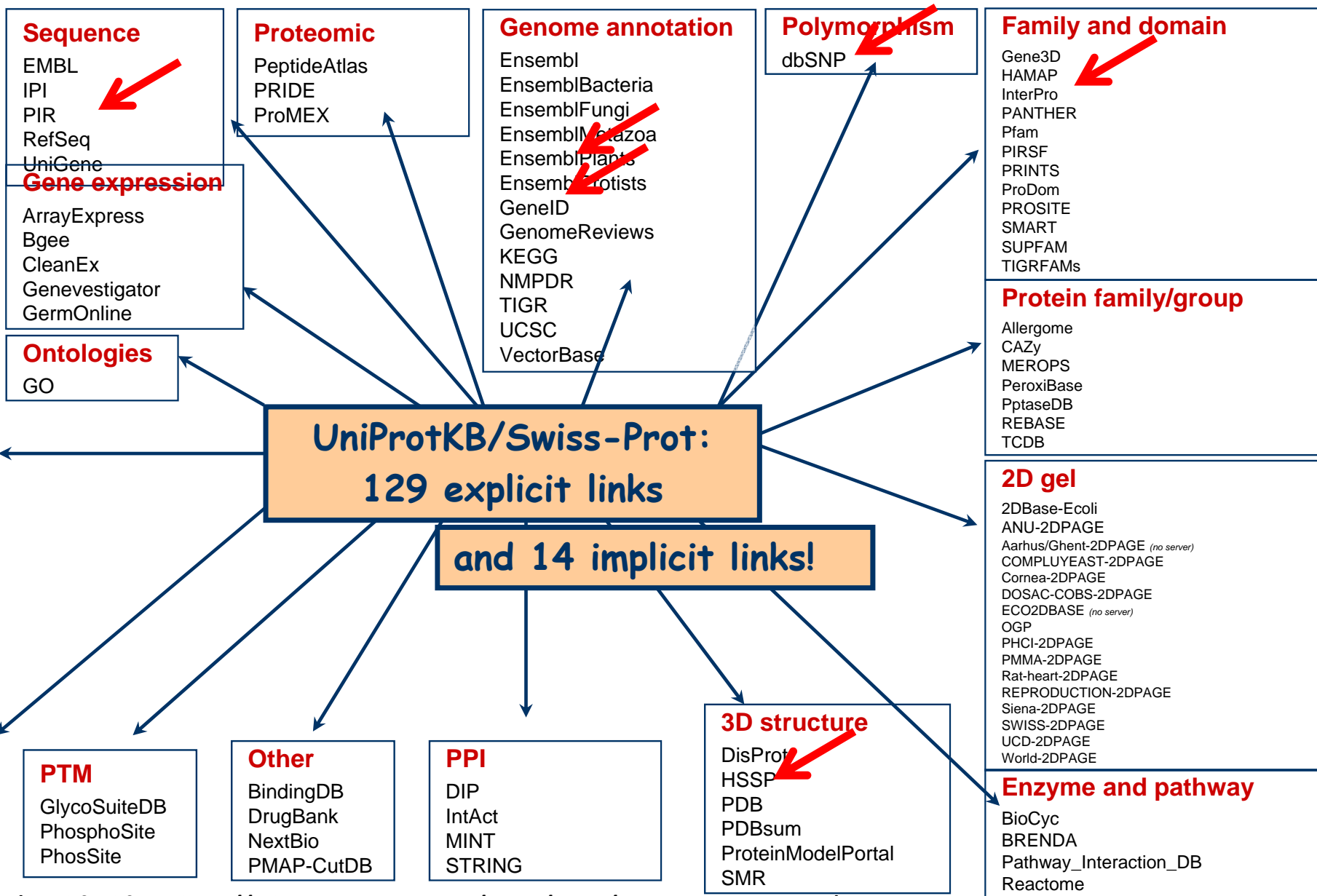
**Organism-specific**

AGD
ArachnoServer
CGD
ConoServer
CTD
CYGD
dictyBase
EchoBASE
EcoGene
euHCVdb
EuPathDB
FlyBase
GeneCards
GeneDB_Spombe
GeneFarm
GenoList
Gramene
H-InvDB
HGNC
HPA
LegioList
Leproma
MaizeGDB
MGI
MIM
neXtProt
Orphanet
PharmGKB
PseudoCAP
RGD
SGD
TAIR
TubercuList
WormBase
Xenbase
ZFIN

**Phylogenomic dbs**

eggNOG
GeneTree
HOGENOM
HOVERGEN
InParanoid
OMA
OrthoDB
PhylomeDB
ProtClustDB

**Sequence**

EMBL
IPI
PIR
RefSeq
UniGene

**Gene expression**

ArrayExpress
Bgee
CleanEx
Genevestigator
GermOnline

**Ontologies**

GO

**Proteomic**

PeptideAtlas
PRIDE
ProMEX

**Genome annotation**

Ensembl
EnsemblBacteria
EnsemblFungi
EnsemblMetazoa
EnsemblPlants
EnsemblProtists
GeneID
GenomeReviews
KEGG
NMPDR
TIGR
UCSC
VectorBase

**Polymorphism**

dbSNP

**Family and domain**

Gene3D
HAMAP
InterPro
PANTHER
Pfam
PIRSF
PRINTS
ProDom
PROSITE
SMART
SUPFAM
TIGRFAMs

**Protein family/group**

Allergome
CAZy
MEROPS
PeroxiBase
PptaseDB
REBASE
TCDB

**2D gel**

2DBase-Ecoli
ANU-2DPAGE
Aarhus/Ghent-2DPAGE *(no server)*
COMPLUYEAST-2DPAGE
Cornea-2DPAGE
DOSAC-COBS-2DPAGE
ECO2DBASE *(no server)*
OGP
PHCI-2DPAGE
PMMA-2DPAGE
Rat-heart-2DPAGE
REPRODUCTION-2DPAGE
Siena-2DPAGE
SWISS-2DPAGE
UCD-2DPAGE
World-2DPAGE

**Enzyme and pathway**

BioCyc
BRENDA
Pathway_Interaction_DB
Reactome

## UniProtKB/Swiss-Prot: 129 explicit links

## and 14 implicit links!

**PTM**

GlycoSuiteDB
PhosphoSite
PhosSite

**Other**

BindingDB
DrugBank
NextBio
PMAP-CutDB

**PPI**

DIP
IntAct
MINT
STRING

**3D structure**

DisProt
HSSP
PDB
PDBsum
ProteinModelPortal
SMR

**(Modified from http://education.expasy.org/cours/Turin/UniProtKB_Turin.ppt)**

Summary

Domain(s)

Hits summary

**Sequences producing significant alignments:**

Select: All None   Selected:0

| Description | Max score | Total score | Query cover | E value | Ident | Accession |
|---|---|---|---|---|---|---|
| ☐ RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 286 | 286 | 100% | 6e-100 | 100% | P69905.2 |
| ☐ RecName: Full=Hemoglobin subunit theta-1; AltName: Full=Hemoglobin theta-1 chain; AltName: Full=Theta-1-globin | 182 | 182 | 100% | 1e-58 | 62% | P09105.2 |
| ☐ RecName: Full=Hemoglobin subunit zeta; AltName: Full=HBAZ; AltName: Full=Hemoglobin zeta chain; AltName: Full=Z | 176 | 176 | 100% | 2e-56 | 60% | P02008.2 |
| ☐ RecName: Full=Hemoglobin subunit mu; AltName: Full=Hemoglobin mu chain; AltName: Full=Mu-globin | 135 | 135 | 99% | 2e-40 | 45% | Q6B0K9.1 |
| ☐ RecName: Full=Hemoglobin subunit delta; AltName: Full=Delta-globin; AltName: Full=Hemoglobin delta chain | 114 | 114 | 97% | 2e-32 | 43% | P02042.2 |
| ☐ RecName: Full=Hemoglobin subunit beta; AltName: Full=Beta-globin; AltName: Full=Hemoglobin beta chain; Contains: | 114 | 114 | 97% | 2e-32 | 43% | P68871.2 |
| ☐ RecName: Full=Hemoglobin subunit gamma-1; AltName: Full=Gamma-1-globin; AltName: Full=Hb F Agamma; AltNam | 113 | 113 | 97% | 6e-32 | 41% | P69891.2 |
| ☐ RecName: Full=Hemoglobin subunit gamma-2; AltName: Full=Gamma-2-globin; AltName: Full=Hb F Ggamma; AltNam | 113 | 113 | 97% | 6e-32 | 41% | P69892.2 |
| ☐ RecName: Full=Hemoglobin subunit epsilon; AltName: Full=Epsilon-globin; AltName: Full=Hemoglobin epsilon chain | 101 | 101 | 95% | 2e-27 | 39% | P02100.2 |
| ☐ RecName: Full=Cytoglobin; AltName: Full=Histoglobin; Short=HGb; AltName: Full=Stellate cell activation-associated pro | 68.9 | 68.9 | 96% | 5e-15 | 28% | Q8WWM9.1 |
| ☐ RecName: Full=Myoglobin | 51.2 | 51.2 | 100% | 5e-09 | 28% | P02144.2 |
| ☐ RecName: Full=Neurobeachin-like protein 1; AltName: Full=Amyotrophic lateral sclerosis 2 chromosomal region candi | 27.7 | 27.7 | 19% | 2.7 | 37% | Q6ZS30.3 |
| ☐ RecName: Full=StAR-related lipid transfer protein 9; AltName: Full=START domain-containing protein 9; Short=StARD9 | 26.6 | 26.6 | 40% | 6.5 | 30% | Q9P2P6.3 |
| ☐ RecName: Full=Intraflagellar transport protein 140 homolog; AltName: Full=WD and tetratricopeptide repeats protein 2 | 26.6 | 26.6 | 75% | 7.0 | 27% | Q96RY7.1 |
| ☐ RecName: Full=Ubiquitin carboxyl-terminal hydrolase 34; AltName: Full=Deubiquitinating enzyme 34; AltName: Full=Ub | 26.2 | 26.2 | 77% | 8.6 | 24% | Q70CQ2.2 |

HBA_HUMAN



HBB_HUMAN

Needleman-Wunsch
Global alignment

Smith-Waterman
Local alignment

BLAST
Local alignment

```
HBA_HUMAN    1 MV-LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-D   48
               || |:|.:|:.|.|.|||    :..|.|.|||.|.::.:|.|.:|..|  |
HBB_HUMAN    1 MVHLTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGD   48

HBA_HUMAN   49 LS-----HGSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLR   93
               ||      .|:.:||.|||||..|.::.:||:|::...:.||:||..||.
HBB_HUMAN   49 LSTPDAVMGNPKVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLH   98

HBA_HUMAN   94 VDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR  142
               |||.||:||.|.:|:..||.|....||||.|.|:..|.:|.|:..|..||.
HBB_HUMAN   99 VDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH  147


HBA_HUMAN    3 LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-DLS-   50
               |:|.:|:.|.|.||||    :..|.|.|||.|.::.:|.|.:|..|  |||
HBB_HUMAN    4 LTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGDLST   51

HBA_HUMAN   51 ----HGSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDP   96
               .|:.:||.|||||..|.::.:||:|::.....:.||:||..||.|||
HBB_HUMAN   52 PDAVMGNPKVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDP  101

HBA_HUMAN   97 VNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKY      141
               .||:||.|.:|:..||.|....||||.|.|:..|.:|.|:..|..||
HBB_HUMAN  102 ENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKY     146


Query    3 LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-DLS-----HGSAQV   56
           L+P +K+ V A WGKV  +   E G EAL R+ + +P T+ +F  F DLS       G+ +V
Sbjct    4 LTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKV   61

Query   57 KGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPA  116
           K HGKKV  A ++ +AH+D++     + LS+LH   KL VDP NF+LL + L+  LA H
Sbjct   62 KAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGK  121

Query  117 EFTPAVHASLDKFLASVSTVLTSKY  141
           EFTP V A+  K +A V+  L  KY
Sbjct  122 EFTPPVQAAYQKVVAGVANALAHKY  146
```

Protein BLAST: sea ×

blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastp&BLAST_PROGRAMS=blastp&PAGE_TYPE=BlastSearch&SHOW

**BLAST®**   Basic Local Alignment Search Tool

Home | Recent Results | Saved Strategies | 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.

▸ NCBI/ BLAST/ blastp suite

**Standard Protein BLAST**

blastn | blastp | blastx | tblastn | tblastx

BLASTP programs search protein databases using a protein query. more...

Reset page    Bookmark

## Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) ⓘ    Clear    Query subrange ⓘ

```
>sp|P69905|HBA_HUMAN
MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVK
GHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHL
PAEFTPAVHASLDKFLASVSTVLTSKYR
```

From [          ]

To [          ]

Or, upload file    [Choose File] No file chosen    ⓘ

Job Title    [sp|P69905|HBA_HUMAN                    ]
Enter a descriptive title for your BLAST search ⓘ

☐ Align two or more sequences ⓘ

## Choose Search Set

Database    ◆ [UniProtKB/Swiss-Prot(swissprot)        ▼] ⓘ

**Title:** Non-redundant UniProtKB/SwissProt sequences.
**Molecule Type:** Protein
**Update date:** 2013/10/10
**Number of sequences:** 456016

Organism
Optional    [Enter organism name or id--completions will be suggested] ☐ Exclude ＋
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 ⓘ

## Program Selection

Algorithm    ⦿ blastp (protein-protein BLAST)
○ PSI-BLAST (Position-Specific Iterated BLAST)
○ PHI-BLAST (Pattern Hit Initiated BLAST)
○ DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)
Choose a BLAST algorithm ⓘ

**BLAST**    Search database UniProtKB/Swiss-Prot(swissprot) 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

Descriptions

**Sequences producing significant alignments:**

Select: All None  Selected:0

Alignments   Download ∨   GenPept  Graphics  Distance tree of results  Multiple alignment

| Description | Max score | Total score | Query cover | E value | Ident | Accession |
|---|---|---|---|---|---|---|
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain >sp|P69 | 286 | 286 | 100% | 9e-99 | 100% | P69905.2 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 282 | 282 | 99% | 5e-97 | 99% | P01923.1 |
| RecName: Full=Hemoglobin subunit alpha-1; AltName: Full=Alpha-1-globin; AltName: Full=Hemoglobin alpha-1 chain | 281 | 281 | 100% | 1e-96 | 99% | Q9TS35.2 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 281 | 281 | 100% | 2e-96 | 98% | P06635.2 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 278 | 278 | 99% | 1e-95 | 98% | P01924.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain >sp|P63 | 278 | 278 | 100% | 1e-95 | 97% | P63107.2 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 277 | 277 | 100% | 5e-95 | 96% | P67817.2 |
| RecName: Full=Hemoglobin subunit alpha-2; AltName: Full=Alpha-2-globin; AltName: Full=Hemoglobin alpha chain | 277 | 277 | 100% | 6e-95 | 96% | Q9TS34.2 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 276 | 276 | 99% | 6e-95 | 97% | P18972.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 276 | 276 | 100% | 2e-94 | 96% | P01926.2 |
| RecName: Full=Hemoglobin subunit alpha-A/Q/R/T; AltName: Full=Alpha-A/Q/R/T-globin; AltName: Full=Hemoglobin alph | 275 | 275 | 99% | 2e-94 | 97% | P21767.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 275 | 275 | 100% | 2e-94 | 96% | P01928.2 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 275 | 275 | 99% | 4e-94 | 96% | P67818.1 |
| RecName: Full=Hemoglobin subunit alpha-1/2/3; AltName: Full=Alpha-1/2/3-globin; AltName: Full=Hemoglobin alpha-1/2 | 275 | 275 | 99% | 5e-94 | 97% | P21766.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 274 | 274 | 99% | 5e-94 | 96% | P01929.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 273 | 273 | 99% | 2e-93 | 96% | P07421.1 |
| RecName: Full=Hemoglobin subunit alpha-1/2; AltName: Full=Alpha-1/2-globin; AltName: Full=Hemoglobin alpha-1/2 ch | 273 | 273 | 99% | 2e-93 | 96% | P21768.1 |
| RecName: Full=Hemoglobin subunit alpha-1/2/3; AltName: Full=Alpha-1/2/3-globin; AltName: Full=Hemoglobin alpha-1/2 | 272 | 272 | 99% | 3e-93 | 96% | P19002.1 |
| RecName: Full=Hemoglobin subunit alpha-1/2; AltName: Full=Alpha-1/2-globin; AltName: Full=Hemoglobin alpha-1/2 ch | 271 | 271 | 99% | 7e-93 | 96% | P07402.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 271 | 271 | 100% | 1e-92 | 95% | P01930.2 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 271 | 271 | 99% | 1e-92 | 96% | Q7M3B6.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 266 | 266 | 99% | 6e-91 | 94% | P01937.1 |
| RecName: Full=Hemoglobin subunit alpha-B; AltName: Full=Alpha-B-globin; AltName: Full=Alpha-II; AltName: Full=Hemo | 266 | 266 | 100% | 8e-91 | 92% | P01939.3 |
| RecName: Full=Hemoglobin subunit alpha-3; AltName: Full=Alpha-3-globin; AltName: Full=Hemoglobin alpha-2 chain | 266 | 266 | 99% | 1e-90 | 89% | P01935.1 |
| RecName: Full=Hemoglobin subunit alpha-3; AltName: Full=Alpha-3-globin; AltName: Full=Hemoglobin alpha-2 chain | 265 | 265 | 99% | 2e-90 | 89% | P01934.1 |
| RecName: Full=Hemoglobin subunit alpha-A; AltName: Full=Alpha-A-globin; Short=Alpha-I; AltName: Full=Hemoglobin al | 265 | 265 | 100% | 3e-90 | 92% | P14259.2 |
| RecName: Full=Hemoglobin subunit alpha-1/2; AltName: Full=Alpha-1/2-globin; AltName: Full=Hemoglobin alpha-1/2 ch | 265 | 265 | 99% | 4e-90 | 94% | P08258.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 264 | 264 | 99% | 6e-90 | 93% | P01938.2 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 264 | 264 | 99% | 6e-90 | 91% | P01940.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 263 | 263 | 99% | 1e-89 | 91% | P11757.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain >sp|P63 | 262 | 262 | 100% | 5e-89 | 92% | P63111.2 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 261 | 261 | 99% | 6e-89 | 93% | P01933.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 261 | 261 | 99% | 9e-89 | 90% | P11753.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 260 | 260 | 99% | 2e-88 | 90% | P01956.1 |
| RecName: Full=Hemoglobin subunit alpha; AltName: Full=Alpha-globin; AltName: Full=Hemoglobin alpha chain | 259 | 259 | 100% | 3e-88 | 90% | P14387.2 |

# Summary Questions

- Why do we need to perform a database searching?

- What's the major challenge/obstacle when searching sequence database?

# 生物信息学：导论与方法
## Bioinformatics: Introduction and Methods



https://www.coursera.org/course/pkubioinfo