



# Sequencing alignment

Ameer Effat M. Elfarash

Dept. of Genetics  
Fac. of Agriculture, Assiut Univ.  
[aelfarash@aun.edu.eg](mailto:aelfarash@aun.edu.eg)

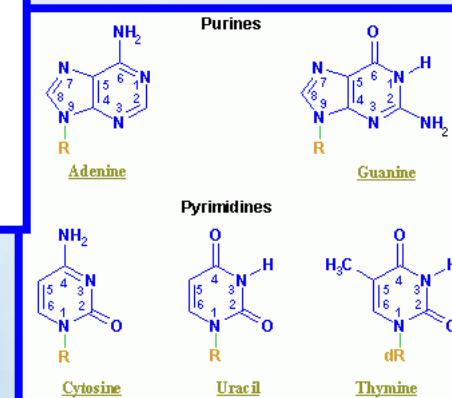
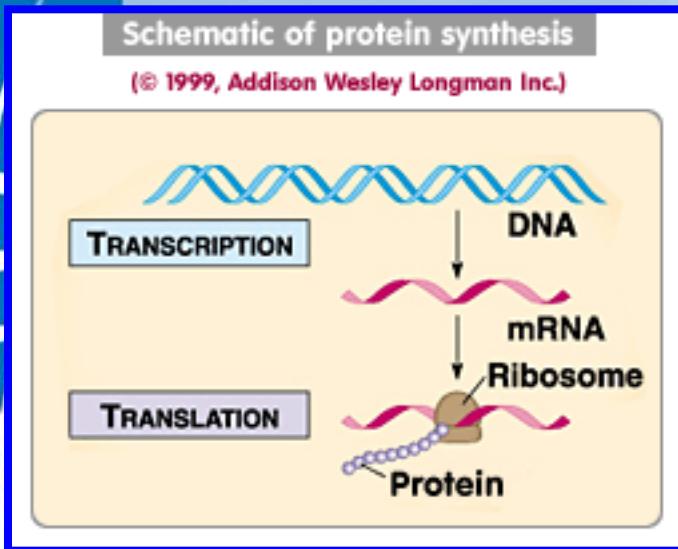


# Why perform a multiple sequence alignment?

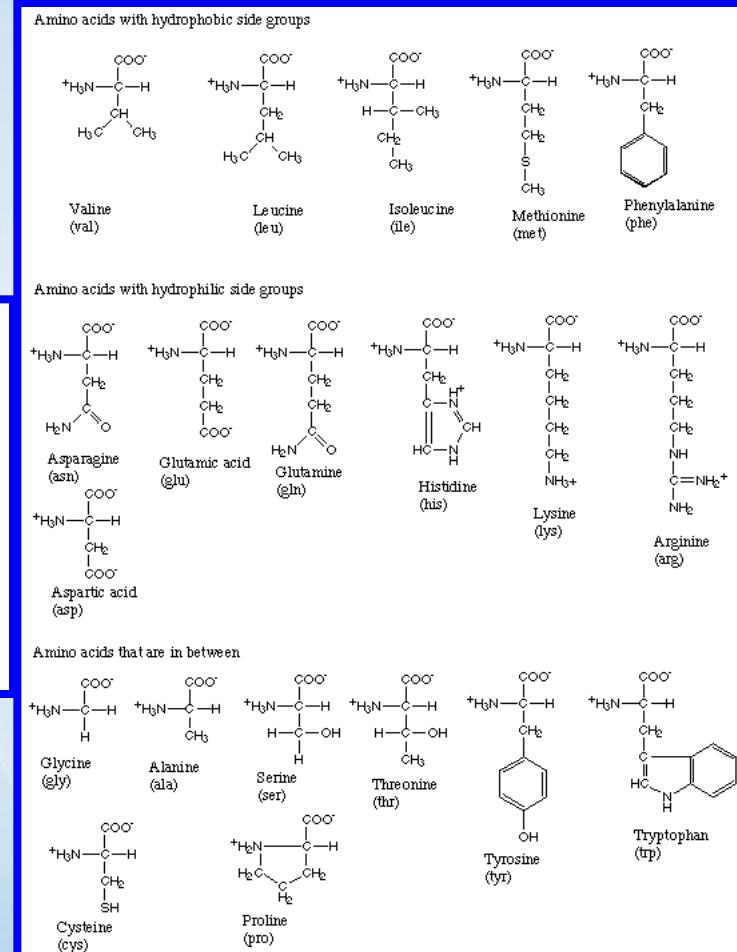
MSAs are at the heart of comparative genomics studies which seek to study evolutionary histories, functional and structural aspects of sequences, and to understand phenotypic differences between species



# The Building Blocks...



ATGC



VLMFNQEDHKRCSTPYW



## Three types of nucleotide changes:

**Substitution** – a replacement of one (or more) sequence characters by another: **AAGA** → **AACA**

**Insertion** - an insertion of one (or more) sequence characters: **AAG A**

**Deletion** – a deletion of one (or more) sequence characters: **AA AGA**

**Alignment:** Comparing two (pairwise) or more (multiple) sequences.

Searching for a series of identical or similar characters in the sequences.



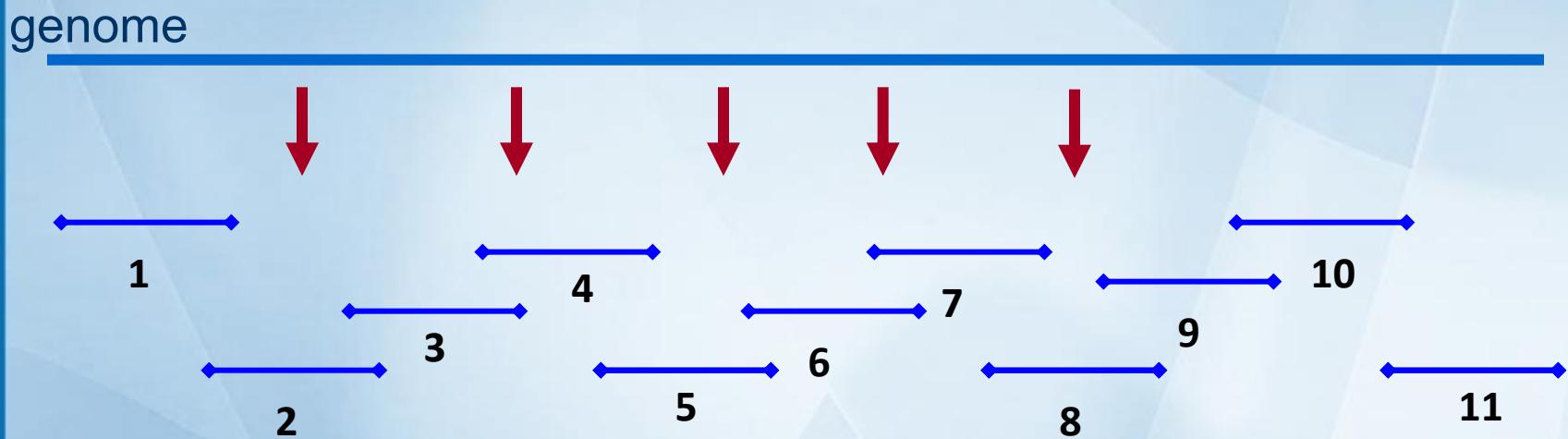
# Why Align Sequences?

- Crucial for genome sequencing
- Discover functional, structural, and evolutionary information because similar Sequences may have similar function
- Identify primers and probes to search for homologous sequences in other organisms



# Why Align Sequences?

- Genome Sequencing Strategy

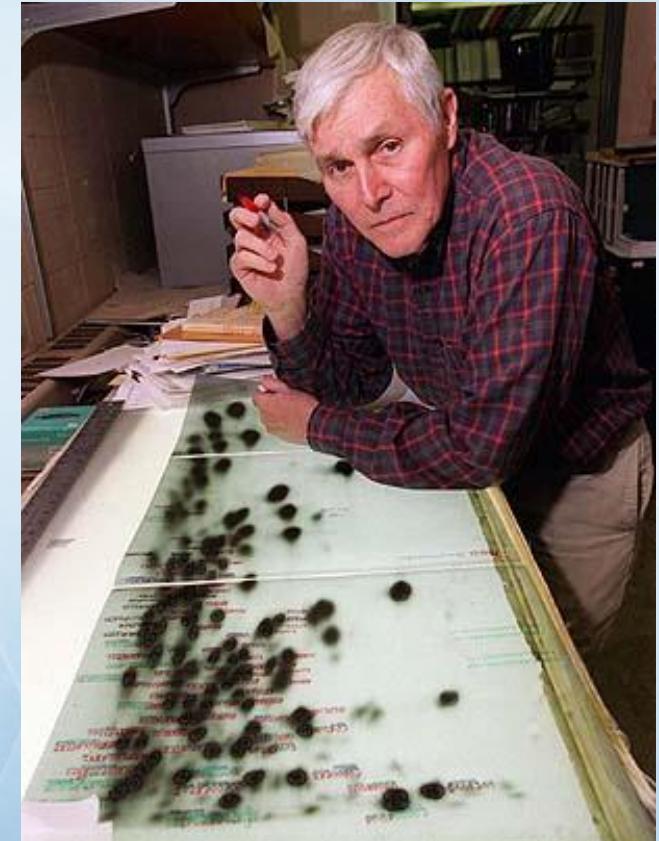


1. Obtain a large collection of BAC clones
2. Map them onto the genome (Physical Mapping)
3. Select a minimum tiling path
4. Sequence each clone in the path with shotgun
5. Assemble
6. Put everything together



# Why Align Sequences?

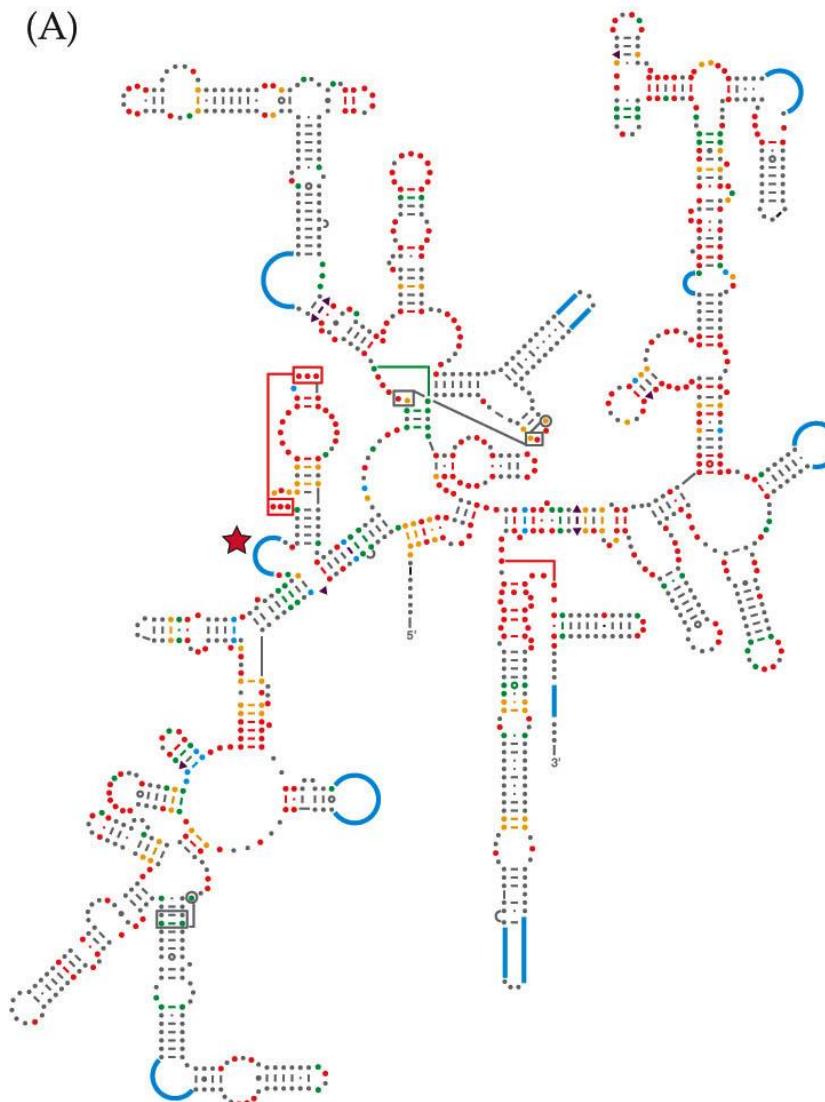
- Homology
  - Similar sequences may have a common ancestor
- In 1990 Carl Woese and colleagues proposed the Tree of Life, using a molecular phylogenetic approach.
- It is based on sequencing rRNA (16S and 18S), which all organisms share.
- All organisms were separated into three domains: *Bacteria*, *Archaea*, *Eukarya*.





# 16S RNA

(A)



- Identical in 98% or more of all organisms
- Conserved only in the *Bacteria*
- Conserved only in the *Archaea*
- Conserved only in the *Eukarya*
- ▲ Conserved within each domain, variable among domains
- Regions that vary structurally among domains



# How DNA Sequence Data is Compared

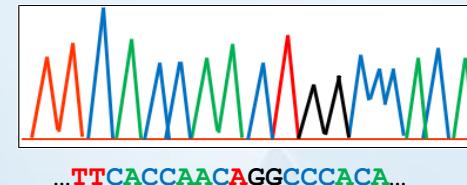
*Obtain Samples:*  
*Blood , Saliva, Hair*  
*Follicles, Feathers, Scales*



## *Genetic Data*

*Extract DNA from Cells*

*Sequence DNA*



***Compare DNA Sequences to One Another***

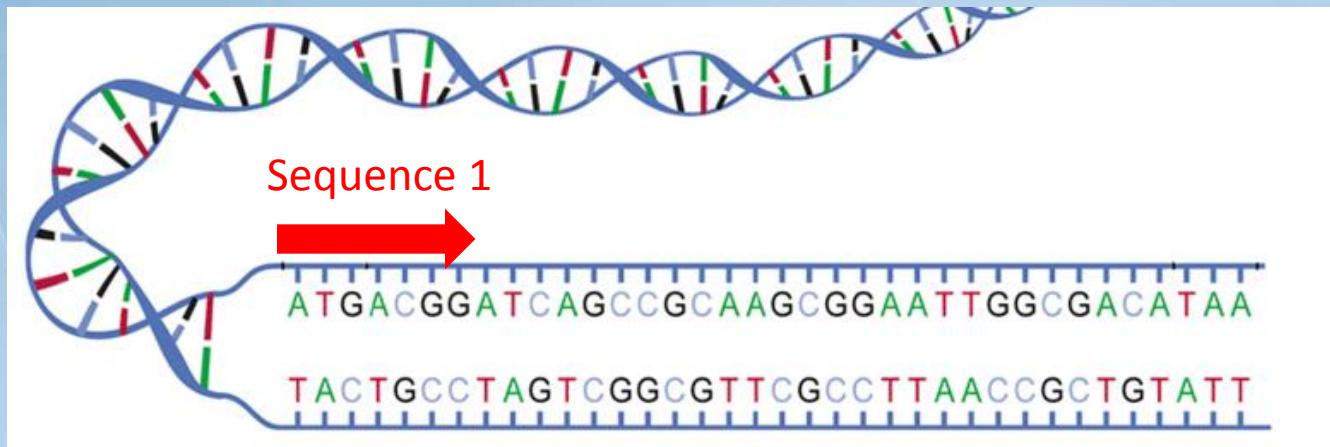
TTCAACAAACAGGCCAC  
TTCACCAACAGGCCAC  
TTCATCAACAGGCCAC

## **GOALS:**

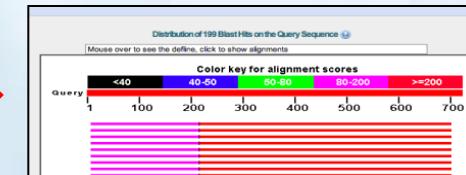
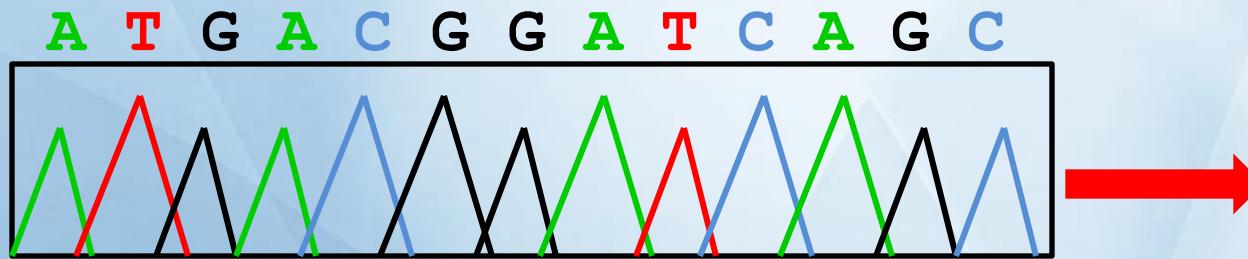
- Identify the organism from which the DNA was obtained.
- Compare DNA sequences to each other.



# confirm insert sequence



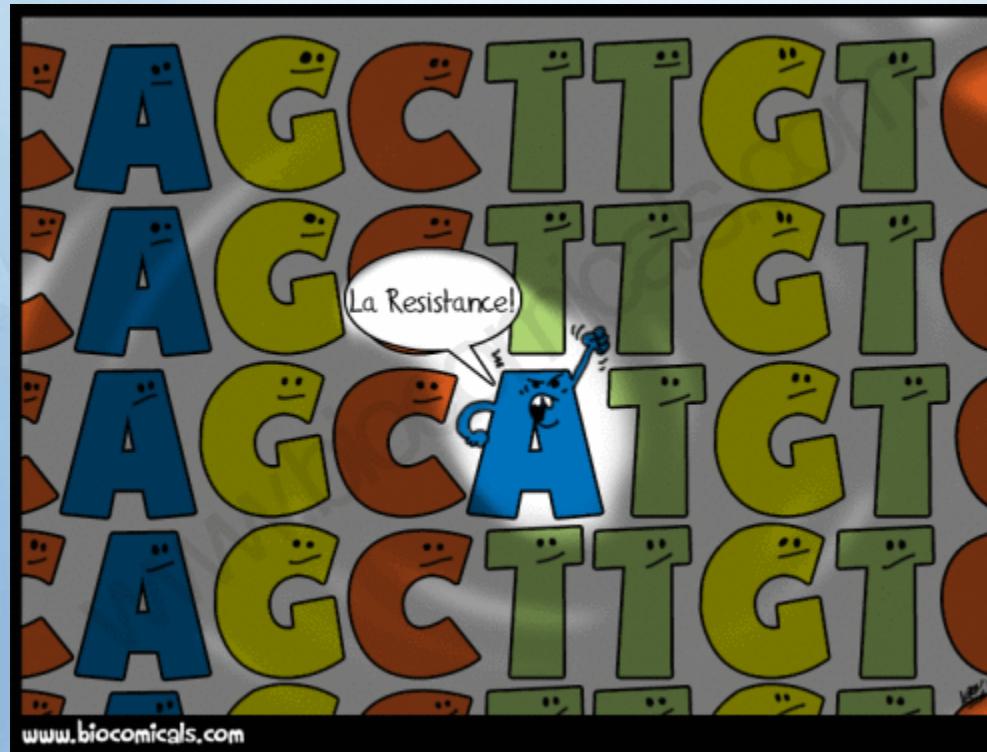
Sequence :





# Comparing two sequences

MVNLT**SDEKTAVLALWNKVDVEDCGGE**  
||| ||| | | | | | | | | | | | | | | |  
MVHLT**PEEKTAVNALWGKVNVDAVGGE**





# Multiple sequence alignment

**Seq1** VTISCTGSSSNIGAG-NHVWKWYQQLPG

**Seq2** VTISCTGTSSNIIGS--ITVNWYQQLPG

**Seq3** LRLSCSSSGFIFSS--YAMYWVRQAPG

**Seq4** LSLTCTVSGTSFDD--YYSTWVRQPPG

**Seq5** PEVTCVVVDVSHEDPQVKFNWYVDG--

**Seq6** ATLVC LISDFYPGA--VTVAWKADS--

**Seq7** AALGCLVKDYFPEP--VTWSWNSG---

**Seq8** VSLTCLVKGFYPSD--IAVEWWNSG--

Each row represents an individual sequence  
Each column represents the ‘same’ position



# Multiple sequence alignment

**Seq1** VTIS**C**TGSSSN**I****GAG**-NHV**KW**YQ**QL****PG**

**Seq2** VTIS**C**TGTSSNI**GS**--ITV**NW**YQ**QL****PG**

**Seq3** LRLS**C**SSSGFIF**SS**--YAM**YW**VR**QA****PG**

**Seq4** LSLT**C**TVSGTSF**DD**--YYST**TW**VR**QP****PG**

**Seq5** PEVT**C**VVVVDVSH**ED**PQVKF**NW**YV**DG**--

**Seq6** ATL**C**LISDFYP**GA**--VTV**AW**KADS--

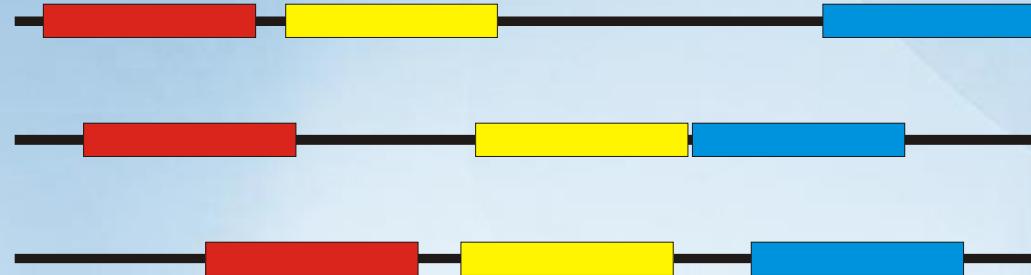
**Seq7** AALG**C**LVKDYFP**EP**--VTV**SW**NSG---

**Seq8** VSLT**C**LVKGFY**PSD**--IAV**EW**WS**NG**--

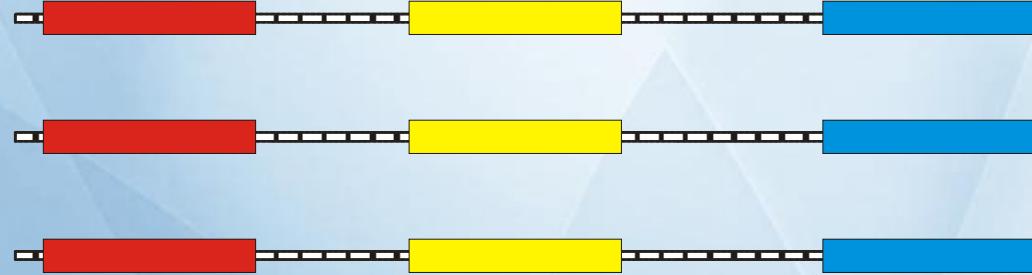


## How to optimize alignment algorithms?

- ◆ Sequences often contain highly conserved regions



These regions can be used for an initial alignment



By analyzing a number of small, independent fragments,  
the algorithmic complexity can be drastically reduced!



# Choosing an alignment:

- Many **different** alignments between two sequences are possible:

AAGCTGAATTCTGAA

AGGCTCATTCTGA

AAGCGAAA**T**TCGAAC  
A-G-GAA-**C**TCGAAC

**A**AGCGAAA**T**TCGAAC  
**AGG**---**AA****C**TCGAAC

**How do we determine which is the best alignment?**



# Assessing the significance of an alignment score

True

AAGCTGAATCGAA  
AGGCTCATTCTGA

AAGCTGAATC-GAA  
AGGCTCATTCTGA-

28.0

Random

AGATCAGTAGACTA  
GAGTAGCTATCTCT

AGATCAGTAGACTA-----  
-----GAGTAG-CTATCTCT

26.0

CGATAGATAGCATA  
GCATGTCATGATTG

CGATAGATAGCATA-----  
-----GCATGTCATGATTG

16.0



# Distance matrix: UPGMA

- UPGMA = unweighted pair group method with arithmetic mean

(A) UPGMA method

Table shows sequence of nine-base region of rRNAs of four strains.

Organism (strain)	Site number								
	1	2	3	4	5	6	7	8	9
a	G	C	G	G	A	C	A	A	A
b	G	A	C	G	C	C	A	A	G
c	G	A	A	A	U	C	U	A	A
d	G	A	A	A	G	C	U	A	G

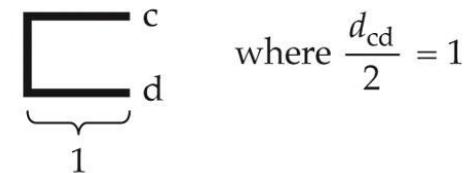
First matrix

1 Construct matrix showing relatedness between strains.

	a	b	c
b	$d_{ab} = 4$	—	—
c	$d_{ac} = 5$	$d_{bc} = 5$	—
d	$d_{ad} = 6$	$d_{bd} = 4$	$d_{cd} = 2$

Beginning tree

2 Diagram relatedness between strains.





1

Courier New

P

84 total sequences

Start  
ruler at: 1

scroll    
speed slow   fast

Sequence alignment across the structure, showing nucleotide positions 140 to 230:

Species	140	150	160	170	180	190	200	210	220	230	
Adriohydrobi	CCAACCGCUCCGACCC			CUG	ACGGG		AAAGAGC	GC	UUUUAUCAGCUCA		
Adrioinsulan	CCAACCGCUCCGACCC			CUUCAACGGG			AAAGAGC	GC	UUUUAUCAGCUCA		
Alvania	CAAACAGCUCCGACCC			UCA			GGGAAAGAGC	GC	UUUUAUUAGUUCA		
Alzoniella 2	CCAACCGCUCCGACCC			UUC	ACGGG		AAAGAGC	GC	UUUUAUCAGCUCA		
Amnicola 106	CUACCAAGCUCCGACCCGGUGGGCCUCGCUUCCGCUUUCGGUUCACAGGGGGAGUCGGGUCCCCA			UUC	ACGGG		AAAGAGC	GC	UUUUAUUAGUUCA		
Amphithalamu	CUACCAAGCUCCGACCC			GUGGU			AAAGCCAGGGAGAGC	GC	UUUUAUUAGUUCA		
Antroselates	CUACCAAGCUCCGACCCGGUGGGCCUCGGUUCCGCUUUCGGUUCACAGGGGGAGUCGGGUCCCCA			UUC	ACGGG		AAAGCCAGGGAGAGC	GC	UUUUAUUAGUUCA		
Ascorhis	AUACAAAGCUCCGACCC			UUC	ACGGG		ACAGAGC	GC	UUUUAUUAGUUCA		
Assiminea	CCACCAAGCUCCGACCC	UGG		UUUCGG			UCAGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Assiminea 16	CCACCAAGCUCCGACCC			GGUCUCU	CGAG		GCCAGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Assiminea 22	CCACCAAGCUCCGACCC			GGUUUC	U		GUCAGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Avenionia 22	CCAACCGCUCCGACCC			UUC	ACGGG		AAAGAGC	GC	UUUUAUCAGCUCA		
Baicalia	UAACAGCUCCGACCC			CCC	UCA		GGGGCCC	GGGAAAGAGC	GC	UUUUAUUAGUUCA	
Barleeria	CCCCCAAGCUCCGACCCG			UUCC			12: Avenionia_2241	GAAAGAGC	GC	UUUUAUUAGUUCA	
Beddomeia	CUAACAGCUCCGACCC			GCA			AGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Belgrandia	CCAACCAGCUCCGACCC			UGC	AAAGG		AAAGAGC	GC	UUUUAUCAGCUCA		
Bithynia	CCAACAGCUCCGACCC			UUC			AACGGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Bythinella 1	CCAACAGCUCCGACCC			CUUCG	CAAGGAGG		GGAAAGAGC	GC	UUUUAUUAGUUCA		
Bythiospeum	CUAACAGCUCCGACCC			UCA			CGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Celopria	UAGUAGCUU	GACCC		UCA			CCGGGAAGAGC	GC	UUUUAUUAGUUCA		
Cecina 2522	CCACCAAGCUCCGACCC	U		GGUCA			AACAGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Clenchiella	CUAACAGCUCCGACCC			UCA			GGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Coxiella	CCACCAAGCUCCGACCC	U		AGUC	UUU	CGAG	GCCGGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Eatonella	GCAAACCUA	UG	UCCG	ACUCC	U		UGGGAAAGAU	UAGAGCC	ACUUUAUUAGUUCA		
Emmericia	CUAACAGCUCCGACCC			CUCAC			GGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Emmericia 30	UAACAGCUCCGACCC			CUCAC			GGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Erhaia 652	CCCCAAGCUCCGACCCG			CUCUUCGCCGG			GCGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Fairbankia	CCAACAGCUCCGACCC			UCA			GGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Fissuria 243	CCAACCGCUCCGACCC			UUG	ACGGG		AAAGAGC	GC	UUUUAUCAGCUCA		
Fluvidiona	CCCCAAGCUCCGACCC			GU			CACGGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Fluvipupa	CCCACAGCUCCGACCC	UG		GGGUUCUGC			CCCGUGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Fontigens	ACCCCAAGCUCCGACCC	UUG		ACCCA			GGGGAAAGAGC	GC	UUUUAUUAGCUCG		
Gammaticula	CCACCAAGCUCCGACCC	U		GUCA			AACAGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Geomelania 8	CCACCAAGCUCCGACCCGG			CCGCC	UGCUU	CACGGGGCAGGUCUGG	AAGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Geomelania 8	CCACCAAGCUCCGACCCGG			CCGCC	UGCUU	CACGGGGCAGGUCUGG	AAGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Graziana 256	CCAACAGCUCCGACCC			CGC	AAAGG		AAAGAGC	GC	UUUUAUCAGCUCA		
Hauffenia 25	CCAACCGCUCCGACCC			UUC	ACAGG		AAAGAGC	GC	UUUUAUCAGCUCA		
Heleobops	UAACAGCUCCGACCC			UCA			GGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Hemistomia	CCACCAAGCUCCGACCC			GU			CCUGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Heterocyclus	CCGCCAGCUCCGACCC	UG		CGCUGAACG			GGGGAGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Horatia 2598	CCACCAAGCUCCGACCC			UGC	AAA	GGG	AAAGAGC	GC	UUUUAUCAGCUCA		
Hydrobia 653	CCACCAAGCUCCGACCC			CGC	AAAGG		AAAGAGC	GC	UUUUAUCAGCUCA		
Hydrococcus	CCACCAAGCUCCGACCC	U		GGUGA			GGGGAAAGAGC	GC	UUUUAUUAGUUCA		
Islamia 2327	CCAACCGCUCCGACCC			CUC	ACGGG		AAAGAGC	GC	UUUUAUCAGCUCA		



## “Optimal” vs. “correct” alignment

For a given group of sequences, there is no single “correct” alignment, only an alignment that is “optimal” according to some set of calculations

Determining what alignment is best for a given set of sequences is really up to the judgment of the investigator

Success of the alignment will depend on the similarity of the sequences. If sequence variation is great it will be very difficult to find an optimal alignment



# Web servers for pairwise alignment



# NCBI



C www.ncbi.nlm.nih.gov

NCBI Resources How To

[Sign in to NCBI](#)



National Center for  
Biotechnology Information

All Databases ▾

Search

NCBI Home

Resource List (A-Z)

All Resources

Chemicals & Bioassays

Data & Software

DNA & RNA

Domains & Structures

Genes & Expression

Genetics & Medicine

Genomes & Maps

Homology

Literature

Proteins

Sequence Analysis

Taxonomy

Training & Tutorials

Variation

## Welcome to NCBI

The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information.

[About the NCBI](#) | [Mission](#) | [Organization](#) | [Research](#) | [NCBI News](#)

## Get Started

- [Tools](#): Analyze data using NCBI software
- [Downloads](#): Get NCBI data or software
- [How-To's](#): Learn how to accomplish specific tasks at NCBI
- [Submissions](#): Submit data to GenBank or other NCBI databases

### NCBI YouTube channel

Learn how to get the most out of NCBI tools and databases with video tutorials on the NCBI YouTube Channel.



|| 1 2 3 4 5 6 7 8

## Popular Resources

[PubMed](#)

[Bookshelf](#)

[PubMed Central](#)

[PubMed Health](#)

**BLAST**

[Nucleotide](#)

[Genome](#)

[SNP](#)

[Gene](#)

[Protein](#)

[PubChem](#)

## NCBI Announcements

Coffee Break tutorial: Brown fat and obesity

Apr 1, 2014

The latest Coffee Break tutorial discusses EUMT1 or SPTAN1

New NCBI YouTube video: Create custom databases for BLAST

Mar 28, 2014

In the newest NCBI video on YouTube, we show you how to create custom



# BLAST – programs



**BLAST** Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

► NCBI/ BLAST Home

BLAST finds regions of similarity between biological sequences. [more...](#)

New Aligning Multiple Protein Sequences? Try the [COBALT Multiple Alignment Tool](#). [Go](#)

### BLAST Assembled Genomes

Choose a species genome to search, or [list all genomic BLAST databases](#).

<input type="checkbox"/> <a href="#">Human</a>	<input type="checkbox"/> <a href="#">Oryza sativa</a>	<input type="checkbox"/> <a href="#">Gallus gallus</a>
<input type="checkbox"/> <a href="#">Mouse</a>	<input type="checkbox"/> <a href="#">Bos taurus</a>	<input type="checkbox"/> <a href="#">Pan troglodytes</a>
<input type="checkbox"/> <a href="#">Rat</a>	<input type="checkbox"/> <a href="#">Danio rerio</a>	<input type="checkbox"/> <a href="#">Microbes</a>
<input type="checkbox"/> <a href="#">Arabidopsis thaliana</a>	<input type="checkbox"/> <a href="#">Drosophila melanogaster</a>	<input type="checkbox"/> <a href="#">Apis mellifera</a>

### Basic BLAST

Choose a BLAST program to run.

<a href="#">nucleotide blast</a>	Search a nucleotide database using a nucleotide query <i>Algorithms:</i> blastn, megablast, discontiguous megablast
<a href="#">protein blast</a>	Search protein database using a protein query <i>Algorithms:</i> blastp, psi-blast, phi-blast
<a href="#">blastx</a>	Search protein database using a translated nucleotide query
<a href="#">tblastn</a>	Search translated nucleotide database using a protein query
<a href="#">tblastx</a>	Search translated nucleotide database using a translated nucleotide query

**nucleotide blast** is circled in orange.

### 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](#) (IgBLAST)
- Search for [SNPs](#) (snp)
- Screen sequence for [vector contamination](#) (vecscren)
- Align two (or more) sequences using BLAST (bl2seq)
- Search [protein](#) or [nucleotide](#) targets in PubChem BioAssay
- Search SRA [transcript libraries](#)
- Constraint Based Protein [Multiple Alignment Tool](#)



# Query type: AA or DNA?

- For coding sequences, AA (protein) data are better
  - Selection operates most strongly at the protein level → the homology is more evident
  - AA – 20 char' alphabet      DNA - 4 char' alphabet



lower chance of random homology for AA



# BLAST – bl2seq

**BLAST** Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

► NCBBI BLAST Home

BLAST finds regions of similarity between biological sequences. [more...](#)

New Aligning Multiple Protein Sequences? Try the COBALT Multiple Alignment Tool. [Go](#)

### BLAST Assembled Genomes

Choose a species genome to search, or [list all genomic BLAST databases](#).

<input type="checkbox"/> <a href="#">Human</a>	<input type="checkbox"/> <a href="#">Oryza sativa</a>	<input type="checkbox"/> <a href="#">Gallus gallus</a>
<input type="checkbox"/> <a href="#">Mouse</a>	<input type="checkbox"/> <a href="#">Bos taurus</a>	<input type="checkbox"/> <a href="#">Pan troglodytes</a>
<input type="checkbox"/> <a href="#">Rat</a>	<input type="checkbox"/> <a href="#">Danio rerio</a>	<input type="checkbox"/> <a href="#">Microbes</a>
<input type="checkbox"/> <a href="#">Arabidopsis thaliana</a>	<input type="checkbox"/> <a href="#">Drosophila melanogaster</a>	<input type="checkbox"/> <a href="#">Apis mellifera</a>

### Basic BLAST

Choose a BLAST program to run.

<a href="#">nucleotide blast</a>	Search a nucleotide database using a nucleotide query <i>Algorithms:</i> blastn, megablast, discontiguous megablast
<a href="#">protein blast</a>	Search protein database using a protein query <i>Algorithms:</i> blastp, psi-blast, phr-blast
<a href="#">blastx</a>	Search protein database using a translated nucleotide query
<a href="#">tblastn</a>	Search translated nucleotide database using a protein query
<a href="#">tblastx</a>	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](#) (IgBLAST)
- Search for [SNPs](#) (snp)
- Screen sequence for [vector contamination](#) (vecscren)
- Align two (or more) sequences using BLAST (bl2seq)
- Search [protein](#) or [nucleotide](#) targets in PubChem BioAssay
- Search [SRA transcript library](#)
- Constraint Based Protein [Multiple Alignment Tool](#)



Sequence alignment - Wiki x Multiple Sequence Alignm x COBALT:Multiple Alignm x

www.ncbi.nlm.nih.gov/tools/cobalt/cobalt.cgi?link\_loc=BlastHomeLink

## COBALT Constraint-based Multiple Alignment Tool

My NCBI [Sign In] [Register]

### Cobalt Constraint-based Multiple Protein Alignment Tool

COBALT computes a multiple protein sequence alignment using conserved domain and local sequence similarity information. [?](#) [Reset page](#)

Enter Query Sequences

Enter at least 2 protein accessions, gis, or FASTA sequences [?](#) [Clear](#)

Or, upload FASTA file [Choose File](#) No file chosen

Job Title

**Align**  Show results in a new window

► [Advanced parameters](#)

Copyright | Disclaimer | Privacy | Accessibility | Contact | Send feedback

NCBI | NLM | NIH | DHHS



# results

**Alignments**  Select All [Get selected sequences NEW](#)

```

>lc0143385 gi|50902697|gb|IAAT86412.1| ABC transporter-associated protein
[Streptococcus pyogenes MGAS10394]
Length=472

Score = 394 bits (1013). Expect = 5e-114. Method: Compositional matrix adjust
Identities = 190/304 (62%). Positives = 232/304 (76%). Gaps = 4/304 (1%)

Query 566 PTYDIQVVGLENFVANGIVAHNSFIYVPPGVHVDIPLQAYFRINTENMGQFERTLIIADT 625
      P D ++ L + V +G +FIYVP GV VDIPLQ YFRIN EN GQFERTLII D
Sbjct 173 PPTDNKLAIRNSAVWSG----GTIYVPKGVKVDIPLQTYFRINNENTGQFERTLIIVDE 228

Query 626 GSYVWYVECTAPIYKSDSLRSAVVEIIIVKPHARVRYTTIQNWSNNVYNLVTKRARVETG 685
      G+ V+YVEC+TAP Y S+SLH+A+VEI A +RYTTIQNWS+NVYNLVTKRAR T
Sbjct 229 GASTHYVEGCTAPTYSSNSLHAATVEIFALDGAYMRYTTIQNWSNDVYNLVTKRARALTD 288

Query 686 AT+EWIDGN+SKVTMKYPAVWMTC+HAKGEVLSVAFAGEGQHQDTGAKMLHLASNTSSN 745
      AT+EWIDGN+G+K TMKYP+V++ G +R+G +L+S+AFA GQHQDTGAKM+H A +TSS+
Sbjct 289 AT+EWIDGN+GKTTMKYPSPVYLDGPC+RTGTMMSIAFANAGQHQDTGAKMIHNAPHTSSS 348

Query 746 IV+KSVARGGGR SYRGLVQVNKGAGHSR+SVKCDALLVDTISRSDTYPVDIRDDVTM 805
      IV+KS+A+A+GG+ YRG V NK + S S +CD +L+D IS+SDT P+ +I V +
Sbjct 349 IV+KSIAKS+GGKV SYRGLVQVNKGAGHSR+SVKCDALLVDTISRSDTYPVDIRDDVTM 408

Query 806 GEAATVSKVSENQLYIIMSRGLAEDEAMAMVVRGFVPIAKELPMYEAYALELNRLIELQME 865
      HEA VSK+S E QLYIIMSRGL+E EA M+V G+VER KELPMYEYA+ELNRLI +ME
Sbjct 409 EHEAKVSKISEEQLYIIMSRGLSESEATEMIVMGEVEPTKELPMYEAYAELNRLISYEME 468

Query 866 QAVG 869
      Q+VG
Sbjct 469 QSVG 472
  
```

**Match**

**Similarity**

**Dissimilarity**

**Gaps**

that is, the substitution of amino acids whose side chains have similar biochemical properties



# MSA input: multiple sequence Fasta file

>gi|4504351|ref|NP\_000510.1| delta globin [Homo sapiens]

MVHLTPEEKAVNALWGKVNVDAVGGEALGRLLVYPWTQRFFESFGDLSSPDAVMGNPKVKAHGKKVLG  
AFSDGLAHLDNLKGTFSQLSELHCDKLHVDPENFRLGNVLVCVLARNFGKEFTPQMQAAYQKVVAGVAN  
ALAHKYH

>gi|4504349|ref|NP\_000509.1| beta globin [Homo sapiens]

MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGKKVLG  
AFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVAN  
ALAHKYH

>gi|4885393|ref|NP\_005321.1| epsilon globin [Homo sapiens]

MVHFTAAEKAATSLWSKMNVEEAGGEALGRLLVYPWTQRFFDSFGNLSSPSAILGNPKVKAHGKKVLT  
SFGDAIKNMDNLKPAFAKLSELHCDKLHVDPENFKLLGNVMVIILATHFGKEFTPEVQAAWQKLVSAAV  
ALAHKYH

>gi|6715607|ref|NP\_000175.1| G-gamma globin [Homo sapiens]

MGHFTEEDKATITSLWGKVNVEDAGGETLGRLLVYPWTQRFFDSFGNLSSASAIMGNPKVKAHGKKVLT  
SLGDAIKHDDLKGTFQLSELHCDKLHVDPENFKLLGNVLTVLAIHFGKEFTPEVQASWQKMVTGVAS  
ALSSRYH

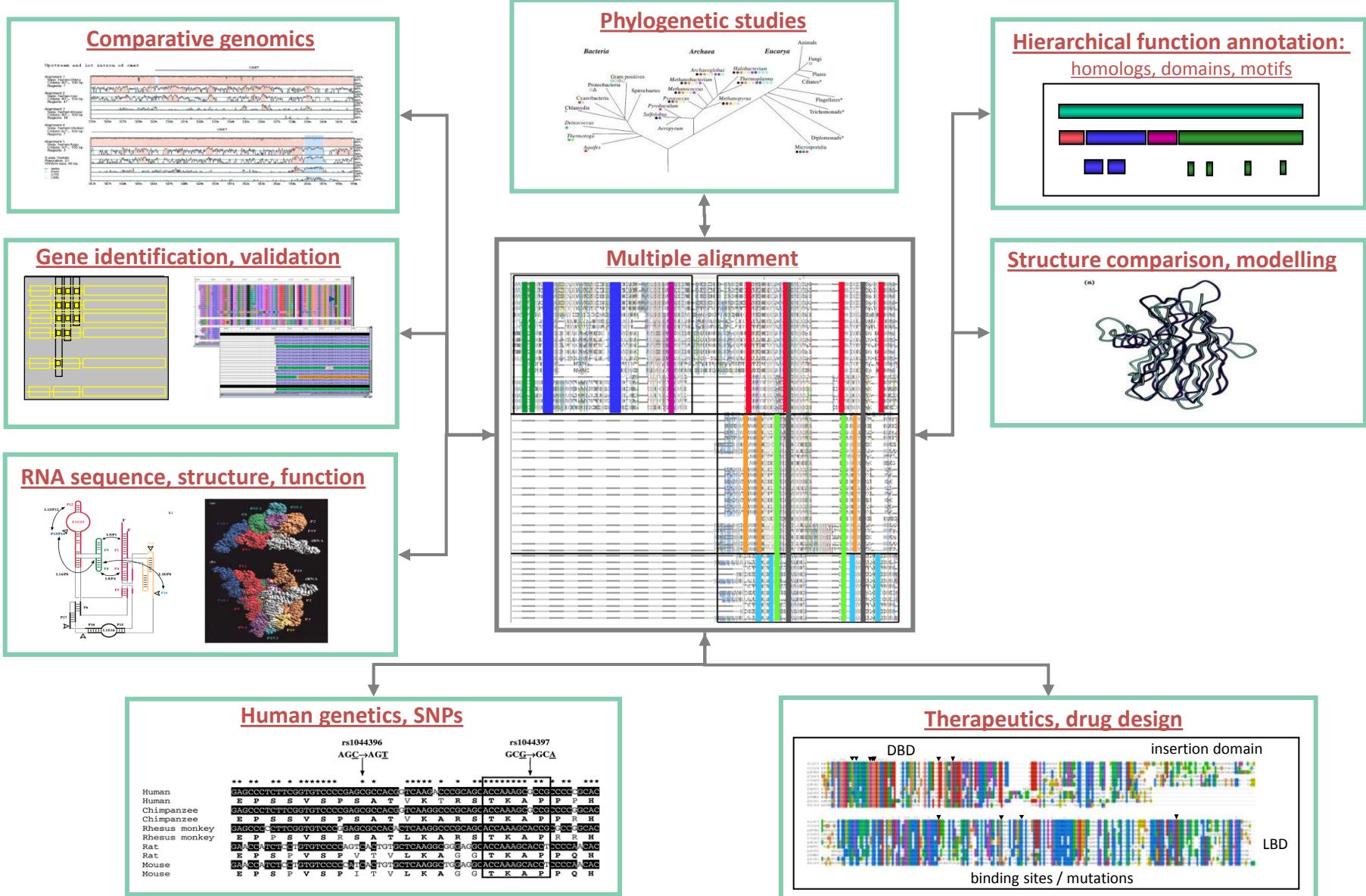
>gi|28302131|ref|NP\_000550.2| A-gamma globin [Homo sapiens]

MGHFTEEDKATITSLWGKVNVEDAGGETLGRLLVYPWTQRFFDSFGNLSSASAIMGNPKVKAHGKKVLT  
SLGDATKHDDLKGTFQLSELHCDKLHVDPENFKLLGNVLTVLAIHFGKEFTPEVQASWQKMVTAVAS  
ALSSRYH

>gi|4885397|ref|NP\_005323.1| hemoglobin, zeta [Homo sapiens]

MSLTKTERTIIVSMWAKISTQADTIGTETLERLFLSHPQTCKTYFPHFDLHPGSAQLRAHGSKVVAAGDA  
VKSIDDIGGALSKLSELHAYILRVDPVNFKLLSHCLLVTLAARFPADFTAEEHAAWDKFLSVVSSVLTEK  
YR

# Central role of multiple alignments



# Most important sequence databases

- **Genbank** – maintained by USA National Center for Biology Information (NCBI)
  - All biological sequences
    - [www.ncbi.nlm.nih.gov/Genbank/GenbankOverview.html](http://www.ncbi.nlm.nih.gov/Genbank/GenbankOverview.html)
  - Genomes
    - [www.ncbi.nlm.nih.gov:80/entrez/query.fcgi?db=Genome](http://www.ncbi.nlm.nih.gov:80/entrez/query.fcgi?db=Genome)
- **Swiss-Prot** - maintained by EMBL- European Bioinformatics Institute (EBI )
  - Protein sequences
    - [www.ebi.ac.uk/swissprot/](http://www.ebi.ac.uk/swissprot/)



# Sequence Similarity Searching

