

Bioinformatics

Bi1X-2010

Part I:
Public databases
Arbel Tadmor

Overview

- Obtaining sequence data from the internet
- Aligning two sequences
- Using the biologist's google: BLAST

We'll use hemoglobin as a case study

Related reading: Ch. 6 of Stryer (Biochemistry, 7th edition)

First stop: wiki

Hemoglobin
From Wikipedia, the free encyclopedia

Hemoglobin (also spelled **haemoglobin** and abbreviated **Hb** or **Hgt**) is the non-containing oxygen-transport metalloprotein in the red blood cells of vertebrates,^[1] and the tissues of some invertebrates.

In mammals, the protein makes up about 97% of the red blood cells dry content, and around 35% of the total content (including water^[2] or m³). Hemoglobin transports oxygen from the lungs or gills to the rest of the body (i.e. the tissues), where it releases the oxygen for cell use.

Hemoglobin has an oxygen binding capacity between 1.36 and 1.37 ml O₂ per gram of hemoglobin,^[3] which increases the total blood oxygen capacity sevenfold.^[4]

Hemoglobin is also found in avian red blood cells and their progenitor lines. Other cells that contain hemoglobin include the 40 adipogenic neurons in the substantia nigra, macrophages, muscle cells, and osteoeryth cells in the kidney. In these tissues, hemoglobin has a non-oxygen-carrying function as an antioxidant and a regulator of iron metabolism.^[5]

Structure

Protein type	metalloprotein, globulin
Function	oxygen transport
Category	heme (4)
Subunit	α, β
Subunit	α, β
Gene	α, β
Chromosomal location	16p11.2-13.1
α	16p11.2-13.1
β	11p15.5-16

Structure of human hemoglobin

- In adults hemoglobin is a tetramer: α₂β₂
- Each subunit contains a non-protein heme group (that holds an iron)
- The iron binds to an oxygen (shifting absorbance from blue to red)
- Multiple subunits give rise to cooperatively in oxygen binding and unbinding allowing this protein to release more oxygen in the tissues (making it a good transporter).

β subunit (146 aa)

α subunit (141 aa)

Iron-containing heme group

Further reading Ch. 7 of Stryer

Let's look up this gene

- Go to the PubMed website: *google pubmed* (<http://www.ncbi.nlm.nih.gov/pubmed/>)
- Search for **HBA1** under **gene** category

Using PubMed
PubMed Quick Start
New and Improved
PubMed Tutorials
Full Text Articles
PubMed FAQs

PubMed Tools
Single Citation Matcher
Batch Citation Matcher
Clinical Queries
Topic-Specific Queries

More Resources
Medline Database
PubMed Database
Clinical Trials
E-Books
LNCMI

Let's focus on the α1 subunit of hemoglobin. This gene is called **HBA1**

On which chromosome is this gene?

We'll start with human hemoglobin

human →

cattle →

frog →

etc.

More about the RefSeq database here... <http://www.ncbi.nlm.nih.gov/bookshelf/br.fc?book=handbook&part=ch18>
<https://www.ncbi.nlm.nih.gov/termsets/termsets/PMC130079/mf604975.pdf>

Which sequence should we pick? Genomic? mRNA? Protein?

Genomic view of HBA1

Genomic context of HBA1

click

Where HBA1 falls on chromosome 16...

chromosome 16

our gene

transcriptome

Genomic context

click

What is the length of ...

- The genomic sequence (NC...)?
- The mRNA sequence (NM...)?
- The protein sequence (NP...)?

Hint: Double click on gene then Right click → properties...

double click

RefSeq codes:
NC_123456 Genomic: Mixed Complete genomic molecules including genomes, chromosomes, organelles, plasmids.
NM_123456 mRNA: Mixed Transcript products; mature messenger RNA (mRNA) transcripts. NP_123456
NP_123456789 Protein: Mixed Protein products; primarily full-length precursor products but may include some partial proteins and mature peptide products.

Genomic vs. mRNA vs. protein sequences

<p>HBA1 total range: NC_000016.9 (226,679..227,520) total length: 842 processed length: 842 gene: 226679..227520 gene="HBA1" notes: Derived by automated computational analysis using gene prediction method: BestRefSeq. db_xref="GeneID:30397" db_xref="NC:4827" db_xref="PMID:1419007"</p> <p>Genomic</p>	<p>NM_000568.3: mRNA-hemoglobin, alpha 1 total range: NC_000016.9 (226,679..227,520) total length: 842 processed length: 571 mRNA product length: 576 join(C226719..226810,226928..227132,227282..227520) /gene="HBA1" /product="hemoglobin, alpha 1" notes: Derived by automated computational analysis using gene prediction method: BestRefSeq. transcript="NM_000568.3" db_xref="GeneID:30397" db_xref="GI:144561114"</p> <p>mRNA</p>
<p>NP_000549.1: alpha 1 globin total range: NC_000016.9 (226,716..227,410) total length: 695 processed length: 429 protein product length: 145 join(226716..226810,226928..227132,227282..227410) /gene="HBA1" notes: Derived by automated computational analysis using gene prediction method: BestRefSeq. accession="NP_000549.1" /product="alpha 1 globin" /protein_id="NP_000549.1" db_xref="NC:5103050399.1" db_xref="GeneID:30397" db_xref="GI:144561114"</p> <p>protein</p>	

Genomic vs. mRNA vs. protein sequences

<p>Genomic</p> <p>HBAT total range: NC_000016.9 (226,679..227,520) total length: 842</p> <p>gene: Z26679..Z27320 (gene="HBA1") note="Derived by automated computational analysis using gene prediction method: BestRefSeq." db_xref="GeneID:3039" db_xref="NCIC:4827" db_xref="MM:141807"</p>	<p>mRNA</p> <p>NM_000558.3: mRNA-hemoglobin, alpha 1 total range: NC_000016.9 (226,679..227,520) total length: 842 processed length: 576 mRNA product length: 576</p> <p>mRNA join(Z26679..Z26810.Z26928..Z27132.Z27282..Z27320) (gene="HBA1") product="hemoglobin, alpha 1" note="Derived by automated computational analysis using gene prediction method: BestRefSeq." transcript_id="NM_000558.3" db_xref="GeneID:3039" db_xref="GI:14456711"</p>
<p>Question 1 Why is the protein <u>product</u> length 142x3 = 426bp shorter than the protein <u>processed</u> length = 429bp?</p>	<p>protein</p> <p>NP_000549.1: alpha 1 globin total range: NC_000016.9 (226,716..227,410) processed length: 429 protein product length: 142</p> <p>CDS join(Z26716..Z26810.Z26928..Z27132.Z27282..Z27410) (gene="HBA1") note="Derived by automated computational analysis using gene prediction method: BestRefSeq." codon_start=1 product="alpha 1 globin" protein_id="NP_000549.1" db_xref="CCDS:CCDS10399.1" db_xref="GeneID:3039" db_xref="GI:4504347"</p>

Genomic vs. mRNA vs. protein sequences

<p>Genomic</p> <p>HBAT total range: NC_000016.9 (226,679..227,520) total length: 842</p> <p>gene: Z26679..Z27320 (gene="HBA1") note="Derived by automated computational analysis using gene prediction method: BestRefSeq." db_xref="GeneID:3039" db_xref="NCIC:4827" db_xref="MM:141807"</p>	<p>mRNA</p> <p>NM_000558.3: mRNA-hemoglobin, alpha 1 total range: NC_000016.9 (226,679..227,520) total length: 842 processed length: 576 mRNA product length: 576</p> <p>mRNA join(Z26679..Z26810.Z26928..Z27132.Z27282..Z27320) (gene="HBA1") product="hemoglobin, alpha 1" note="Derived by automated computational analysis using gene prediction method: BestRefSeq." transcript_id="NM_000558.3" db_xref="GeneID:3039" db_xref="GI:14456711"</p>
<p>Question 1 Why is the protein <u>product</u> length 142x3 = 426bp shorter than the protein <u>processed</u> length = 429bp?</p> <p>The stop codon was removed</p>	<p>protein</p> <p>NP_000549.1: alpha 1 globin total range: NC_000016.9 (226,716..227,410) processed length: 429 protein product length: 142</p> <p>CDS join(Z26716..Z26810.Z26928..Z27132.Z27282..Z27410) (gene="HBA1") note="Derived by automated computational analysis using gene prediction method: BestRefSeq." codon_start=1 product="alpha 1 globin" protein_id="NP_000549.1" db_xref="CCDS:CCDS10399.1" db_xref="GeneID:3039" db_xref="GI:4504347"</p>

Genomic vs. mRNA vs. protein sequences

<p>Genomic</p> <p>HBAT total range: NC_000016.9 (226,679..227,520) total length: 842</p> <p>gene: Z26679..Z27320 (gene="HBA1") note="Derived by automated computational analysis using gene prediction method: BestRefSeq." db_xref="GeneID:3039" db_xref="NCIC:4827" db_xref="MM:141807"</p>	<p>mRNA</p> <p>NM_000558.3: mRNA-hemoglobin, alpha 1 total range: NC_000016.9 (226,679..227,520) total length: 842 processed length: 576 mRNA product length: 576</p> <p>mRNA join(Z26679..Z26810.Z26928..Z27132.Z27282..Z27320) (gene="HBA1") product="hemoglobin, alpha 1" note="Derived by automated computational analysis using gene prediction method: BestRefSeq." transcript_id="NM_000558.3" db_xref="GeneID:3039" db_xref="GI:14456711"</p>
<p>Question 2 Why is the mRNA length after splicing (576bp) longer than the protein processed length (429bp)?</p>	<p>protein</p> <p>NP_000549.1: alpha 1 globin total range: NC_000016.9 (226,716..227,410) processed length: 429 protein product length: 142</p> <p>CDS join(Z26716..Z26810.Z26928..Z27132.Z27282..Z27410) (gene="HBA1") note="Derived by automated computational analysis using gene prediction method: BestRefSeq." codon_start=1 product="alpha 1 globin" protein_id="NP_000549.1" db_xref="CCDS:CCDS10399.1" db_xref="GeneID:3039" db_xref="GI:4504347"</p>

Genomic vs. mRNA vs. protein sequences

<p>Genomic</p> <p>HBAT total range: NC_000016.9 (226,679..227,520) total length: 842</p> <p>gene: Z26679..Z27320 (gene="HBA1") note="Derived by automated computational analysis using gene prediction method: BestRefSeq." db_xref="GeneID:3039" db_xref="NCIC:4827" db_xref="MM:141807"</p>	<p>mRNA</p> <p>NM_000558.3: mRNA-hemoglobin, alpha 1 total range: NC_000016.9 (226,679..227,520) total length: 842 processed length: 576 mRNA product length: 576</p> <p>mRNA join(Z26679..Z26810.Z26928..Z27132.Z27282..Z27320) (gene="HBA1") product="hemoglobin, alpha 1" note="Derived by automated computational analysis using gene prediction method: BestRefSeq." transcript_id="NM_000558.3" db_xref="GeneID:3039" db_xref="GI:14456711"</p>
<p>Question 2 Why is the mRNA length after splicing (576bp) longer than the protein processed length (429bp)?</p> <p>The protein sequence is the sequence between the start and stop codons. The mRNA includes an additional untranslated 5' region and an untranslated 3' regions</p>	<p>protein</p> <p>NP_000549.1: alpha 1 globin total range: NC_000016.9 (226,716..227,410) processed length: 429 protein product length: 142</p> <p>CDS join(Z26716..Z26810.Z26928..Z27132.Z27282..Z27410) (gene="HBA1") note="Derived by automated computational analysis using gene prediction method: BestRefSeq." codon_start=1 product="alpha 1 globin" protein_id="NP_000549.1" db_xref="CCDS:CCDS10399.1" db_xref="GeneID:3039" db_xref="GI:4504347"</p>

Let's look at the sequence

Homo sapiens chromosome 16, GRCh37 primary reference assembly

click

Let's look at the sequence

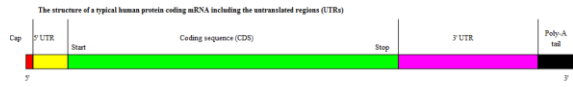
Homo sapiens chromosome 16, GRCh37 primary reference assembly

Untranslated 5' region of the mRNA (blue) ATG Start codon (red) Exon (red) 5' splice site (GU) (green) Intron (green) 3' splice site (AG) (green)

UAA stop codon Untranslated 3' region of the mRNA (blue)

Total length mRNA = blue + red = 576 bp
Total length protein = red = 429 bp

Typical Eukaryotic mRNA



We would like to find similar proteins in nature

Which sequence should we use for the search?

- Genomic?
- mRNA ?
- Protein?

OK let's grab the aa sequence of the protein

and the nt sequence of the protein

Nucleotide versus amino acid sequences

> NP_000549.1 Homo sapiens nt (429)
ATGGTGCTGTCTCTCCGCCACAAACCAAGCCGCTGGGGTAAGTGGGGCGCACCGCTGGGAGATGGTGGCGAGG
 CCTCGAGAGAGATGTTCTGCTCTCCACCACCAAGACTCTCCGGACTCTGACCTACCCACCGCGCTCCGCCAGCTTAAGGGCA
 CGGCAAGAGGTGGCGAGCCCTGACACCGCTGGCCAGTGGAGCAATGCCAGCGCTGTCGCCCTGAGGCACTGCAC
 GGGCAACAGCTTGGTGGACCGGTCACCTAGCTCTGCTGGTACCGTGGCCACCTCCCGCCACCTCCCGCGGATC
 ACCCTCGGGTGCAGCCTCCCTGGACAACTGCTGCTCTGTGAGCAAGCTGCTCACTCAATACCGT**TAA**

> NP_000549.1 Homo sapiens aa (142)
MVLSPADKTVKAAWGVGAAHGEYGALEAFMISPTTKYTFPHDLHSGAQVGHGKVKVADALTNVAVHDDMPNALSALSDHA
 HKLRVDFVFKLSHCLLVLAALPAEFTFAVWASDKFLASVSTLSKRY

		Second letter							
		G	A	C	T	G	A	C	T
First letter	U	UCU Phenylalanine	UUA Leucine	UCU Tyrosine	UGU Cysteine	UUA Leucine	UUA Leucine	UUA Leucine	UUA Leucine
	C	CCU Proline	CCA Proline	CCU Proline	CCU Proline	CCA Proline	CCA Proline	CCA Proline	CCA Proline
	A	AUU Methionine	AUA Methionine	AUU Methionine	AUU Methionine	AUA Methionine	AUA Methionine	AUA Methionine	AUA Methionine
	G	GUU Valine	GUA Valine	GUU Valine	GUU Valine	GUA Valine	GUA Valine	GUA Valine	GUA Valine

Nucleotide versus amino acid sequences

> NP_000549.1 Homo sapiens nt (429)
ATGGTGCTGTCTCTCCGCCACAAACCAAGCCGCTGGGGTAAGTGGGGCGCACCGCTGGGAGATGGTGGCGAGG
 CCTCGAGAGAGATGTTCTGCTCTCCACCACCAAGACTCTCCGGACTCTGACCTACCCACCGCGCTCCGCCAGCTTAAGGGCA
 CGGCAAGAGGTGGCGAGCCCTGACACCGCTGGCCAGTGGAGCAATGCCAGCGCTGTCGCCCTGAGGCACTGCAC
 GGGCAACAGCTTGGTGGACCGGTCACCTAGCTCTGCTGGTACCGTGGCCACCTCCCGCCACCTCCCGCGGATC
 ACCCTCGGGTGCAGCCTCCCTGGACAACTGCTGCTCTGTGAGCAAGCTGCTCACTCAATACCGT**TAA**

> NP_000549.1 Homo sapiens aa (142)
MVLSPADKTVKAAWGVGAAHGEYGALEAFMISPTTKYTFPHDLHSGAQVGHGKVKVADALTNVAVHDDMPNALSALSDHA
 HKLRVDFVFKLSHCLLVLAALPAEFTFAVWASDKFLASVSTLSKRY

- Which sequence should we use to search with, the amino acid sequence or the nucleotide sequence?

Nucleotide versus amino acid sequences

> NP_000549.1 Homo sapiens nt (429)

ATGGTGCTGTCTCCTGCCGACAAGACCAACGTCAAGGCCCTGGGGTAAGTTCGGCGCCACGCTGGGGAGTATGGTGGGAGG
 CCTGTGAGAGAGATGTTCTGTCTCCCTCCACCAAGACTACTTCCCGCACTTCACCTGAGCCACGGCTGTGCCAGGGTTAAGGGCCA
 CGGCAAGAGAGGTTGGCGCCCTGTGACAGCCCTGTGGCGAGCTGAGCAATGTCACACCGCTGTCTCCCTGGAGCGACTGCAAC
 GGCAACAAGCTTGGGTGACCCGGTCAACTTCAAGCTCCTAAGCACTGCTGCTGGTGAACCTGGCCGCCCACTCCCGCCGAGTTC
 ACCCTTGGGGTGCACGCTCCCTGGACAGTTCTGGCTCTGTGAGCACCGTCTGCACTTCAAATACCGT**TAA**

> NP_000549.1 Homo sapiens aa (142)

MVLSADKTNWAAWIKVGAHAGEYGALEEMFISFPTTKTYFPHFDLSHESAQVKGHGKIVADLTNAVAVHDDMPNLSALSOLDHA
 HKLRVDPNPKLLSHLLVLAHLPAETRAVHNSLDKFLASVSTVTSKYR

- Which sequence should we use to search with, the amino acid sequence or the nucleotide sequence?

It depends on your goal, but generally to find homologs, aa sequences is the way to go:

- Selection pressure on amino acid sequence is much stronger than on nt sequence
- Two random nt sequences share 25% (50% with gaps) identical characters whereas amino acids sequences share only 5% (10-15% with gaps), improving the signal to noise considerably <we'll see this later...>

Bioinformatics

Bi1X-2010

Part II: Sequence alignment and BLAST

Arbel Tadmor

Overview

- Sequence alignment - basic concepts
- Local vs. global alignment
- BLAST –a local alignment tool
- Exercise: alignment of random sequences
- Exercise: finding homologs of HBA1 with BLAST

Finding homologues sequences

- **Homology** ≠ similarity

Homology: two sequences are descended from a common ancestor, therefore in an alignment identical residues at a site are identical by descent

Similarity: merely reflects proportion of sites that are identical

- What changes could occur over time in a sequence?

Finding homologues sequences

- **Homology** ≠ similarity

Homology: two sequences are descended from a common ancestor, therefore in an alignment identical residues at a site are identical by descent

Similarity: merely reflects proportion of sites that are identical

- What changes could occur over time in a sequence?

- Changes that conserve length:
 - Substitutions (one nt mutation)
 - Inversions
- Changes that do not conserve length :
 - Deletions (e.g. DNApol replication error, unequal crossover, transposon)
 - Insertions (e.g. DNApol replication error, unequal crossover, transposon HGT via transposable elements)

Matches, Mismatches and Indels

Two aligned, identical characters in an alignment are a **match**.

Two aligned, unequal characters are a **mismatch**.

A character aligned with a **gap**, represents an **indel** (insertion/deletion).

```
A A C T A C T - C C T A A C A C T - -
- - C T C C T A C C T - - T A C T T T
```

the bars show mismatches

10 matches, 2 mismatches, 7 gaps
 Total number of characters 10+2+7=19
Percent identity = 10/19 = 53%

Scoring scheme for an alignment – example:

$w(\text{match}) = +2$ or $4X4$ nt substitution matrix

$w(\text{mismatch}) = -1$ or $4X4$ nt substitution matrix

$w(\text{gap}) = -3$ or some other heuristic penalty

Find alignment with maximum score

Substitution matrix for amino acids BLOSUM: Blocks Substitution Matrix (used to calculate an alignment score)

- Values in matrix are **empirical**. Based on a large sample of verified aa pairwise alignments
- i, j element = Log of probability (p_{ij}) that amino acid i mutates into amino acid j in a homologous sequence normalized by the probability of this match by chance given the frequency of the amino acids (q_i, q_j) -> **positive: better than chance, negative: worst than chance**

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
C	9																			
S	-1	4																		
T	-1	1	5																	
P	-3	-1	-1	7																
A	0	1	0	-1	4															
G	-3	0	-2	-2	0	6														
N	-3	1	0	-2	-2	0	6													
D	-3	0	-1	-1	-2	-1	1	6												
E	-4	0	-1	-1	-1	-2	0	1	6											
Q	-3	0	-1	-1	-1	-2	0	2	5	6										
H	-3	-1	-2	-2	-2	-2	1	-1	0	8	6									
R	-3	-1	-1	-2	-1	-2	0	2	1	5	5	6								
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5							
M	-1	-1	-2	-1	-3	-2	-2	0	-2	-1	-1	5								
I	-1	-2	-1	-3	-1	-4	-2	-3	-3	-3	1	4	5							
L	-1	-2	-1	-3	-1	-4	-2	-3	-2	-3	2	2	4	5						
V	-1	-2	-2	-2	0	-3	-3	-2	-3	-2	1	3	1	4	5					
F	-2	-2	-2	-4	-2	-3	-3	-3	-1	-3	0	0	4	-1	6					
Y	-2	-2	-3	-2	-3	-2	-3	-2	-1	-2	-2	-1	-1	-1	3	7				
W	-2	-3	-2	-4	-3	-2	-4	-3	-2	-2	-3	-3	-1	-3	-2	1	2	11		

$$S_{ij} = \left(\frac{1}{\lambda} \right) \log \left(\frac{p_{ij}}{q_i * q_j} \right)$$

Positive for chemically similar substitutions (more likely than chance)

Unlikely substitution compared with chance

Common aa have a low score

Rare aa have a high score

• We will come back to substitution matrices later when we talk about phylogeny

Local alignment vs. global alignment

- **Global alignment** – attempts to align every residue in every sequence
P-ELICAN--
COELICANTH
- **Local alignment** – find best subsequence alignment (useful for finding similar exons in two genomes, finding similar functional regions in a protein, etc.)
ELICAN
ELICAN

BLAST- Basic Local Alignment Search Tool

The Google of biology (google blast)

Basic BLAST

Choose a BLAST program to run

- nt vs. nt** nucleotide blast Search a nucleotide database using a nucleotide query
Algorithms: blastn, megablast, discontinuous megablast
- aa vs. aa** protein blast Search protein database using a protein query
Algorithms: blastp, psi-blast, phi-blast
- tras nt vs. aa database** blastx Search protein database using a translated nucleotide query
- aa vs. trans nt database** tblastn Search translated nucleotide database using a protein query
- tras nt vs. trans nt database** 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 (cdt)
- 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** (vecscreen)
- **Align two (or more) sequences** using BLAST (BL2seq)
- Search **proteins** or nucleotide targets in PubChem BioAssay
- Search **SRA transcript and genomic libraries**
- **Constraint Based Protein Multiple Alignment Tool**
- Needleman-Wunsch **Global Sequence Alignment Tool**

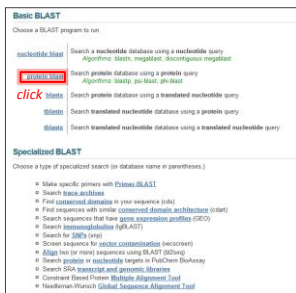
Global alignment

BLAST Scores and Statistics

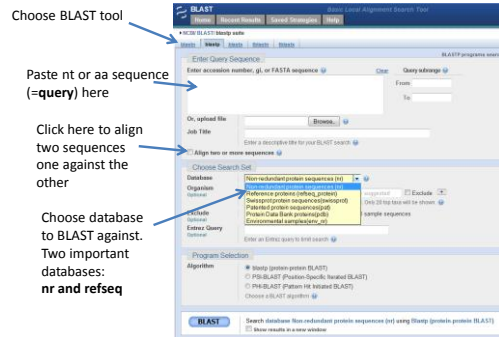
- **Percent identity** is the fraction of identical characters
- **Percent similarity** (for amino acids) is the fraction of identical or chemically similar residues
- **Bit score** – A normalized alignment score (S'). The bit score gives an indication of how good the alignment is; **the higher the score, the better the alignment.**
- **E value** – The E-value is a test statistic that gives an indication of the statistical significance of a given pairwise alignment by comparing it to a model of random sequences. It's the average number of sequences with this level of similarity (i.e. raw alignment score S) or better expected to be in the database by chance. Reflects the size of the database and the scoring system. **Lower is better.** The threshold is usually placed at 10^{-3} .
- **P value** = The probability of finding at least one such sequence in the database by chance = $1 - e^{-E}$ (The E value is just the λ parameter in a Poisson distribution $\lambda e^{-\lambda}/n!$...)

Further reading: <http://www.ncbi.nlm.nih.gov/blast/tutorial/>

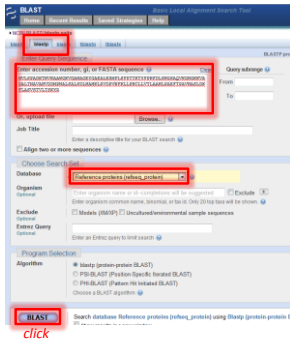
Main BLAST window



Main BLAST window



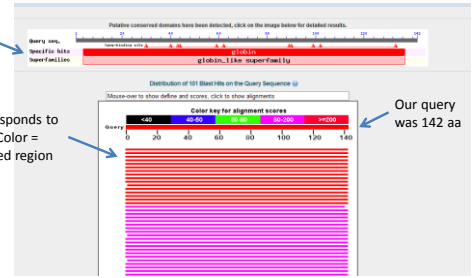
Let's search now for homologs of human HBA1 in the refseq_protein database



Results

Conserved domains in the protein

Each line corresponds to a positive hit. Color = score for aligned region



Example of local alignment result

Link to gene record

Bit score E value

Our seq.

Database seq.

mismatch

+ indicates similar aa

106 identities + 14 "+" = 120 positive hits

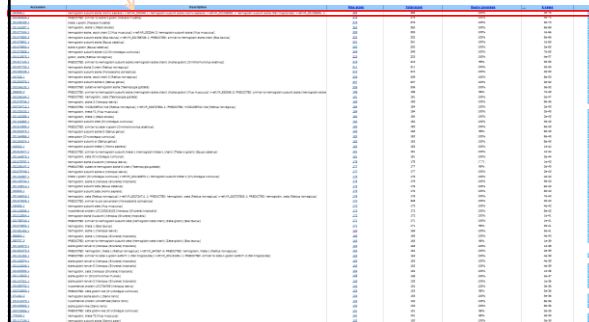
```

>ref|NP_001039756.1|UC| globin, alpha [Rattus norvegicus]
Length=142
GENE ID: 287167 LOC287247 | globin, alpha [Rattus norvegicus]
(10 or fewer PubMed links)
Score = 223 bits (568), Expect = 4e-57, Method: Compositional matrix adjust.
Identities = 106/142 (74%), Positives = 120/142 (84%), Gaps = 0/142 (0%)

Query 1  MVLSPADKINVKAAWGKVCAGAGEYVAEALERMIFPPITKITYFFPHDLSHGSAQVKGHG 60
Sbjct 1  MVLSDRNRKAMKQVIGNVAALVARTIQRIIFPSTKIPFPHFSSGQVKAHG 60
Query 61  KKVDALINAVAVDMDINALSLSLHAKLAVDVPVNFYLSHCLLVITLAHLPDEFIP 120
Sbjct 61  KKVDALINA +H+DD+P ALS LSDLHAKLAVDVPVNFK LSHCLLVITLA+H P +FIP 120
Query 121 AVHASLDKFLSVSTVLSKYR 142
Sbjct 121 A+HASLDKFLSVSTVLSKYR 142
    
```

Results

Here's our sequence NP_000549.1



Results

Here's our sequence NP_000549.1

Human NP_000549.1

Orive baboon NP_001162287.1

Chicken NP_001004376.1

African clawed frog NP_001004376.1

Atlantic salmon NP_001177354.1

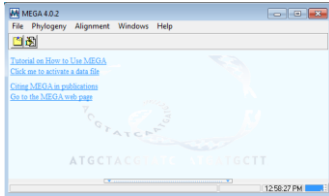
Zebra fish NP_571232.2

Before we align...

- Should we align the amino acid sequences or the nucleotide sequences of a protein coding gene?

Multiple Alignment in MEGA4

- Google MEGA4 and install
- Launch MEGA4
- Drag the FASTA file into MEGA4



Multiple Alignment in MEGA4

nt view



aa view



ClustalW: Popular multiple alignment algorithm

Algorithm overview:

- Global alignment on all sequence pairs to find the **distance** between all pairs of sequences
- Uses distances to create a **guide tree**
- Align the closest sequences in the guide tree, followed by adding more sequences to the initial alignment

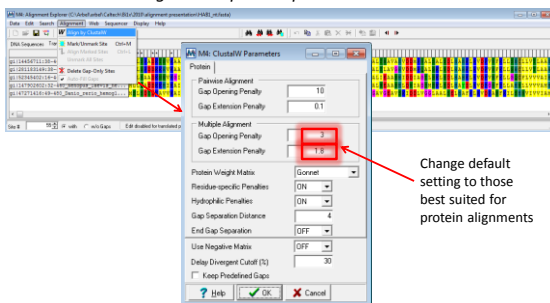
Read more at <http://www.ebi.ac.uk/2can/tutorials/nucleotide/clustalw.html>

What can we do with a multiple alignment?

- Identify conserved regions within protein
 - Signifies conserved function
 - Useful for primer design
- Identify variable regions within protein
 - Functionally not important
 - Important but under positive selection pressure or rapidly changing
 - Identify non-silent mutations (e.g. leading to disease, due to adaptation, etc.)
- Construct a phylogenetic tree (discuss later)

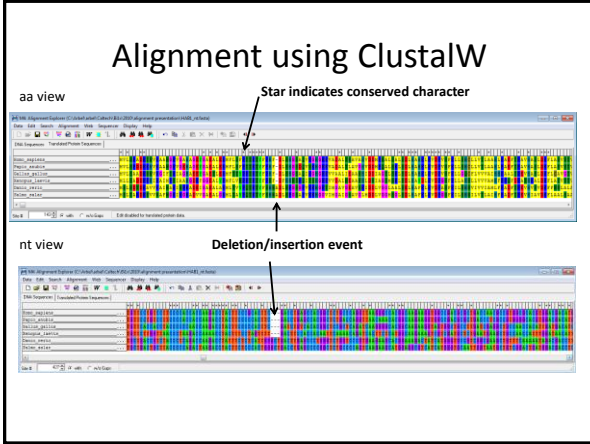
Alignment using ClustalW

Hint: be sure to align in the protein pane



Manual inspection of alignment

- Pay attention to the edges
- Are there any obviously wrong sequences that did not align well?
 - Sequences too divergent? (must have $\geq 20\%$ aa identity)
 - Reverse complement?
 - Frame shift?

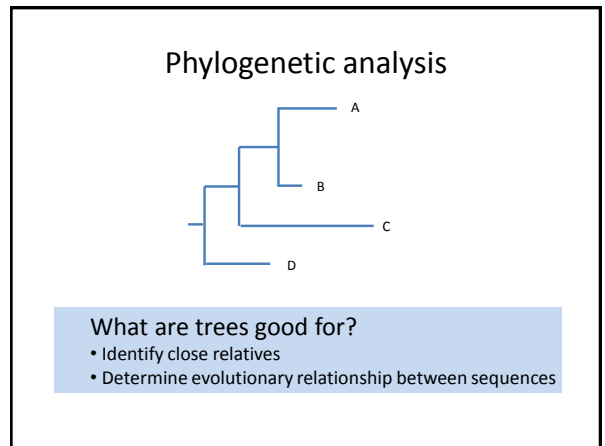
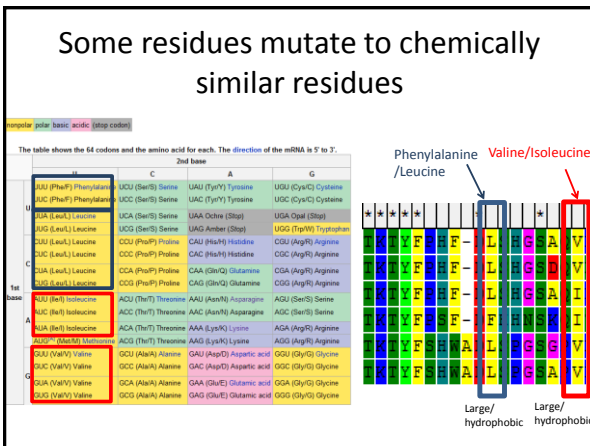
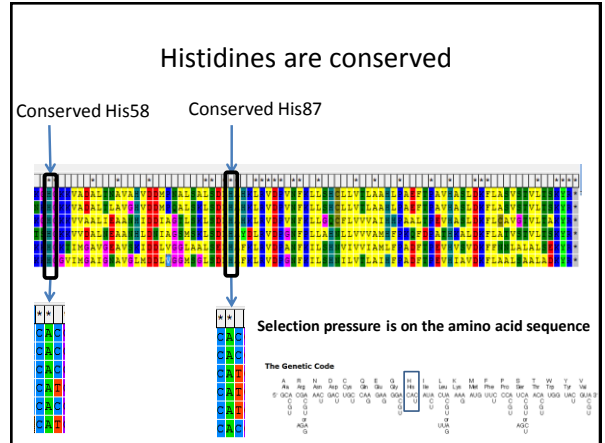
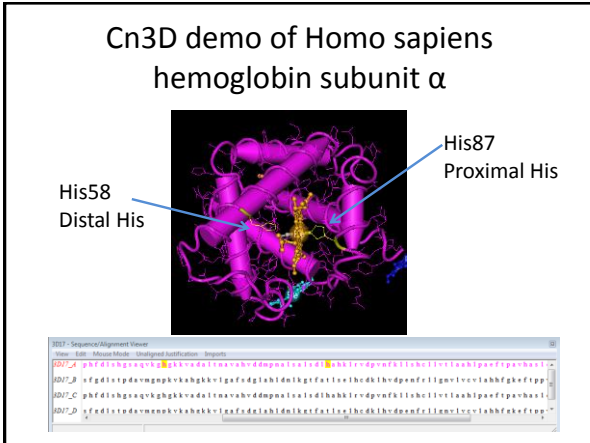


Example of residues conserved due to function: His87 and His58 maintain the heme and oxidation state of the iron

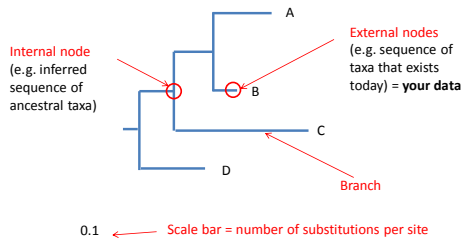
Proximal His: Anchoring of the heme is facilitated by a nitrogen from ϵ histidine that binds to the iron.

Distal His: The bound oxygen can be in two states, dioxygen (bound to Fe^{2+}) and superoxide (bound to Fe^{3+}). Oxygen must be released in the former because the later is both harmful and leaves the iron in a state that cannot bind oxygen. The distal histidine binds more strongly to superoxide and the oxygen is therefore less likely to be released.

Mathews et al. 2000. Biochemistry 3 rd edition

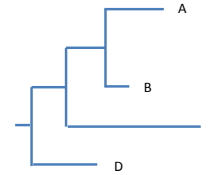


Some tree terminology



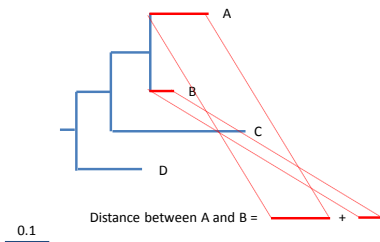
We will discuss only **bifurcating** trees: each node has two immediate descendant lineages, i.e. we assume evolutionary speciation is a binary process

Trees have two elements: branch lengths and branching order



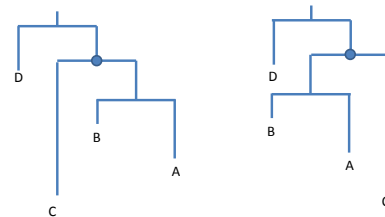
- **Topology** = branching pattern
- To estimate a tree you need to estimate
 - Branch lengths (simple problem)
 - Branching order (difficult for many seq)

Branch lengths



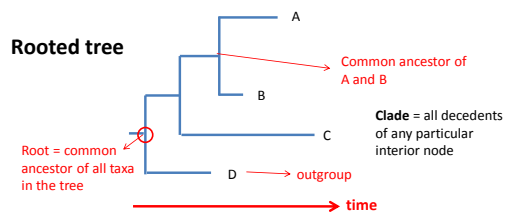
- **Evolutionary distance** = accumulated horizontal distance between two external nodes = estimated number of substitution per site that differ between the two sequences = **d**

Branch order: note that trees are like mobiles...



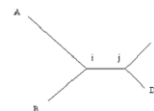
These two trees are equivalent in every way

Rooted trees



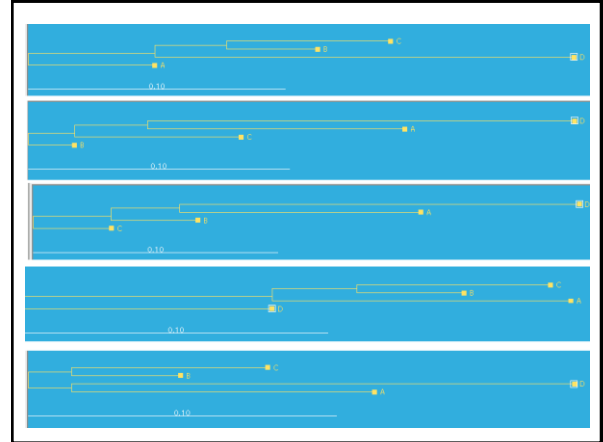
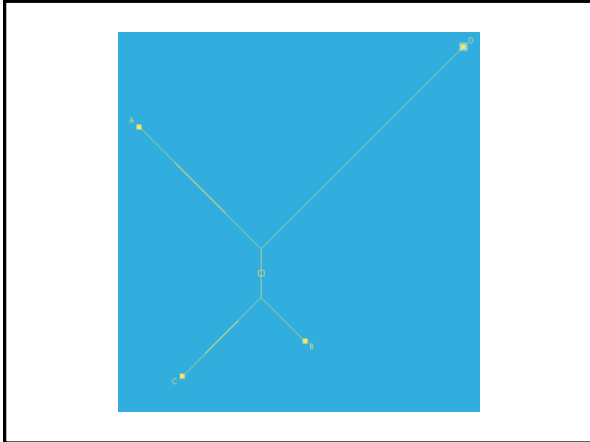
- If a tree is rooted, root = left most internal node
- By selecting a node to be a **root** you set a time arrow
- **The more recently species share a common ancestor, the more closely related they are** (e.g. A and C are more closely related than A and D)
- Which node is the root? You "break the symmetry" by adding additional information: e.g. you know D is more distantly related to the ingroup sequences than the ingroup sequences are related to each other

Unrooted trees



5 ways to root this tree

- If it is not explicitly said that the tree is rooted assume it is unrooted
- unrooted trees do not specify an evolutionary pathway (who descended from whom) only relationships among taxa



How to construct a tree?

Algorithmic methods (distance based)	Tree-searching methods (character based)
Use algorithm to construct a single tree from the data	Construct many trees then use some criterion to decide which is the best tree
<ol style="list-style-type: none"> Multiple alignment Calculate distance matrix = matrix of evolutionary distances between all pairs of aligned sequences Calculate tree topology 	<ol style="list-style-type: none"> Multiple alignment Compare characters at each column in the alignment and give each topology a score Choose the topology with the best score
Example: Neighbor Joining (many others)	Examples: <ul style="list-style-type: none"> Maximum Likelihood methods Maximum Parsimony methods Bayesian methods

Measuring evolutionary distance between two sequences

- The evolutionary distance between two sequences d is the total number of number of aa/nt substitutions per site between the two sequences
- $d=2rt$
t=time in years, *r*= substitution rate per year per site
- Branch length in tree = d
- How can we estimate d from the sequences?

Measuring evolutionary distance between two sequences

- p distance:** $p=n_d/n$ = number of different aa/nt between two aligned sequences of length n
 - Doesn't account for multiple hits: A → C → T
 - Doesn't account for back mutations: A → C → A
 - Doesn't account for parallel mutations: A → C; A → C
 - Underestimates d
 - Saturates at $p=0.75$ (not a good estimate of d when p is high)
 - Can result in wrong topology
- Estimation of d based on a stochastic model:** aa/nt substitutions are modeled as a stochastic process.
 - Different stochastic models make different assumptions regarding the probability of aa/nt substitutions
 - Different models assume different **substitution matrices**

Calculating distances in MEGA4: nt p distance

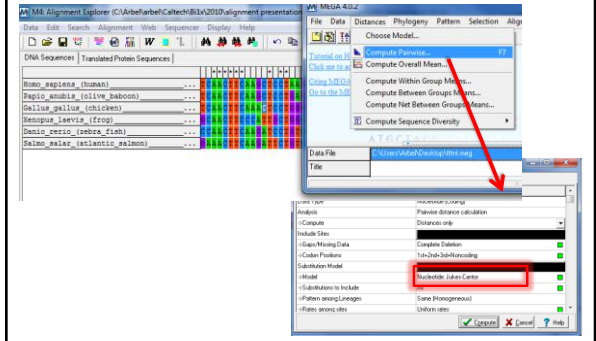
Calculating distances in MEGA4: p distance

p distance = 100 - percent identity/100

	1	2	3	4	5	6
1. Homo sapiens (human) ...						
2. Papio anubis (olive baboon) ...	0.0513					
3. Gallus gallus (chicken) ...	0.2587	0.2681				
4. Xenopus laevis (frog) ...	0.4056	0.3869	0.3660			
5. Danio rerio (zebra fish) ...	0.3846	0.4009	0.3776	0.4336		
6. Salmo salar (atlantic salmon) ...	0.3963	0.3963	0.3986	0.4522	0.2960	

Note that p distances < 0.75

Calculating distances in MEGA4: JC correction



Calculating distances in MEGA4: Evolutionary distance

JC method

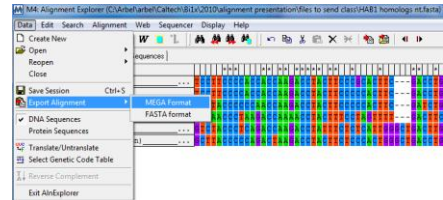
	1	2	3	4	5	6
1. Homo sapiens (human) ...						
2. Papio anubis (olive baboon) ...	0.0531					
3. Gallus gallus (chicken) ...	0.3173	0.3317				
4. Xenopus laevis (frog) ...	0.5837	0.5441	0.5020			
5. Danio rerio (zebra fish) ...	0.5393	0.5736	0.5251	0.6472		
6. Salmo salar (atlantic salmon) ...	0.5637	0.5637	0.5686	0.6928	0.3765	

p distance

	1	2	3	4	5	6
1. Homo sapiens (human) ...						
2. Papio anubis (olive baboon) ...	0.0513					
3. Gallus gallus (chicken) ...	0.2587	0.2681				
4. Xenopus laevis (frog) ...	0.4056	0.3869	0.3660			
5. Danio rerio (zebra fish) ...	0.3846	0.4009	0.3776	0.4336		
6. Salmo salar (atlantic salmon) ...	0.3963	0.3963	0.3986	0.4522	0.2960	

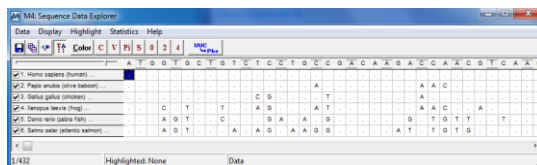
Building a Neighbor joining tree

Export nt alignment in MEGA format



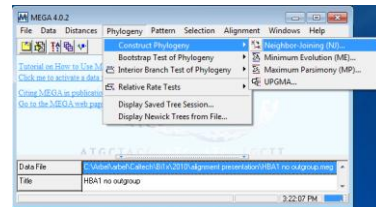
Building a Neighbor joining tree

Close MEGA4 and open your MEGA file.
You should get this window with a nt sequence:



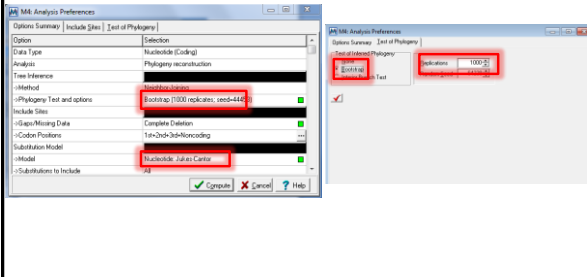
Building a Neighbor joining tree

Now let's build a NJ tree

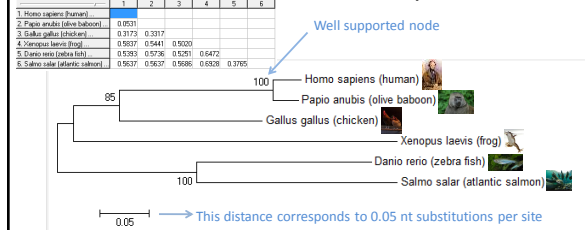


Building a Neighbor joining tree

- Use JC nt model
- Calculate bootstrap support with 1000 replications



Building a Neighbor joining tree (unrooted nt tree)



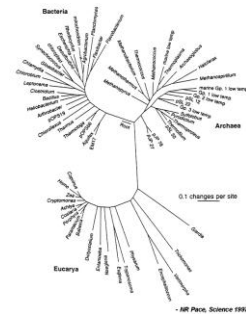
- Human and primates group together
- Fish group together

Bioinformatics Bi1X-2010

Part IV: Phylogenetic analysis of rRNA sequences - an exercise in class

Arbel Tadmor

Universal phylogenetic tree based on small-subunit (SSU) rRNA sequences



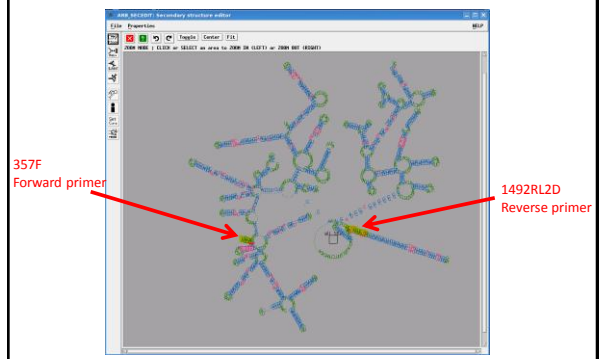
How are SSU rRNA sequences different?

- rRNA genes are universal genes that are highly conserved
- Used as phylogenetic markers for species

Some technical points

- No aa sequence
- Selection is directly on the nt sequence
- Alignment should take into account secondary structure of rRNA molecule
- We will therefore use a dedicated website called **green genes** to align the sequences and analyze the alignment in MEGA

rRNA secondary structure



Program for today

- Learn more about the rRNA gene and the nature of the amplicon you generated
- Read the **green genes** website tutorial
- Convert your traces to nt using **Sequence Scanner**
- Align sequences with green genes
- Check for chimeras using green genes
- Import alignment into MEGA
- Calculate distance matrix
- Build a NJ tree
- Identify closest relatives of your sequences
- Is it likely that you found your phylotypes in the pond?

Additional slides

Example 1: Global alignment of two random 300bp nucleotide sequences

We will generate random sequences of 300 nt in Matlab:

```
>> rand_int = floor(4*rand([1, 300])+1);
>> rand_nt = int2nt(rand_int)

rand_nt =
TCGAAGGGCGCTCGGTAGAGTACGTGTCCCAACTGTTGCCTAAGCGGCC
GTACAGTAGGGCGAGGCACGCTACTGTTACGAGATTCCACCGAAGAA
AAGTTAAGCCCTCGAAAGGTAACCATCGAGGCCGTGATCTGGCATG
AAATACTACGGGCTTCCCCAACATAAGGCAACTCATCGGGGATC
ACATGGCCCTCGGTCCGATGATTTGCCGATTTTCACGGTTGCCTCA
TCAAGCCCGCCAACGGGTTAGTGGAAACGAATATGAGGCAGACTCTCAC
ATCGCTATCTGT
```

Example 1: Global alignment of two random 300bp nucleotide sequences

```
>Random seq 1
TCGAAGGGCGCTCGGTAGAGTACGTGTCCCAACTGTTGCCTAAGCGGCCGTACGAGTGGCGAGGCAGCGTACTG
TTACGAGATTCCACCGAAGAAAGTTAAGCCCTCGAAAGGTAACCATCGAGGCCGTGATCTGGCATGAAATA
CTACGGGCTTCCCCAACATAAGGCAACTCATCGGGGATCAGCATCGCCCTCGGTCCGATGATTCGGCATT
TTACGGTTGCCATCAAGCCCGCCAACGGGTTAGTGGAAAGCAATATGAGGCAGACTCTCACATCGCTACTG
```

```
>Random seq 2
CCGTAAGCGTCTACAGAGGGCGTGGACCAACTGTCTATCGTCACTGTGTAGGTGACCAAGCTG
AGAGCTCTATCTTCTTAATAGTACCTTAAAGTGTGTTGACCTACAAGCAATCGTCGCTC
CTCAACTTCTCGCTCATCACTAGAAATGTGTAAGCCGACAGCCGCAACGGTCAATGGCGTTCAGG
ATTACGACATTATCAAGGCCAATCGTCGGATGGCGAGGGCCCTTCAACGAAAGGGGGTGTGATA
```

Example 1: Global alignment of two random 300bp nucleotide sequences

```
>Random seq 1
TCGAAGGGCGCTCGGTAGAGTACGTGTCCCAACTGTTGCCTAAGCGGCCGTACGAGTGGCGAGGCAGCGTACTG
TTACGAGATTCCACCGAAGAAAGTTAAGCCCTCGAAAGGTAACCATCGAGGCCGTGATCTGGCATGAAATA
CTACGGGCTTCCCCAACATAAGGCAACTCATCGGGGATCAGCATCGCCCTCGGTCCGATGATTCGGCATT
TTACGGTTGCCATCAAGCCCGCCAACGGGTTAGTGGAAAGCAATATGAGGCAGACTCTCACATCGCTACTG
```

```
>Random seq 2
CGTAAAGTGTACAGAGGGCGTGGACCAACTGTCTATCGTCACTGTGTAGGTGACCAAGCTG
AGAGCTCTATCTTCTTAATAGTACCTTAAAGTGTGTTGACCTACAAGCAATCGTCGCTC
CTCAACTTCTCGCTCATCACTAGAAATGTGTAAGCCGACAGCCGCAACGGTCAATGGCGTTCAGG
ATTACGACATTATCAAGGCCAATCGTCGGATGGCGAGGGCCCTTCAACGAAAGGGGGTGTGATG
```

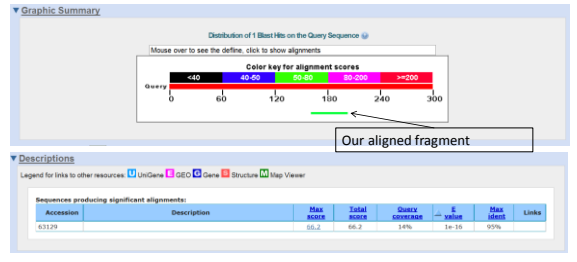
Example 2: Local alignment of a 300bp sequence and an internal fragment containing a single insertion and a single mismatch using BLAST

```
>Random seq 1
TCGAAGGGCGCTCGGTAGAGTACGTGTCCCAACTGTTGCCTAAGCGGCCGTACGAGTGGCGAGGCAGCGTACTG
TTACGAGATTCCACCGAAGAAAGTTAAGCCCTCGAAAGGTAACCATCGAGGCCGTGATCTGGCATGAAATA
CTACGGGCTTCCCCAACATAAGGCAACTCATCGGGGATCAGCATCGCCCTCGGTCCGATGATTCGGCATT
TTTACGGTTGCCATCAAGCCCGCCAACGGGTTAGTGGAAAGCAATATGAGGCAGACTCTCACATCGCTACTG
GT
```

```
>Random seq 2
CGTAAAGTGTACAGAGGGCGTGGACCAACTGTCTATCGTCACTGTGTAGGTGACCAAGCTG
AGAGCTCTATCTTCTTAATAGTACCTTAAAGTGTGTTGACCTACAAGCAATCGTCGCTC
CTCAACTTCTCGCTCATCACTAGAAATGTGTAAGCCGACAGCCGCAACGGTCAATGGCGTTCAGG
ATTACGACATTATCAAGGCCAATCGTCGGATGGCGAGGGCCCTTCAACGAAAGGGGGTGTGATA
```

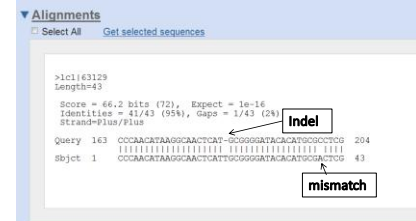
Example 2: Local alignment of a 300bp sequence and an internal fragment containing a single insertion and a single mismatch using BLAST

```
>Random seq 1
TCGAAGGCGCTCGGTAGAGTACGTGTCCCAACTGTTGCCTAAGCGCGGTACAGTAGGCGGAGGACGCTACTG
TTACGAGATTCTACCGAAGAAAGTTAAGCCCTCGAAAGGTAACCTCGAGCCGCTGATCGGCATGAATA
CTACGGCCCTCCCCCAACATAAGGCAACTCATCGGGGATACACATCGCACTCGGTCGGATATGATGCCGCA
TTTTACGGTTGCCTCAATCAAGCCCGCAACGGGTTAGTGAACAAGATGAGGCGACACTCACATCGCTATCT
GT
```



Example 2: Local alignment of a 300bp sequence and an internal fragment containing a single insertion and a single mismatch using BLAST

```
>Random seq 1
TCGAAGGCGCTCGGTAGAGTACGTGTCCCAACTGTTGCCTAAGCGCGGTACAGTAGGCGGAGGACGCTACTG
TTACGAGATTCTACCGAAGAAAGTTAAGCCCTCGAAAGGTAACCTCGAGCCGCTGATCGGCATGAATA
CTACGGCCCTCCCCCAACATAAGGCAACTCATCGGGGATACACATCGCACTCGGTCGGATATGATGCCGCA
TTTTACGGTTGCCTCAATCAAGCCCGCAACGGGTTAGTGAACAAGATGAGGCGACACTCACATCGCTATCT
GT
```



Example 3: Global alignment of two random 300 residue amino acid sequences

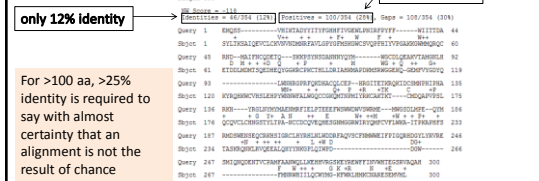
We will generate random sequences of 300 aain Matlab:

```
>> rand_int = floor(20*rand([1,300])+1);
>> rand_aa = int2aa(rand_int)
```

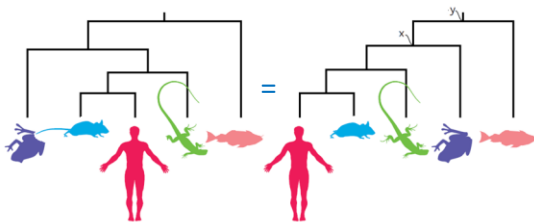
```
rand_aa =
EMQSSVHIKTADYYITYFGHHFIVGEWLPNIRFPYFFWIIITDARDNM
AIFNCQDEBTQSKKPSYNSDANNYQYMWGCDLQEAKV TAMGNLHLWNH
RGPFRFQKDHACQLCEPHRGI TETKRQKIDCSMNPHI PHARKHYRGLNY
MYMAENMR FIELPTEEEFWSWDDVSWRMEMWGS DLMPEQYMRMSWE
NSEQQRKHSIGRCLHYRHLNLWDDRFAQVSCFNMWWEI FPIGRQHDGY
LYRVRESMIQNQDENTVCPAMFAANWQLLKEHHVRSKEYREWFIFINV
WHTEGSRVAQAH
```

Example 3: Global alignment of two random 300 residue amino acid sequences

```
>Random seq 1
EMQSSVHIKTADYYITYFGHHFIVGEWLPNIRFPYFFWITTDARDNMAIFNCQDEBTQSKKPSYNSDANNYQYMW/
GCDLQEAKV TAMGNLHLWNHRGPRFKDHAQDLCEPHRGI TETKRQKIDCSMNPHI PHARKHYRGLNYMYMAE
NMRFIELPTEEEFWSWDDVSWRMEMWGS DLMPEQYMRMSWENSEQQRKHSIGRCLHYRHLNLWDDRFA
QVSCFNMWWEIFPIGRQHDGYLYRVRESMIQNQDENTVCPAMFAANWQLLKEHHVRSKEYREWFIFINVWHT
E GSRVAQAH
>Random seq 2
SYLTKSAIQEVCLCKVNVDMNRFALVPGYFGMSKWCSCVOPHYVPGA KGWMMRQRCETD DDLMDMTSQE
DHEQGGKRCPKCTHLDRIAHKMAPDKMSRWGEGKEGMFVGDYQYRQHKWCVHSLHPYWNWFWALW
GQCCGKQMTNPMYRKA KTKTMDQAPVPSLQCVLCHNGSYLTPANCDCCQEQHESGNMGGGRWRYQ
MFCVFLWKAITPAKPHSTASKRQNLKRVQEALQHYNKGPLQWPDGGWFMNMRWHILQEQWYMGKFWRLH
MKCNARESEMVML
```



Tree challenge



Is the frog more closely related to the fish or the human?

Suggested reading: **The Tree-Thinking Challenge**, Baum et al. Science 2005