

# BCB 5200 Introduction to Bioinformatics

**Sequence Analysis Lab**

Bioinformatics and Computational Biology  
Saint Louis University

# Batch BLAST jobs through web

- There are two methods to do batch BLAST jobs.
  - The first is through the web interface
  - the second is using the standalone BLAST binaries and downloaded NCBI databases.
- If you are going to submit a batch BLAST search on the web we recommend that you do not submit a file of more than 40 sequences.

**BLAST**® » **blastn suite** » **RID-XRCJVK08014**[Home](#)[Recent Results](#)[Saved Strategies](#)[Help](#)**BLAST Results**[Edit and Resubmit](#)[Save Search Strategies](#)[▶ Formatting options](#)[▶ Download](#)[How to read this page](#)[Blast report description](#)**Job title: MH0068\_GL0058265 (4176 letters)****Results for:**

4:|cl|Query\_133185 MH0068\_GL0055100(2892bp) ▼

**RID****Query ID****Description****Molecule type****Query Length**

Other reports: ▶

**Graphic Summary**

1:|cl|Query\_133182 MH0068\_GL0058265(4176bp)  
2:|cl|Query\_133183 MH0068\_GL0019004(2217bp)  
3:|cl|Query\_133184 MH0068\_GL0010430(1812bp)  
4:|cl|Query\_133185 MH0068\_GL0055100(2892bp)  
5:|cl|Query\_133186 MH0068\_GL0055096(2814bp)  
6:|cl|Query\_133187 MH0074\_GL0034786(2679bp)  
7:|cl|Query\_133188 MH0068\_GL0055841(3474bp)  
8:|cl|Query\_133189 MH0068\_GL0027893(1392bp)  
9:|cl|Query\_133190 MH0068\_GL0026512(4713bp)  
10:|cl|Query\_133191 MH0068\_GL0012576(1923bp)  
11:|cl|Query\_133192 MH0006\_GL0176780(1536bp)  
12:|cl|Query\_133193 MH0177\_GL0078034(2880bp)  
13:|cl|Query\_133194 MH0095\_GL0085358(2475bp)  
14:|cl|Query\_133195 MH0074\_GL0048502(4395bp)  
15:|cl|Query\_133196 MH0068\_GL0037198(1983bp)  
16:|cl|Query\_133197 MH0095\_GL0040817(2124bp)  
17:|cl|Query\_133198 MH0095\_GL0025062(2637bp)  
18:|cl|Query\_133199 MH0006\_GL0031342(1584bp)  
19:|cl|Query\_133200 MH0074\_GL0030661(3480bp)  
20:|cl|Query\_133201 MH0019\_GL0006447(2280bp)

**Database Name** nr**Description** Nucleotide collection (nt)**Program** BLASTN 2.7.0+ ▶ [Citation](#)[tree of results](#)**0 Blast Hits on 1000 subject sequences** ⓘ

e the title, click to show alignments

**y for alignment scores**

■ 50-80

■ 80-200

■ ≥ 200

**Query**

NCBI/ BLAST/ blastp suite/ Formatting Results - U6PWDXXWJ014

**Your search is limited to records matching entrez query: txid4932 [ORGN].**

[Edit and Resubmit](#)
[Save Search Strategies](#)
[Formatting options](#)
[Download](#)

Download					
Alignment			Search Strategies	Bioseq	
<a href="#">Text</a>	<a href="#">XML</a>	<a href="#">ASN.1</a>	<a href="#">Hit Table(text)</a>	<a href="#">Hit Table(csv)</a>	
			<a href="#">ASN.1</a>	<a href="#">ASN.1</a>	

## 20 sequences (YAL015C Chr1 complement(126904..128103)...

**Results for:** 1:lc|63955 YAL015C Chr1 complement(126904..128103) [1200 bp, 399 aa] DNA N-glycosylase and apurinic/aprimidin...(400aa) What's this?

**Query ID** lc|63955

**Description** YAL015C Chr1 complement(126904..128103) [1200 bp, 399 aa] DNA N-glycosylase and apurinic/aprimidin (AP) lyase involved in base excision repair, localizes to the nucleus and mitochondrion

**Molecule type** amino acid

**Query Length** 400

**Database Name** refseq\_protein

**Description** NCBI Protein Reference Sequences

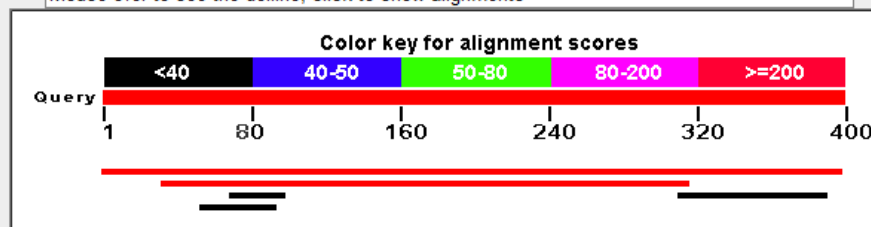
**Program** BLASTP 2.2.25+ [Citation](#)

Other reports: [Search Summary](#) [Taxonomy reports](#) [Distance tree of results](#) [Multiple alignment](#)

## Graphic Summary

### Distribution of 5 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments



## Descriptions

Legend for links to other resources: [U](#) UniGene [E](#) GEO [G](#) Gene [S](#) Structure [M](#) Map Viewer [P](#) PubChem BioAssay

### Sequences producing significant alignments:

Accession	Description	Max score	Total score	Query coverage	E value	Links
<a href="#">NP_009387.1</a>	Ntg1p [Saccharomyces cerevisiae S288c]	<a href="#">818</a>	818	99%	0.0	<a href="#">UG</a>
<a href="#">NP_014500.1</a>	Ntg1p [Saccharomyces cerevisiae S288c]	<a href="#">805</a>	805	98%	1e-03	<a href="#">UG</a>

# Command-line BLAST+

Visit the BLAST site at NCBI (“help” tab) to find the URL for the BLAST+ download.

```
$ mkdir database # this creates a new directory
$ cd database/ # we navigate into that directory
# Enter the following, without arguments, to see a help document.
$ update_blastdb.pl
# Next get a list of all available databases
$ update_blastdb.pl --showall
$ update_blastdb.pl --showall | less
```

```
$ update_blastdb.pl refseq_protein
```

```
$ tar -zxvf refseq_protein.00.tar.gz
```

# Command-line BLAST+

- Go to NCBI-Genome database
- Search “Escherichia coli [orgn]”
- Find NC\_002695.1

Reference genomes: [\[see all organisms\]](#)

◦ [Escherichia coli O157:H7 str. Sakai](#)

Submitter: GIRC

Human Pathogen

Morphology: Gram:Negative, Shape:Bacilli, Motility:Yes

Environment: OxygenReq:Facultative, OptimumTemperature:37, Temperat

Phenotype: Disease:Hemorrhagic colitis

Type	Name	RefSeq	Size (Mb)
Chr	-	NC_002695.1	5.5
Plasm	pOSAK1	NC_002127.1	0
Plasm	pO157	NC_002128.1	0.09

◦ [Escherichia coli IA139](#)

Submitter: Genoscope

Human Pathogen

Morphology: Gram:Negative, Shape:Bacilli, Motility:No

Environment: OxygenReq:Facultative, TemperatureRange:Mesophilic, Hemorrhagic Escherichia coli O157:H7

Phenotype: Disease:Urinary tract infection

Type	Name	RefSeq	INSDC	Size (Mb)
Chr	-	NC_011750.1	CU928164.2	5.13

◦ [Escherichia coli str. K-12 substr. MG1655](#)

Submitter: Univ. Wisconsin

## Related information

[Assembly](#)

[BioProject](#)

[Components \(Core\)](#)

[Full text in PMC](#)

[Gene](#)

[Genome](#)

[Identical GenBank Sequence](#)

[Probe](#)

[Protein](#)

[PubMed](#)

[PubMed \(Weighted\)](#)

[Reference Genome BioProject](#)

[Representative Genome BioProject](#)

[Taxonomy](#)

# Command-line BLAST+

- Download 5200 proteins in the fasta format

Protein  [Advanced](#) [Search](#)

Summary ▾ 20 per page ▾ Sort by Default order ▾ [Send to:](#) [Filters: Manage Filters](#)

**Items: 1 to 20 of 5200**

<< First < Prev Page

- ☐ [crotonobetaine/carnitine-CoA ligase \[Escherichia coli O157:H7 str. Sakai\]](#)  
1. 517 aa protein  
Accession: NP\_308067.2 GI: 1134749681  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)
- ☐ [leucine transcriptional activator \[Escherichia coli O157:H7 str. Sakai\]](#)  
2. 314 aa protein  
Accession: NP\_308107.3 GI: 1134749680  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)
- ☐ [poly\(A\) polymerase \[Escherichia coli O157:H7 str. Sakai\]](#)  
3. 465 aa protein  
Accession: NP\_308174.3 GI: 1134749679  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)
- ☐ [glutamyl-Q tRNA\(Asp\) synthetase \[Escherichia coli O157:H7 str. Sakai\]](#)  
4. 308 aa protein  
Accession: NP\_308175.3 GI: 1134749678

**Choose Destination**

☒ File ☐ Clipboard  
☐ Collections

Download 5200 items.

Format

- FASTA ▾
- Summary
- GenPept
- GenPept (full)
- FASTA
- ASN.1
- XML
- INSDSeq XML
- TinySeq XML
- Feature Table
- Accession List
- GI List
- GFF3

complete genome  
Escherichia coli  
Escherichia coli  
Protein Links for (323425)

# Command-line BLAST+

- Download 5200 proteins in the fasta format
- Go to terminal, log in with your account
- `mkdir course`
- `cd course`
- `cat> Ecoli.O157.fasta`
- Paste the sequences
- Ctrl+D for ending terminal line input



# Command-line BLAST+

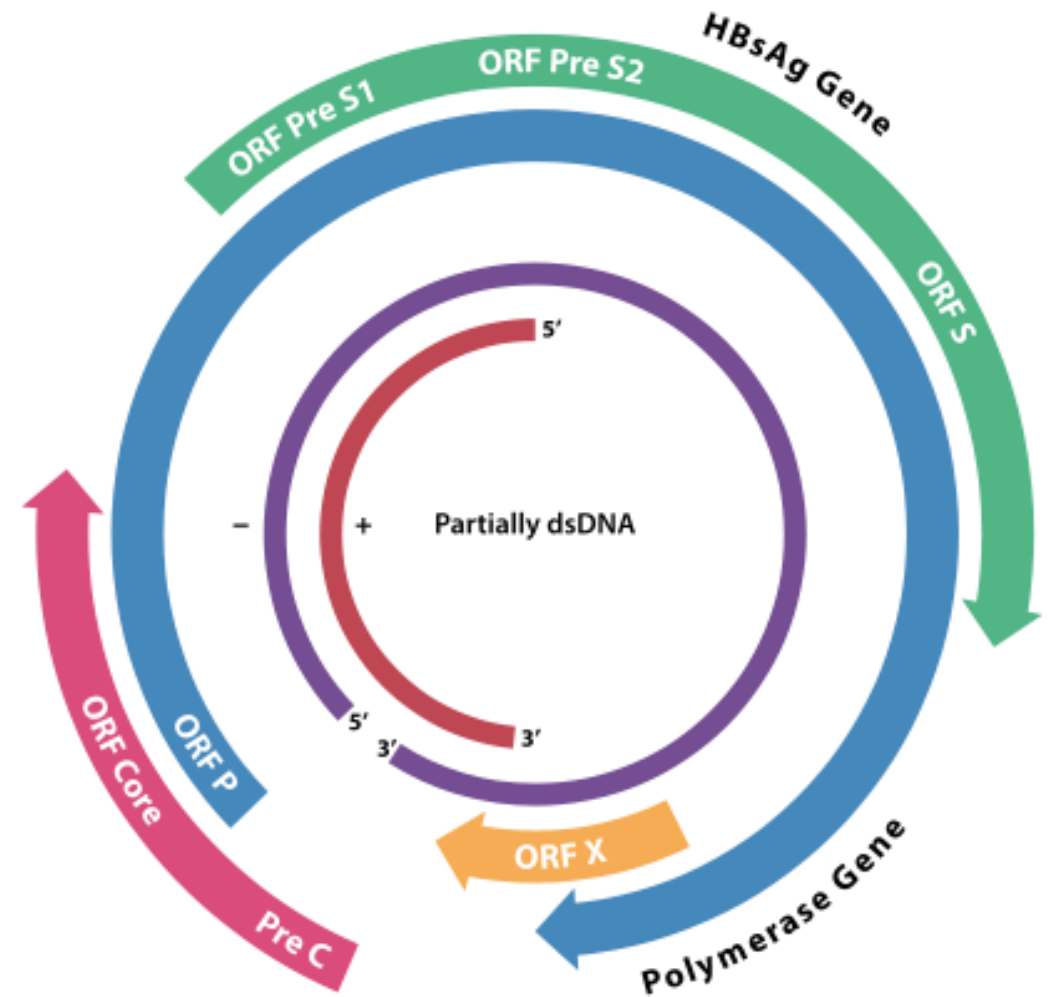
***makeblastdb to make BLAST database (-help)***

- `makeblastdb -in Ecoli.O157.fa -parse_seqids -dbtype prot`
- `cat>RHS.fasta`
- `blastp -query RHS.fasta -db Ecoli.O157.fasta -out RHS.Ecoli.O157.result`
- `less RHS.Ecoli.O157.result`

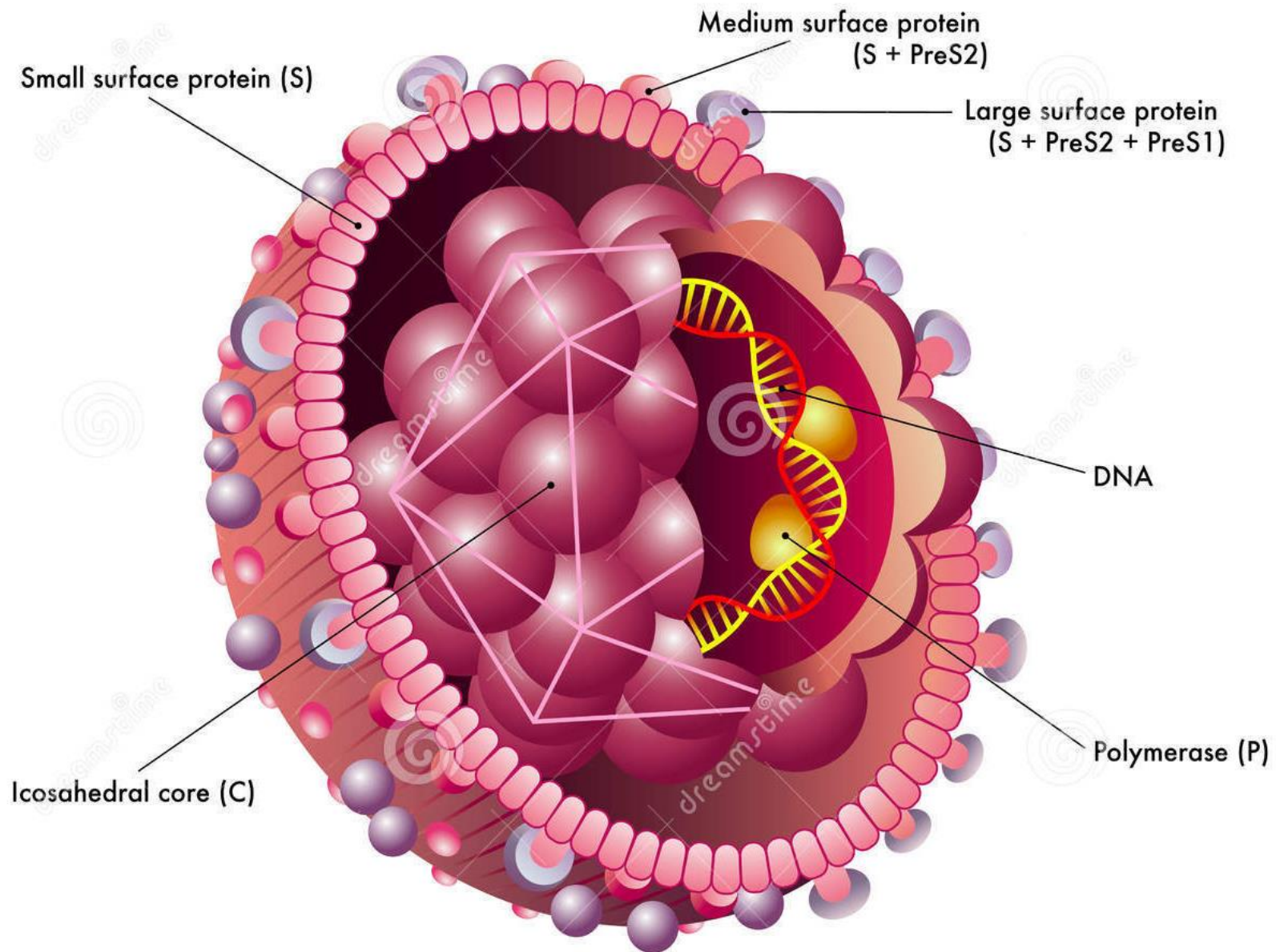
# Understanding function of HBV surface antigen (CAA53344)

# Hepatitis B virus (HBV)

- A double-stranded DNA virus belonging to the Hepadnaviridae family
- Causes human liver disease , hepatitis B
- Genome size: 3182 and 3248 bp depending on genotypes
- The genome encodes four overlapping open reading frames (ORFs):
  - Core protein (HBcAg, 2 isoforms)
  - Polymerase/reverse transcriptase (RT)
  - Surface proteins (HBsAg, 3 isoforms)
  - HBx (regulatory unit, transcriptional transactivator)



# Hepatitis B virus



# HBV genome

- [https://www.ncbi.nlm.nih.gov/nuccore/NC\\_003977.2](https://www.ncbi.nlm.nih.gov/nuccore/NC_003977.2)

# Fasta format

```
>YP_009173869.1 large envelope protein [Hepatitis B virus]
MGQNLSTSNPLGFFPDHQLDPAFRANTANPDWDFNPNKDTWPDANKVGAGAFGLGFTPPHGGLLGWSPQA
QGILQTLPANPPPASTNRQSGRQPTPLSPPLRNTHPQAMQWNSTTFHQTLQDPRVRGLYFPAGGSSSGTV
NPVLTTASPLSSIFSRI GDPALNMENITSGFLGPLLVLQAGFFLLTRILTIPQSLDSWWTSLNFLGGTTV
CLGQNSQSPTSNHSPTSCPPTCPGYRWMCLRRFIIFFILLCLIFLLVLLDYQGMLPVCPLIPGSSTTS
TGPCRTCMTTAQGTSMPSCCCTKPSDGNCTCIPIPSSWAFGKFLWEWASARFSWLSLLVPFVQWFVGLS
PTVWLSVIWMMWYWGPSLYSILSPFLPLLPIFFCLWVYI
```

# How to understand its function?

- PubMed search?
- BLAST search
  - Homologs?
  - Species distribution?
  - Any likely features?

# Choice of BLAST programs?

- BLASTP
- PSI-BLAST
- DELTA-BLAST
- RPS-BLAST
- Tips on parameters
  - Database (NR vs RefSeq)
  - Max targets (20,000)
  - Threshold (0.05 vs 0.001)
  - Filter (low complexity)
  - Taxonomy report (restrict organism)



# Iteration 1 of PSI-BLAST

1. Review the BLAST result from iteration 1 of PSI-BLAST
  - Go to Format options to change the parameters
2. Too many sequences
  - If you are interested in the human HBV mutations
  - If you are interested in a more broad functional features
3. Restrict hits from certain lineages, clades or species
  - What can you get?

# Iteration 2 of PSI-BLAST


Other reports: [Search Summary](#) [Taxonomy reports](#) [Distance tree of results](#) [Multiple alignment](#) [MSA viewer](#)

**New** Analyze your query with [SmartBLAST](#)

## Graphic Summary

☒ Show Conserved Domains

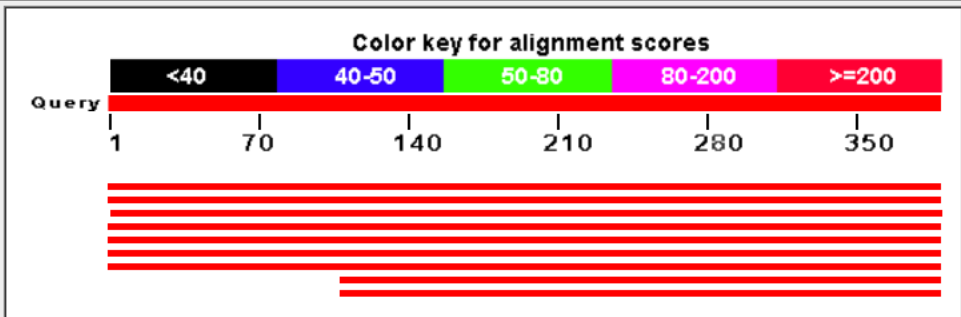
Putative conserved domains have been detected, click on the image below for detailed results.



The diagram shows a query sequence of 389 amino acids. A light blue bar labeled 'vMSA' spans from position 1 to 389. Below it, a yellow bar labeled 'Superfamilies' is shown, with a gap between position 1 and 50, and then a continuous bar from 50 to 389.

Distribution of the top 9 Blast Hits on 9 subject sequences ⓘ

Mouse over to see the define, click to show alignments



The color key for alignment scores is as follows:

Score Range	Color
<40	Black
40-50	Blue
50-80	Green
80-200	Pink
>=200	Red

The query sequence is shown with positions 1, 70, 140, 210, 280, and 350 marked. Below the query sequence, there are 9 horizontal bars representing the top 9 blast hits. The first bar is red, indicating a score of >=200. The other 8 bars are also red, indicating scores of >=200.

## Descriptions

Legend for links to other resources: [U](#) UniGene [E](#) GEO [G](#) Gene [S](#) Structure [M](#) MIM [P](#) PubChem BioAssay

**NEW** - alignment score below the threshold on the previous iteration

● - alignment was checked on the previous iteration

Run PSI-Blast iteration 2 with max

**Sequences producing significant alignments with E-value BETTER than threshold**

Accession	Description	Max score	Total score	Query coverage	E value	Ident	Links
<a href="#">AOG18115.1</a>	PreS1/PreS2/S [Cloning vector pAAV2neo-1.3HBV(ayw)]	<a href="#">708</a>	708	100%	0.0	99%	
<a href="#">ABB04016.1</a>	PreS1 [synthetic construct]	<a href="#">708</a>	708	100%	0.0	100%	
<a href="#">APA28776.1</a>	large S protein [synthetic construct]	<a href="#">656</a>	656	99%	0.0	91%	
<a href="#">AOG18108.1</a>	PreS1/PreS2/S [Cloning vector pAAV2neo-1.3HBV(adr)]	<a href="#">646</a>	646	100%	0.0	89%	
<a href="#">AAL13124.1</a>	pre-S1/pre-S2 protein [synthetic construct]	<a href="#">646</a>	646	100%	0.0	89%	
<a href="#">AAK58874.1</a>	PreS1+PreS2+HBsAg [synthetic construct]	<a href="#">646</a>	646	100%	0.0	89%	
<a href="#">ALV66577.1</a>	large S protein [Heron hepatitis B virus]	<a href="#">646</a>	646	100%	0.0	89%	
<a href="#">AOG18116.1</a>	PreS2/S [Cloning vector pAAV2neo-1.3HBV(ayw)]	<a href="#">646</a>	646	100%	0.0	89%	
<a href="#">AOG18109.1</a>	PreS2/S [Cloning vector pAAV2neo-1.3HBV(adr)]	<a href="#">646</a>	646	100%	0.0	89%	
<a href="#">AAT37474.1</a>	Hepatitis B surface antigen/human papillomavirus EE7-FLAG hybrid protein	<a href="#">357</a>	357	58%	1e-120	90%	
<a href="#">AAT37475.1</a>	Hepatitis B surface antigen/human papillomavirus EE7 (D1-35) hybrid protein	<a href="#">358</a>	358	99%	2e-118	56%	
<a href="#">YP_009175035.1</a>	preS1 surface protein [Woolly monkey hepatitis B Virus]	<a href="#">355</a>	355	99%	5e-117	55%	
<a href="#">AAC64339.1</a>	surface antigen subtype ayw [Expression vector pUK-S]						
<a href="#">KRY25800.1</a>	Large envelope protein [Trichinella spiralis]						
<a href="#">AIT39699.1</a>	DSV4 [synthetic construct]						
<a href="#">APA28777.1</a>	S protein [synthetic construct]						
<a href="#">AGH10171.1</a>	surface antigen [Bat hepatitis virus]						
<a href="#">WP_080445144.1</a>	hypothetical protein [Bacillus cereus]						
<a href="#">AAM28153.1</a>	surface antigen [synthetic construct]						
<a href="#">AOG18110.1</a>	S [Cloning vector pAAV2neo-1.3HBV(adr)]						
<a href="#">AAL13120.1</a>	S protein [synthetic construct]						
<a href="#">ALV66578.1</a>	S protein [Heron hepatitis B virus]						
<a href="#">ARM20234.1</a>	surface antigen [Bat hepatitis virus]						
<a href="#">AAL62458.1</a>	preS/S [Duck hepatitis B virus]						
<a href="#">AHC08491.1</a>	S protein [synthetic construct]						
<a href="#">AGH10175.1</a>	surface antigen [Bat hepatitis virus]						
<a href="#">YP_007678000.1</a>	surface antigen [Bat hepatitis virus]						

Pay attention to:

1. New sequences are identified

- Woolly monkey HBV
- Bacillus cereus
- Bat HBV
- Duck HBV

2. E-value changed

# Search reaches convergence in iteration 4

**BLAST Results**

ⓘ Your results are being filtered to match entrez query: all [filter] NOT(txid10407 [ORGN]).  
[Edit and Resubmit](#) [Save Search Strategies](#) [▶ Formatting options](#) [▶ Download](#) [YouTube](#) [How to read this page](#) [Blast report c](#)

PSI blast Iteration 4

**Job title: (2) - Protein Sequence (389 letters)**

<b>RID</b>	<a href="#">XRWC3V1C015</a> (Expires on 10-11 06:00 am)
<b>Query ID</b>	Id Query_175982
<b>Description</b>	None
<b>Molecule type</b>	amino acid
<b>Query Length</b>	389

<b>Database Name</b>	nr
<b>Description</b>	All non-redundant GenBank CDS translations+PDB+SwissProt+PIR+PRF excluding environmental samples from WGS projects
<b>Program</b>	BLASTP 2.7.0+ <a href="#">▶ Citation</a>

ⓘ No new sequences were found above the 0.005 threshold

Other reports: [▶ Search Summary](#) [\[Taxonomy reports\]](#) [\[Distance tree of results\]](#) [\[Related Structures\]](#) [\[Multiple alignment\]](#) [\[MSA viewer\]](#)

**New** Analyze your query with [SmartBLAST](#)

**Graphic Summary**

**Distribution of the top 34 Blast Hits on 34 subject sequences** ⓘ

Mouse-over to show define and scores, click to show alignments

**Color key for alignment scores**

<40	40-50	50-80	80-200	>=200
-----	-------	-------	--------	-------

**Query**

1	70	140	210	280	350

# Now let's use another query

```
>ALV66578.1 S protein [Heron hepatitis B virus]  
MENTTSGFLGPLLVLQAGFFSLTRILTIPQSLDSWWTSLNFLGGAPTCPGQNSQSPTSNHSPTSCPPICP  
GYRWMCLRRFIIIFLFIILLCLIFLLVLLDYQGMLPVCPLLPGTSTTSTGPCKTCTIPAQGTSMEFPSCCCT  
KPSDGNCTCIPIPSSWAFARFLWEWASVRFSWLSLLVPFVQWFVGLSPTVWLSVIWMMWYWGPSLYNILS  
PFLPLLPIFFYLWVYI
```

# Many more hits

NEW	✓ WP_080445144.1	hypothetical protein [Bacillus cereus]	368	368	100%	1e-127	92%	
NEW	✓ AAK58874.1	PreS1+PreS2+HBsAg [synthetic construct]	374	374	100%	4e-127	92%	
NEW	✓ APA28776.1	large S protein [synthetic construct]	373	373	100%	5e-127	91%	
NEW	✓ AAT37474.1	Hepatitis B surface antigen/human papillomavirus EE7-FLAG hybrid protei	368	368	100%	1e-125	92%	
NEW	✓ AAT37475.1	Hepatitis B surface antigen/human papillomavirus EE7 (D1-35) hybrid pro	365	365	100%	2e-125	92%	
NEW	✓ AAC64339.1	surface antigen subtype ayw [Expression vector pUK-S]	362	362	100%	5e-125	92%	
NEW	✓ AOG18116.1	PreS2/S [Cloning vector pAAV2neo-1.3HBV(ayw)]	364	364	100%	5e-125	92%	
NEW	✓ AAL62458.1	preS/S [Duck hepatitis B virus]	360	360	100%	2e-124	93%	
NEW	✓ ABB04016.1	PreS1 [synthetic construct]	366	366	100%	3e-124	92%	
NEW	✓ AOG18115.1	PreS1/PreS2/S [Cloning vector pAAV2neo-1.3HBV(ayw)]	365	365	100%	4e-124	92%	
NEW	✓ YP_009045992.1	surface protein [Roundleaf bat hepatitis B virus]	269	269	100%	2e-88	68%	G
NEW	✓ YP_009175036.1	surface protein [Woolly monkey hepatitis B Virus]	269	269	92%	2e-88	77%	G
NEW	✓ AGH10171.1	surface antigen [Bat hepatitis virus]	275	275	98%	2e-88	65%	
NEW	✓ YP_009175035.1	preS1 surface protein [Woolly monkey hepatitis B Virus]	271	271	92%	6e-87	77%	G
NEW	✓ YP_009045996.1	surface protein [Horseshoe bat hepatitis B virus]	260	260	100%	8e-85	65%	G
NEW	✓ AQT40957.1	surface protein [Horseshoe bat hepatitis B virus]	262	262	100%	4e-83	65%	
NEW	✓ AGH10175.1	surface antigen [Bat hepatitis virus]	258	258	98%	7e-82	64%	
NEW	✓ YP_007678000.1	surface antigen [Bat hepatitis virus]	257	257	98%	2e-81	65%	G
NEW	✓ ARM20226.1	surface antigen [Bat hepatitis virus]	256	256	98%	6e-81	65%	
NEW	✓ ARM20234.1	surface antigen [Bat hepatitis virus]	256	256	98%	7e-81	65%	
NEW	✓ ARM20230.1	surface antigen [Bat hepatitis virus]	256	256	98%	9e-81	64%	
NEW	✓ AIW47284.1	surface antigen [Bat hepatitis virus]	254	254	98%	7e-80	65%	
NEW	✓ ARM20218.1	surface antigen [Bat hepatitis virus]	252	252	100%	3e-79	62%	
NEW	✓ NP_955537.1	surface antigen [Ground squirrel hepatitis virus]	239	239	93%	2e-76	63%	G
NEW	✓ NP_040995.1	hypothetical protein [Ground squirrel hepatitis virus]	239	239	93%	6e-76	63%	G
NEW	✓ AAB08035.1	small envelope protein [Arctic ground squirrel hepatitis B virus]	235	235	93%	6e-75	64%	
NEW	✓ AAB08034.1	middle envelope protein [Arctic ground squirrel hepatitis B virus]	235	235	93%	2e-74	64%	
NEW	✓ ABI31620.1	major surface antigen [Woodchuck hepatitis virus]	233	233	93%	3e-74	62%	
NEW	✓ P03144.1	RecName: Full=Large envelope protein; AltName: Full=L glycoprotein; Alt	239	239	93%	8e-74	63%	
NEW	✓ AGW01291.1	surface protein [Tent-making bat hepatitis B virus]	232	232	99%	1e-73	59%	
NEW	✓ YP_009046000.1	surface protein [Tent-making bat hepatitis B virus]	231	231	99%	1e-73	59%	G
NEW	✓ NP_944491.1	surface protein [Woodchuck hepatitis virus]	231	231	93%	1e-73	62%	G
NEW	✓ AAA46775.1	surface protein [Woodchuck hepatitis virus]	230	230	93%	4e-73	61%	
NEW	✓ AAA46773.1	surface antigen [Woodchuck hepatitis virus]	229	229	93%	6e-73	61%	

... <a href="#">Duck hepatitis B virus isolate DHBVQCA34</a>	<a href="#">viruses</a>	47.4	1	<a href="#">Duck hepatitis B virus isolate DHBVQCA34 hits</a>
... <a href="#">Duck hepatitis B virus strain China</a>	<a href="#">viruses</a>	47.4	1	<a href="#">Duck hepatitis B virus strain China hits</a>
... <a href="#">Stork hepatitis B virus</a>	<a href="#">viruses</a>	47.0	4	<a href="#">Stork hepatitis B virus hits</a>
... <a href="#">Hepatitis B virus duck/DHBV-16</a>	<a href="#">viruses</a>	45.8	1	<a href="#">Hepatitis B virus duck/DHBV-16 hits</a>
... <a href="#">Duck hepatitis B virus brown Shanghai duck/S5</a>	<a href="#">viruses</a>	45.4	1	<a href="#">Duck hepatitis B virus brown Shanghai duck/S5 hits</a>
.. <a href="#">Roundleaf bat hepatitis B virus</a>	<a href="#">viruses</a>	269	5	<a href="#">Roundleaf bat hepatitis B virus hits</a>
.. <a href="#">Woolly monkey hepatitis B virus</a>	<a href="#">viruses</a>	269	6	<a href="#">Woolly monkey hepatitis B Virus hits</a>
.. <a href="#">Bat hepatitis virus</a>	<a href="#">viruses</a>	275	13	<a href="#">Bat hepatitis virus hits</a>
.. <a href="#">Hepatitis B virus Woolly monkey/Louisville</a>	<a href="#">viruses</a>	271	1	<a href="#">Hepatitis B virus Woolly monkey/Louisville hits</a>
.. <a href="#">Horseshoe bat hepatitis B virus</a>	<a href="#">viruses</a>	260	3	<a href="#">Horseshoe bat hepatitis B virus hits</a>
.. <a href="#">Ground squirrel hepatitis virus</a>	<a href="#">viruses</a>	239	4	<a href="#">Ground squirrel hepatitis virus hits</a>
.. <a href="#">Arctic ground squirrel hepatitis B virus</a>	<a href="#">viruses</a>	235	3	<a href="#">Arctic ground squirrel hepatitis B virus hits</a>
.. <a href="#">Woodchuck hepatitis virus</a>	<a href="#">viruses</a>	233	34	<a href="#">Woodchuck hepatitis virus hits</a>
.. <a href="#">Tent-making bat hepatitis B virus</a>	<a href="#">viruses</a>	232	5	<a href="#">Tent-making bat hepatitis B virus hits</a>
.. <a href="#">Woodchuck hepatitis virus w64 (ISOLATE PWS23)</a>	<a href="#">viruses</a>	229	1	<a href="#">Woodchuck hepatitis virus w64 (ISOLATE PWS23) hits</a>
.. <a href="#">Woodchuck hepatitis virus 7</a>	<a href="#">viruses</a>	231	1	<a href="#">Woodchuck hepatitis virus 7 hits</a>
.. <a href="#">Woodchuck hepatitis virus 59</a>	<a href="#">viruses</a>	231	1	<a href="#">Woodchuck hepatitis virus 59 hits</a>
.. <a href="#">Woodchuck hepatitis virus 8</a>	<a href="#">viruses</a>	231	1	<a href="#">Woodchuck hepatitis virus 8 hits</a>
.. <a href="#">Woodchuck hepatitis virus 1</a>	<a href="#">viruses</a>	229	1	<a href="#">Woodchuck hepatitis virus 1 hits</a>
.. <a href="#">Woodchuck hepatitis virus 2</a>	<a href="#">viruses</a>	225	1	<a href="#">Woodchuck hepatitis virus 2 hits</a>
.. <a href="#">Bluegill hepadnavirus</a>	<a href="#">viruses</a>	107	2	<a href="#">Bluegill hepadnavirus hits</a>
.. <a href="#">Ross's goose hepatitis B virus</a>	<a href="#">viruses</a>	66.2	6	<a href="#">Ross's goose hepatitis B virus hits</a>
.. <a href="#">White sucker hepatitis B virus</a>	<a href="#">viruses</a>	66.2	2	<a href="#">White sucker hepatitis B virus hits</a>
.. <a href="#">Snow goose hepatitis B virus</a>	<a href="#">viruses</a>	52.4	12	<a href="#">Snow goose hepatitis B virus hits</a>
.. <a href="#">Sheldgoose hepatitis B virus</a>	<a href="#">viruses</a>	51.2	6	<a href="#">Sheldgoose hepatitis B virus hits</a>
.. <a href="#">Tibetan frog hepadnavirus</a>	<a href="#">viruses</a>	51.2	4	<a href="#">Tibetan frog hepadnavirus hits</a>

Check the taxonomy report

# Search converged at Iteration 4

## BLAST Results

**Your search is limited to records that exclude: Human hepatitis B virus (taxid:10407)** [Full Entrez Query](#)

[Edit and Resubmit](#) [Save Search Strategies](#) [Formatting options](#) [Download](#)

[YouTube](#) [How to read this page](#) [Blast report description](#)

PSI blast Iteration 4

**Job title: (2) - ALV66578:S protein [Heron hepatitis B virus]**

**RID** [XS57YPUF015](#) (Expires on 10-11 08:31 am)

**Query ID** [ALV66578.1](#)

**Description** S protein [Heron hepatitis B virus]

**Molecule type** amino acid

**Query Length** 226

**Database Name** nr

**Description** All non-redundant GenBank CDS translations+PDB+SwissProt+PIR+PRF excluding environmental samples from WGS projects

**Program** BLASTP 2.7.0+ [Citation](#)

**No new sequences were found above the 0.005 threshold**

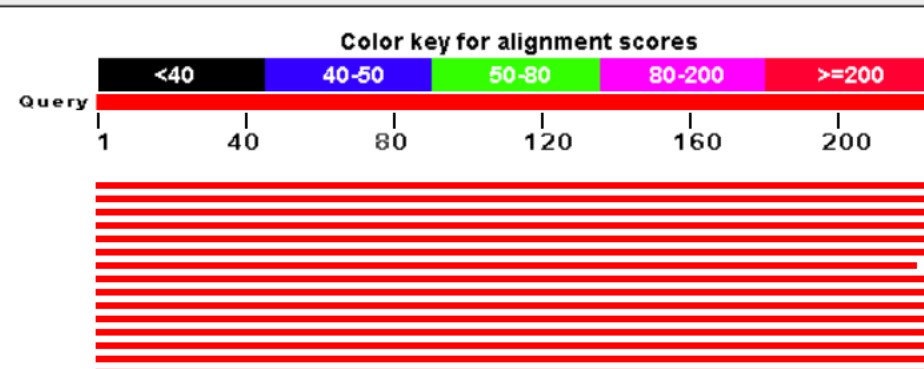
Other reports: [Search Summary](#) [Taxonomy reports](#) [Distance tree of results](#) [Related Structures](#) [Multiple alignment](#) [MSA viewer](#)

**New** Analyze your query with [SmartBLAST](#)

### [Graphic Summary](#)

**Distribution of the top 100 Blast Hits on 199 subject sequences**

Mouse over to see the define, click to show alignments





# PSI-BLAST searches (HBV Surface protein)



- Heron hepatitis B virus
- Bat hepatitis B virus
- Woodchuck hepatitis virus



- Ground squirrel hepatitis B virus
- Bluegill hepadnavirus
- Tibetan frog hepadnavirus



- Duck hepatitis B virus
- Woolly monkey hepatitis B Virus
- Snow goose hepatitis B virus



- Trichinella spiralis



- Download the Text or Hit Table files

**BLAST**® » **blastp suite** » RID-XRWC3V1C015 Home Recent Results

### BLAST Results

ⓘ Your results are being filtered to match entrez query: all [filter] NOT(txid10407 [ORGN]).

[Edit and Resubmit](#) [Save Search Strategies](#) [▶ Formatting options](#) [▼ Download](#) YouTube [How to read](#)

Download					
<a href="#">Text</a> <a href="#">XML</a> <a href="#">ASN.1</a> <a href="#">JSON Seq-align</a> <a href="#">Hit Table(text)</a> <a href="#">Hit Table(csv)</a>		<a href="#">Search Strategies</a> <a href="#">ASN.1</a>		<a href="#">PSSM to restart search</a> <a href="#">PSSM</a>	
<a href="#">Multiple-file XML2</a> <a href="#">Single-file XML2</a> <a href="#">Multiple-file JSON</a> <a href="#">Single-file JSON</a>					

PSI blast Iteration 4

**Job title: (2) - Protein Sequence (389 letters)**

<b>RID</b>	<a href="#">XRWC3V1C015</a> (Expires on 10-11 06:00 am)	<b>Database Name</b>	nr
<b>Query ID</b>	Id Query_175982	<b>Description</b>	All non-redundant GenBank CDS translations+PDB+Swiss environmental samples from WGS projects
<b>Molecule type</b>	amino acid	<b>Program</b>	BLASTP 2.7.0+ <a href="#">▶ Citation</a>
<b>Query Length</b>	389		

ⓘ No new sequences were found above the 0.005 threshold

Other reports: [▶ Search Summary](#) [\[Taxonomy reports\]](#) [\[Distance tree of results\]](#) [\[Related Structures\]](#) [\[Multiple alignment\]](#) [\[MSA viewer\]](#)

**New** Analyze your query with [SmartBLAST](#)

[Graphic Summary](#)

Distribution of the top 34 Blast Hits on 34 subject sequences

# From sequence ids to fasta format

- Put the sequence id list to a text file in Unix  
cat>te  
paste your list  
Ctrl+D for ending terminal line input
- cat te | cut -f 1 -d ' ' > HBVS.ac.list
  - (NOTE: Remove “prf| |1803562A, prf| |1305266A” in your file)
- cat HBVS.ac.list | epost -db protein -format acc | efetch -format fasta > HBVS.ac.fasta
  - Or use “grep -v '\|' HBVS.ac.list | epost -db protein -format acc | efetch -format fasta > HBVS.ac.fasta
- Or Batch Entrez
  - Given a file of Entrez accession numbers or other identifiers, Batch Entrez downloads the corresponding records
  - [www.ncbi.nlm.nih.gov/sites/batchentrez](http://www.ncbi.nlm.nih.gov/sites/batchentrez)

# Multiple sequence alignment

- ClustalW or Omega: <https://www.ebi.ac.uk/Tools/msa/clustalo/>
- Kalign: <https://www.ebi.ac.uk/Tools/msa/kalign/>
- T-coffee: <http://tcoffee.crg.cat/>
- Muscle: <https://www.ebi.ac.uk/Tools/msa/muscle/>
- Promals: <http://prodata.swmed.edu/promals/promals.php>

# BLASTClust to remove the highly-similar sequences

- BLASTClust is a program within the standalone BLAST package used to cluster either protein or nucleotide sequences. The program begins with pairwise matches and places a sequence in a cluster if the sequence matches at least one sequence already in the cluster. In the case of proteins, the blastp algorithm is used to compute the pairwise matches; in the case of nucleotide sequences, the Megablast algorithm is used.
- `blastclust -help`
- `cat HBVS.ac.fasta | blastclust -S 1.75 -L .9 > HBVS.bc`
- `echo "$(tail -n +2 te)" | cut -f 1 -d ' ' | less` (remove the first line + get the first column ids)

- Generate multiple sequence alignment with these sequences
- Clustal Omega
- Promals3D:
  - Advantage: profile-profile comparison and secondary structural information
  - Note: PSI-BLAST iteration number set to 1
  - [http://prodata.swmed.edu/promals3d/getResult2.php?name=QUERY\\_CBI\\_uhw&email=empty\\_email&target\\_name=](http://prodata.swmed.edu/promals3d/getResult2.php?name=QUERY_CBI_uhw&email=empty_email&target_name=)

# Multiple sequence alignment editors

- BioEdit - MS-Windows
- Genedoc - MS-Windows
- EditSeq/MegAlign - Lasergene - Mac or MS-Windows
- DNA Strider - Macintosh
- Seq-AL - Macintosh
- ASAD - Excel - Macintosh or MS-Windows
- Chroma - Windows
- SeqPup - Mac. MS-Windows, X-Windows
- AliView - Mac. MS-Windows, Linux
- Boxshade: [www.ch.embnet.org/software/BOX\\_form.html](http://www.ch.embnet.org/software/BOX_form.html)

# MSA-Visualization and improvement

- GeneDoc (Windows)
- Download: <http://genedoc.software.informer.com/download/>
  - Arranging and Editing
    - GeneDoc's Grab and Drag arrangement mode allows you to move residues around like beads on a string
  - Shading Alignments
  - Reports: Stats, Score, Composition
  - Exporting and Copying Figures



# GeneDoc: Conservation Mode

- GeneDoc (Windows): use import to load alignment file

```
GeneDoc - [Gene2]
File Project Edit Arrange Shade Groups Score Tree Reports Plot Window Help
C S G
C Q P E S H I L D G
gi|1154593 : AKEHESVWFKGRFTPNVWANTPGELQIKQIKPAIDSKGRKVGEEWTTIKVENALTRYHEACDNAKRKVLLELRGSSSELQDKI-NVLVFCSTMIIITKALFG---- : 702
gi|1180947 : SFSSVRGFFIQNADC-----SALPNGKLPSSTTKVTKTRNTYSSTADLIKNNERCQESLRIYHMTYLVCKLNEIYEHI-HCLYKLSDIYSLMLLS-FAH : 509
gi|1572857 : NFRDEVLYLIEVKNKNSQ-----IKDLPPDIKVNNTKMVSRITTPRTQKLTQKLEYKDLLIRESELQYKEFLNKKITAAY-TELKRTITLNAQYTCILS---- : 733
d 6 s

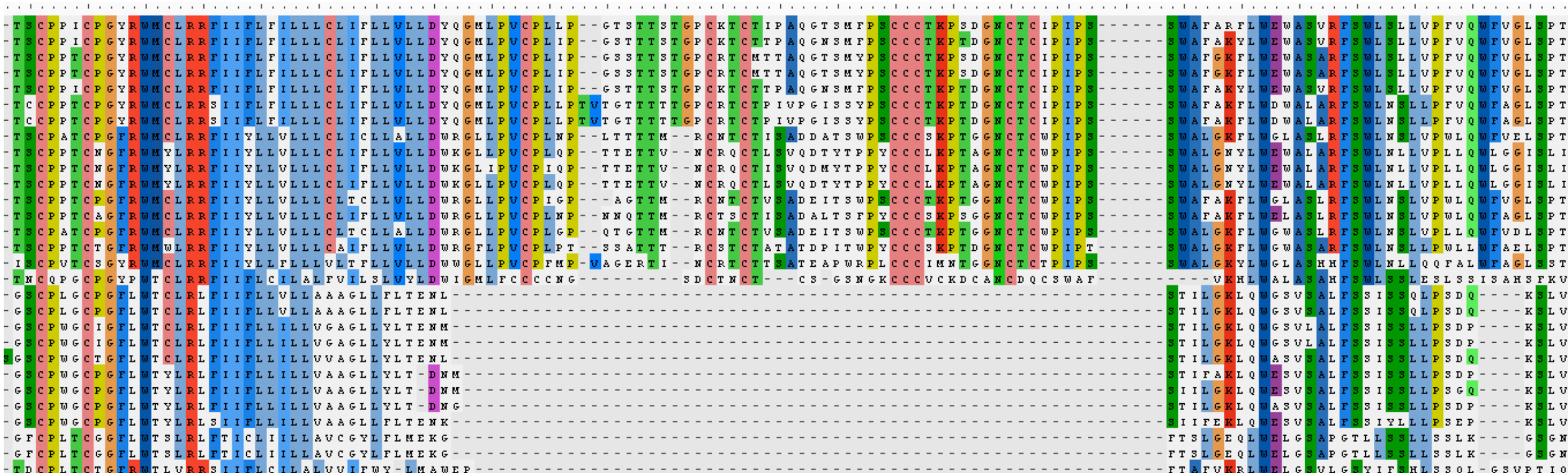
*      1200      *      1220      *      1240      *      1260      *      1280      *
gi|172000| : FAVLANERNLVCKKVDENK-----DEVNGRRLMVEEGLSARSLETFTANNCEIAKDN-----LWVIITGPNMGGKSTFLRCNAIIVILAQIG : 794
gi|1321950 : FAQVSAERNYABEQVLEDEC-----IIEINGRHALYET-----FLDNYTENSTMIDGGLFSELSWCEQNKGRIVVTCANASGKSVYITQNGLIVYLAQIG : 666
gi|1490521 : KVAK--QGDIYCRFTVQEER-----KIVIKNGRHPVIDVLL--GEQDQYVFNNTDLSED-----SERVMIIITGPNMGGKSSYIKQVALITIMAQIG : 910
gi|2463653 : ACTLS---DYVRE---EFTDT-----IAKQGWHPPILEK---ISAEPKANNFYVTEG-----SN-FLIITGPNMSGKSTYLKQIALCQIMAQIG : 703
gi|4139230 : LAASCP-TPYCRBEITSL-DAG-----DIVIEGSRHPCQEA---QDVVNFENDCRLMRG-----KSWFCIVTGNMGGKSTFIQCVGVIVLMAQVG : 689
gi|4504191 : NYSRGGDGPMPQREVILLPED--T-----PPFELIKGRHPCITKTF--FG-DDFENDIILIGCEEEQE---NGKAYCVLVITGNMGGKSTLMRCAGLIAVMAQMG : 1157
gi|4529899 : LASAARDYGYSRERYSPQVLG-----VRQNGRHPILVEL-----CARTEVENSTECGGD-----KGRVKVITGNPSSGKSIYLRQVGLITFMAIVG : 615
gi|4557761 : VSN GAP-VPYVREFAILEK-GQG-----RIILKASRHACVEV---QDEIAFENDVYFEKD-----KQMFHIIITGPNMGGKSTYIRQCTGVIVLMAQIG : 692
gi|6320302 : RTSEYLGAPSCRETIIVDEVDSKTNTQ--LNGFIFKFKSLRHPCFNLGA--TTAKDEFENDIILIGKE-----QPRLELLTGNNAAGKSTILMACIHAVIMAQMG : 1005
gi|6324482 : TSSYAP-IPYIRBKLEHMDSER-----RTHIISRHPVILEG-----QDDISEFNDVILESG-----KGDFLIITGPNMGGKSTYIRQCVGVISLMAQIG : 711
gi|1264384 : IAAASLSAGSMAREVIFPESEATDQNKTKGPIKIQGLWHPFAVA---ADGQLRVENDIILGEARRSSG---SIHPRSILLTGNMGGKSTILLRATCAVIFAQIG : 876
gi|1523522 : FASDSYEGVRCREVISGSTS-DG-----VPHISATGLGHPVIRGDS--LGRGSEVENNVKIGGA-----EKASFILLTGNMGGKSTILLRQVCLAVILAQIG : 1106
gi|3031380 : HACEGRRRKWVFETLVGFSLDE-----GAKPIDGASRMKTTGLS--PYWFDVSSGTAHVNTVD-----MQSLFLLTGNPGCGKSSLLRSICAAALIGISC : 791
gi|3068692 : TLSR--NKNYVREPFVDDCE-----PVEINQSGRHPVLELIL--QDN--EVENDTILHAE-----GEYCCIIITGPNMGGKSCYIRQVALISIMAQVG : 840
gi|3248864 : TLAK--QNKYVRENFVRENE-----ASQIHKDGGRHPVLESLL--GVN--EVENDTILHAN-----SEYCCIVTGNMGGKSCYIRQVALITLMAQVG : 829
gi|3433012 : VSN GAP-VPYVREVVLEK-GQG-----RIVIKGARHPCIEV---QDEVAFENDVTFEKG-----KQMFHIIITGPNMGGKSTYIRQCTGVIVLMAQIG : 620
gi|4847521 : LAIVARQNNYVREILTEDSI-----LEIQNGRHALQEM-----TVDTFENDTKIR-S-----SGRINIITGNPNYSKSIYIKQVALIVVFLAHIG : 595
gi|5109150 : IASDFFEGPTCCFIILIKESYGPD-----TPTIHARNLGHPTVRSDS--LGSGSEVENDIKMGGP-----GNASFIVLTGNMGGKSTILLRQVCTIILAQIG : 1027
gi|6232072 : LACVAHQNNYVREVLTVESL-----IDIRNGRHPVLELIL--QDN--EVENDTILHAE-----NGRIHIIITGNPNYSKSIYIKQVALIVVFLAHIG : 330
gi|7947696 : TISTKPVDRYSRE---ELTDSG-----PIADAGRHPVILEG---IHNDF--VSNSTFMSSEA-----TN-MLVVMGPNMSGKSTYLQCVCLVIVILAQIG : 576
gi|1154593 : HVSEGRRRGWVLEFETISPLCKDN-----VTEETIS--SEMEISGTF--PYWLDTNQGNAILNDVH-----MHSLFILITGNPGCGKSSMLRSVCAAAALIGIG : 788
gi|1180947 : ACTIS---DYVRE---EFTDT-----IAKQGWHPPILEK---IAMEKPVSNNAYLTEG-----SN-FLIITGPNMSGKSTYIKQIALCQIMAQIG : 584
gi|1572857 : LAATSCNVNYSRETFVNGQQ-----AIIAKNARNPILES---LDVHYVENDINMSPE-----NGKINIITGPNMGGKSSYIRQVALITIMAQIG : 814
rP      h      n      66tGpN GKS 6 q      aq G

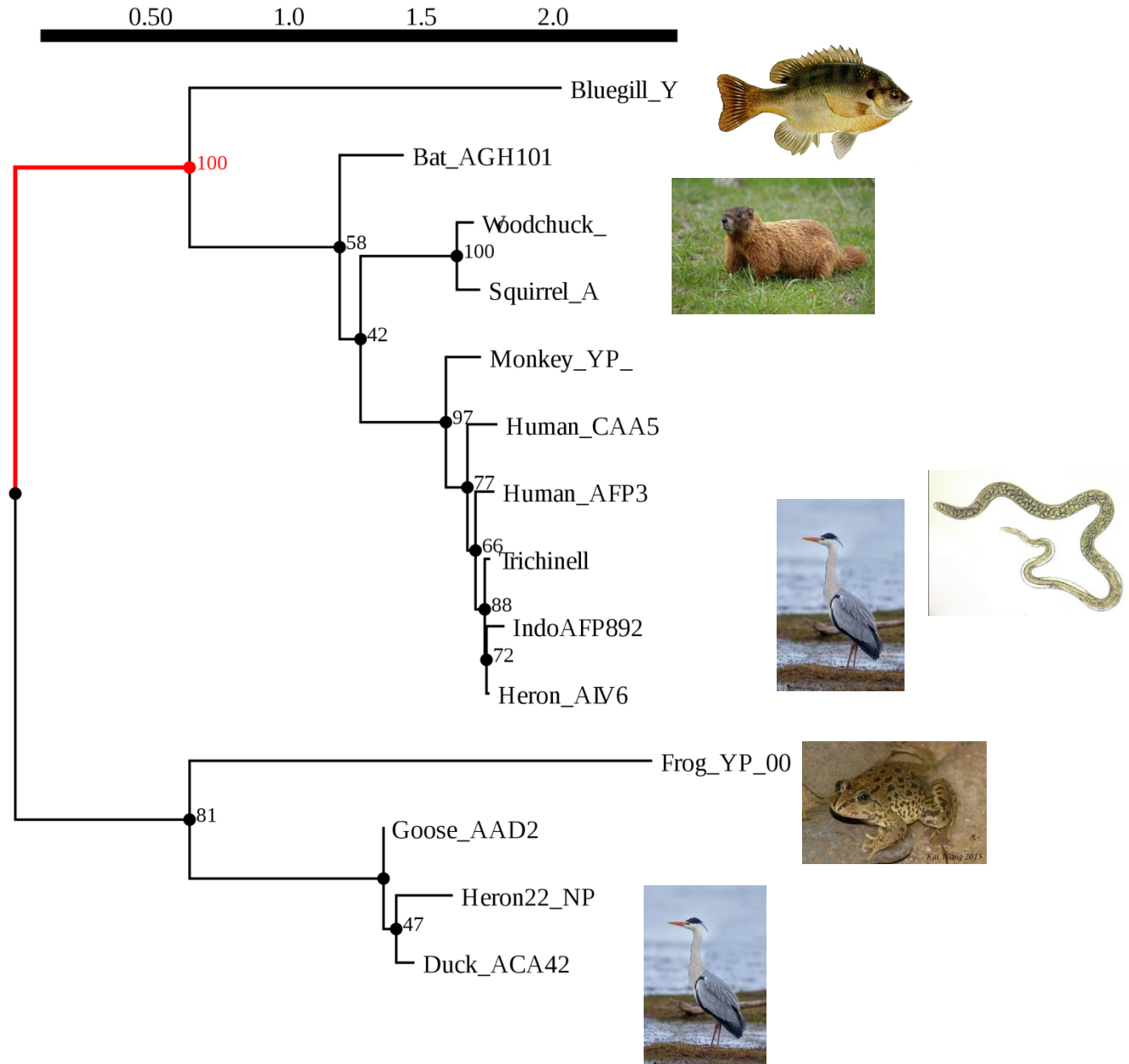
*      1300      *      1320      *      1340      *      1360      *      1380      *      1400
gi|172000| : CFVBCSKARVGVVDKLESRVGSADLYNENSTFMVEMETSTFIIQGATERSTAILDEIGRGTSGKEGISIAYATLKYLENNQCRTIFATHFGQELKQIIDNKCSKGM : 902
NUM
```

# GeneDoc: Property Mode



# AliView alignment of HBV-S





- Problem:
  - Trichinella spiralis (pork worm)
  - Heron/Human?
  - Root

# Our prediction on HBV surface protein

- A protein family can be found in many viruses found in mammals, bird, frog and even fish
- HBV-S has 4 TM region
- Conserved residues:
  - Conserved hydrophobic Trp (W) for TM
  - a Cys-rich domain between TM2-TM3
- The Cys-rich domain insert event during evolution
- Gene transfers
- Potential function:
  - The insert domain has multiple conserved Cys residues which potentially form disulfide bonds
  - Peptide hormone like function?
  - will contribute to binding host receptor during infection



# ClustalW

CLUSTAL W (1.83) multiple sequence alignment

```

beta globin  -----MVHLTPEEKSAVTALWGKVNVDD--EVGGEALGRLLVVYPWTQRFFESFG- 47
myoglobin   -----MGLSDGEWQLVLNVWVGKVEADIPGHGQEVLIQLFKGHPEETLEKFDKFK- 48
neuroglobin -----MERPEPELIRQSWRAVSRSPLEHGTVLFARLFALPDLLPLFQYNCR 47
soybean      -----MVAFTTEKQDALVSSSFEEAFKANIPQYSVVFYTSILEKAPAAKDLFSFLA- 49
rice         MALVEDNNAVAVSFSEEQEALVLKSWAILKKDSANIALRFFLKIFEVAPSASQMFSFLR- 59
              :   :   :   :   .   .   .   :   *   *
              ▽
beta globin  DLSTPDAVMGNPKVKAHGKKVLGAFSDGLAHLDNLKGTFATLS-----ELHCDKLHVDPE 102
myoglobin   HLKSEDEMKALEDLKKHGATVLTALGGILKKKGHHAEIKPLA-----QSHATKHKIPVK 103
neuroglobin QFSSPEDCLSSPEFLDHIRKVMLVIDAAVTNVEDLSSLEEYLA---LGRKHRAVGVKLS 104
soybean      --NGVDPT--NPKLTGHA EKLFALVRDSAGQLKASGTVVADAA---LGSVHAQKAVTDP 101
rice         --NSDVPLEKNPKLKTHAMSVFVMTCEAAQAQLRKAGKVTVRDTTLKRLGATHLKYGVGDA 117
              .   .   .   *   .::   :   :   :
              :   :   :   :   :   :   :   :
              :   :   :   :   :   :   :   :
beta globin  NFRLLGNVLVLCVLAHFF-GKEFTPPVQAAYQKVVAGVANALAHKYH----- 147
myoglobin   YLEFISECIIQVLQSKH-PGDFGADAQGAMNKALELFRKDMASNYKELGFQG 154
neuroglobin SFSTVGESLLYMLEKCL-GPAFTPATRAAWSQLYGAVVQAMSRGWDGE---- 151
soybean      QFVVVKEALLKTIKAAV-GDKWSDLSRAWEVAYDELA AAIKKA----- 144
rice         HFEVVKFALLDTIKEEVPADMWS PAMKSAWSEAYDHLVAAIKQEMKPAE--- 166
              :   :   ::   :   :   *   .   .   :

```

Note how the region of a conserved histidine (▼) varies depending on which of five prominent algorithms is used

# Praline

(a) Praline multiple sequence alignment

```

beta globin      .....MVHLTPEEKSAVTALWGKV..NVDEVGGEALGRLLVVYPWTQRFFES.FG
myoglobin        .....MGLSDGEWQLVLNVWGKVEADIPGHGQEVLRIRLFKGH PETLEKFDK.FK
neuroglobin      .....MERPEPELIQSWRAVSRSPLEHGTVL FARLFALEPDLLPLFQYNCR
soybean          .....MVAFT EKQDALVSSSF EAFKANIPQYSVVFYTSILEKAPAAKD LFS..FL
rice             MALVEDNNAVAVSFSEEQEALVLKSWAILKKDSANIALRFFLKIFEVAPSASQMFS..FL
Consistency      000000000014265438257934573463364343624453686433*35344*50063

beta globin      DLSTPDAVMGNPKVKAHGKKVLGAFSDG LAHLDNLKGT FATLSEL..HCDKLH....VDP
myoglobin        HLKSEDEMKA SEDLKKHGATVLTALGGILKKKGHHEAEIKPLAQS..HATKHK....IPV
neuroglobin      QFSSPEDCLSSPEFLDHIRKVMLVIDAAVTNVEDLSSLEEYLASLGRKHRAVG....VKL
soybean          A.NGVDP..TNPKLTGHA EKLFALVRDSAGQL.KASGTVVADAA....LGSVHAQKAVTD
rice             R.NSDVPLEKNPKLKTHAMSVFVMTCEAAAQL.RKAGKVTVRDTTLKRLGATHLKYGVGD
Consistency      3166354224776653*4368635424454451335634333542003335440000922

beta globin      ENFRLLGNVLVCVLAHHF.GKEFTPPVQAAYQKV VAGVANALAHKYH.....
myoglobin        KYLEFISECIIQVLQSKH.PGDFGADAQGAMNKALELFRKDMASNYKELGFQG
neuroglobin      SSFSTVGESLLYMLEKCL.GPAFTPATRAAWSQLYGAVVQAMSRGWD..GE..
soybean          PQFVVVKEALLKTIKAAV.GDKWSELSRAWEVAYDELAAAIKKA.....
rice             AHFEVVKFALLDTIKEEVPADMWS PAMKSAWSEAYDHLVAAIKQEMKPAE...
Consistency      43744844498258542305336554454*55465426446754322001000
```

Note also the changing pattern of gaps within the boxed region in these five different alignments.

# MUSCLE

(b) MUSCLE (3.6) multiple sequence alignment

```

beta globin  -----MVHLTPEEKSAVTALWGKVNV--EVGGEALGRLLVVYPWTQRFFES-FG
myoglobin   -----MGLSDGEWQLVLNVWVKVEADIPGHGQEVLIIRLFKGH PETLEKFDK-FK
neuroglobin -----MERPEPELIHQSWRAVSRSPLEHGTVLFARLFALPDLLPLFQYNCR
soybean      -----MVAFTKQDALVSSSF EAFKANIPQYSVVFYTSILEKAPAAKDLFSF-LA
rice         MALVEDNNAVAVSFSEEQEALVLKSWAILKKDSANIALRFFLKIFEVAPSASQMFSF-LR
              :   :   :   :   .   .   .   :   *   *
              ∇                               ▽

beta globin  DLSTPDAVMGNPKVKAHGKKVLGAF---SDG LAHLDNLKGT FATLSELHCDKLH--VDPE
myoglobin   HLKSEDEMKA SEDLKKHGATVLTAL---GGILKKKGHHEAEIKPLAQSHATKHK--IPVK
neuroglobin QFSSPEDCLSSPEFLDHIRKVMLVI---DAAVTNVEDLSSLEEYLASLGRKHRAVGKLS
soybean      NGVDP----TNPKLTGHA EKLFALVRDSAGQLKASGTVVAD----AALGSVHAQKAVTDP
rice         NSDVP--LEKNPKLKTHAMSVFVMTCEAAAQLRKAGKVTVRDTTLKRLGATHLKYGVGDA
              . . . * .::                :      :

beta globin  NFRLLGNVLVLCVLAHFGKE-FTPPVQAAYQKV VAGVANALAHKYH-----
myoglobin   YLEFISECIIQVLQSKHPGD-FGADAQGAMNKALELFRKDMASNYKELGFQG
neuroglobin SFSTVGESLLYMLEKCLGPA-FTPATRAAWSQLYGAVVQAMSRGWDGE----
soybean      QFVVVKEALLKTIKAAVGDK-WSDELSRAW EVAYDELA AAIKKA-----
rice         HFEVVKFALLDTIKEEVPADMWS PAMKSAWSEAYDHLVAAIKQEMKPAE---
              :   :   ::   :               :   *   .   .   :
  
```



# Probcons

(c)  
PROBCONS

beta globin	M-----VHLTPEEKSAVTALWGKVNVD--EVGGEALGRLLVVYPWTQRFFES-FG	
myoglobin	M-----GLSDGEWQLVLNVWGKVEADIPGHGQEVLIIRLFKGH	PETLEKFDK-FK
neuroglobin	M-----ERPEPELIRQSWRAVSRS	PLEHGTVLFARLFALPDLLPLFQYNCR
soybean	M-----VAFTEKQDALVSSSF	EAFKANIPQYSVVFYTSILEKAPAAKDLFSF-LA
rice	MALVEDNNAVAVSFS	EEQEALVLKSWAILKKDSANIALRFFLKI
	* : : : : . . . . : : *	*
beta globin	DLSTPDAVMGNPKVKAHGKKVLGAFSDGLAHLD---NLK---GTFATLSELHCDKLHVDP	
myoglobin	HLKSEDEMKA	SEDLKKHGATVLTALGGI---LKKKGHHE---AEIKPLAQSHATKHKIPV
neuroglobin	QFSSPEDCLSS	PEFLDHIRKVMLVIDAAVTNVEDLSSLE---EYLASLGRKHRAV-GVKL
soybean	NGVDP----	TNPKLTGHAEKLFALVRDSAGQLKASGTVV---ADAALGSVHAQK-AVTD
rice	NSDVP--LEKN	PKLKTTHAMSVFVMTCEAAAQLRKAGKVTVRDTTLKRLGATHLKY-GVGD
	. : . . * . : :	: : . * . *
beta globin	ENFRLLGNVLVCVLAHHF-GKEFTPPVQAAYQKVVAGVANALAHK-----YH	
myoglobin	KYLEFISECIIQVLQSKH-PGDFGADAQGAMNKALELFRKDMASNYKELGFQG	
neuroglobin	SSFSTVGESLLYMLEKCL-GPAFTPATRAAWSQLYGAVVQAMSRG---W-DGE	
soybean	PQFVVVKEALLKTIKAAV-GDKWSELSRAWEVAYDELA	AAIK-----KA
rice	AHFEVVKFALLDTIKKEVPADMWS	PAMKSAWSEAYDHLVAAIKQE---MKPAE
	: : : : : *	. . :

# TCoffee

(d)

CLUSTAL FORMAT for T-COFFEE Version\_5.13

```

beta globin  -----MVHLTPEEKSAVTALWGKVNV--EVGGEALGRLLVVYPWTQRFFESFG
myoglobin   -----MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIIRLFKGHPEKLEKFD-KFK
neuroglobin -----MERPEPELIQSWRAVSRSPLHGTVLFARLFALPDLLPLFQYNCR
soybean      -----MVAFTKQDALVSSSFQAFKANIPQYSVVFYTSILEKAPAAKDLFS-FLA
rice         MALVEDNNAVAVSFS EEQEALVLKSWAILKSDSANIALRFFLKIFEVAPSASQMFS-FLR
              :   :   :   :   . . .   .   : :   *   * .

              ▽                               ▴

beta globin  DLSTPDAVMGNPKVKAHGKKVLGAFSDGLAHLDNL---KGTF---ATLSELHCDKLHVDP
myoglobin   HLKSEDEMKA SEDLKKHGATVLTAL---GGILKKKGHHEAE---IKPLAQSHATKHKIEV
neuroglobin QFSSPEDCLSSPEFLDHIRKVMLVIDAAVTNVEDL---SSLEEYLA SLGRKH-RAVGVKL
soybean     NGVDP----TNPKLTGHA EKL FALVRDSAGQLKASGTVVAD----AALGSVHAQKAVTDIP
rice        NSDVP--LEKNPKLKTHAMSVFVMTCEAAQLRKAGKVTVRD TTKRLGATHLKYGVGDA
              .   . . . * . : :   :   :   * . *

beta globin  ENFRLLGNVLVCVLAHHF-GKEFTPPVQAAYQKVVAGVANALAHKYH-----
myoglobin   KYLEFISECIIQVLQSKH-PGDFGADAQGAMNKALELFRKDMASNYKELGFQG
neuroglobin SSFSTVGESLLYMLEKCL-GPAFTPATRAAWSQLYGAVVQAMSRGWDG----E
soybean     Q-FVVVKEALLKTIKAAV-GDKWSD ELSRAWEVAYDELA AAIKKA-----
rice        H-FEVVKFALLDTIKEEVPA DMWSPAMKSAWSEAYDHLVAAIKQE---MKPAE
              :   :   : :   :   :   * . .   :
  
```

# HMMER

- **HMMER3** is a package to build and use HMMs developed by Sean Eddy (<http://hmmer.wustl.edu/>).
- Software available in HMMER3:
  - hmmbuild to build an HMM from a multiple alignment;
  - hmmscan to align sequences to an HMM model;
  - hmmsearch to search a sequence database with an HMM model;
  - jackhmmer to iteratively search sequence(s) against a protein database;
  - hmmscan to search protein sequence(s) against a protein profile database;
  - hmmeemit to get sample sequences from a profile HMM;
  - hmmsfetch to retrieve profile HMM(s) from a file

# HMMER installation

- Visit <http://linuxbrew.sh>
- Follow the install instructions
- `brew tap homebrew/science`
- `brew update`
- `brew install hmmer`

# Pfam download

- `mkdir ~/data`
- `cd ~/data`
- `wget -c -t0 ftp://ftp.ebi.ac.uk/pub/databases/Pfam/current\_release/Pfam-A.hmm.gz`
- `gunzip Pfam-A.hmm.gz`
- `hmmcompress Pfam-A.hmm`

- `cat RHS.fasta | muscle > RHS.muscle`
- `more RHS.muscle`
- Copy the result to a text file
- Open it using any MSA viewer

Let make hmm profile using this alignment:

- `hmmbuild RHS.hmm RHS.muscle`

Search Ecoli genome using this HMM model:

- `hmmsearch -E 1 RHS.hmm Ecoli.O157.fa > RHS.Ecoli.hmm.result`
- Using more to see the result

Get the hit id and fasta sequences

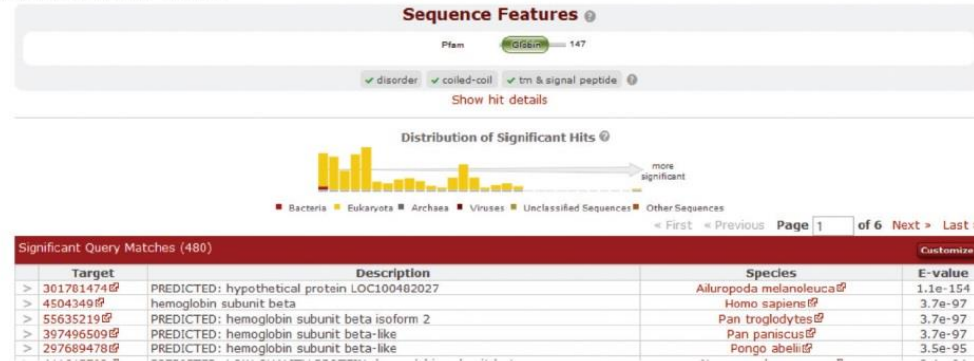
- `hmmsearch -E 1 RHS.hmm Ecoli.O157.fa | grep '>>' | cut -f 2 -d ' ' | epost -db protein -format acc | efetch -format fasta > RHS.Ecoli.fasta`

Identifying pfam domains in these sequences

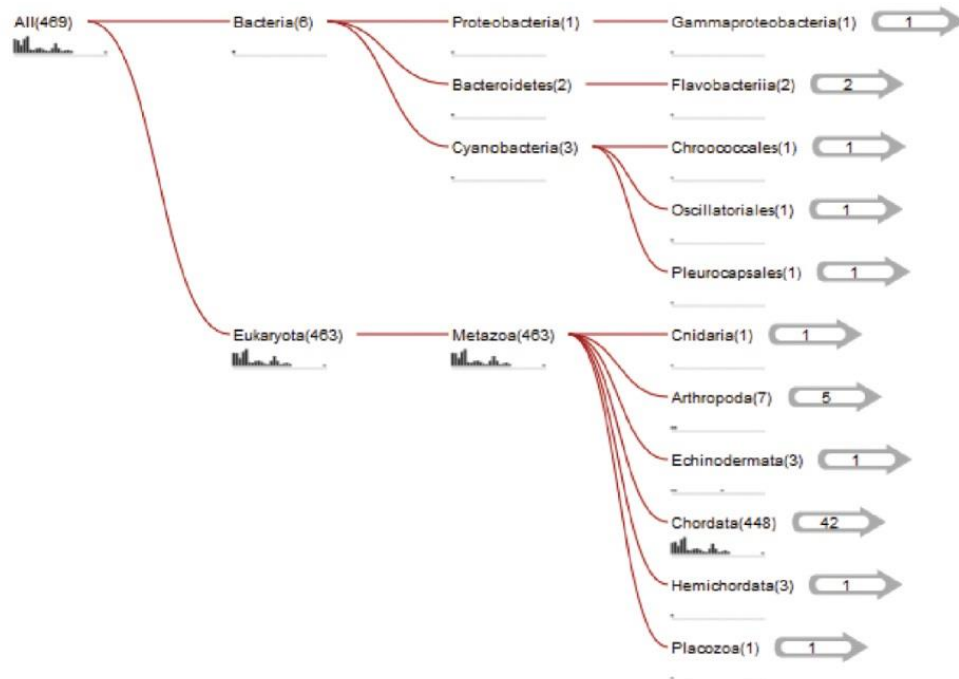
- `hmmsearch ~/data/Pfam-A.hmm RHS.Ecoli.fasta > RHS.Ecoli.pfam.result`

# HMMER is available online

(a) HMMER web output



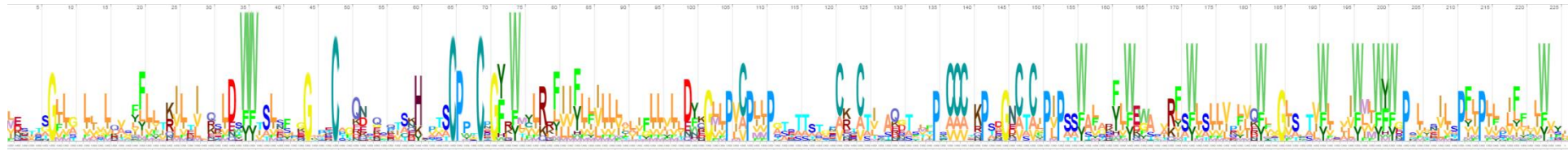
(b) HMMER phylogenetic output



- Domain architecture
  - Taxonomy
- Iterative manner

# Surface protein (HBsAg)

- HMMER searches:
  - HBV S protein
  - Ground Squirrel hepatitis virus S protein
  - Duck hepatitis B virus S protein





# HMMER software: build profiles, complement BLAST

Build a profile HMM (input is a multiple sequence alignment)

```
$ ./hmmbuild -h # provides brief help documentation
$ ./hmmbuild globins4.hmm ../tutorial/globins4.sto
```

Download a database to search (e.g. human RefSeq proteins)

```
$ wget ftp://ftp.ncbi.nlm.nih.gov/refseq/H_sapiens/mRNA_Prot/human.protein
.faa.gz
$ gunzip human.protein.faa.gz
$ wc -l human.protein.faa
302761 human.protein.faa
```

Search an HMM against a database

```
$ ./hmmsearch globins4.hmm human.protein.faa > globins4.out
```

# HMMER results

```
# hmmsearch :: search profile(s) against a sequence database
# HMMER 3.1b1 (May 2013); http://hmmerr.org/
# Copyright (C) 2013 Howard Hughes Medical Institute.
# Freely distributed under the GNU General Public License (GPLv3).
# - - - - -
# query HMM file:                globins4.hmm
# target sequence database:      /mnt/reference/human.protein.faa
# - - - - -

Query:        globins4  [M=149]
Scores for complete sequences (score includes all domains):
--- full sequence ---
   E-value   score   bias    Sequence                        Description
   -----
   3.3e-64   216.6    0.0    ref|NP_000509.1|      hemoglobin subunit beta [Homo sa
       7e-61   205.8    0.0    ref|NP_000510.1|      hemoglobin subunit delta [Homo s
   2.3e-60   204.2    1.3    ref|NP_000508.1|      hemoglobin subunit alpha [Homo s
   2.3e-60   204.2    1.3    ref|NP_000549.1|      hemoglobin subunit alpha [Homo s
   6.2e-60   202.8    0.3    ref|NP_976311.1|      myoglobin [Homo sapiens]
   6.2e-60   202.8    0.3    ref|NP_976312.1|      myoglobin [Homo sapiens]
   6.2e-60   202.8    0.3    ref|NP_005359.1|      myoglobin [Homo sapiens]
   4.8e-55   186.9    0.0    ref|NP_000175.1|      hemoglobin subunit gamma-2 [Homo
   1.4e-54   185.4    0.4    ref|NP_005321.1|      hemoglobin subunit epsilon [Homo
   2.1e-54   184.8    0.1    ref|NP_000550.2|      hemoglobin subunit gamma-1 [Homo
   4.9e-48   164.2    0.2    ref|NP_005323.1|      hemoglobin subunit zeta [Homo sa
   1.7e-40   139.7    0.1    ref|NP_005322.1|      hemoglobin subunit theta-1 [Homo
   1.8e-39   136.4    0.2    ref|NP_599030.1|      cytoglobin [Homo sapiens]
       5e-35   121.9    0.3    ref|NP_001003938.1|    hemoglobin subunit mu [Homo sapi
       3e-08    35.0    0.0    ref|NP_067080.1|      neuroglobin [Homo sapiens]
----- inclusion threshold -----
       0.14    13.4    0.0    ref|NP_001371.1|      dedicator of cytokinesis protein
       0.25    12.6    0.8    ref|NP_006737.2|      sex comb on midleg-like protein
       0.28    12.4    0.8    ref|NP_001032629.1|    sex comb on midleg-like protein
```

HMMER output includes scores, E values