

# Genome Annotation

## MODULE 2 PART -II

**Rama Dey-Rao, PhD**

Clinical Assistant Professor

Biotechnical and Clinical Lab Sciences

Senior Scientist

Department of Microbiology & Immunology, SUNY at Buffalo

[dey@buffalo.edu](mailto:dey@buffalo.edu)

# Sequence-based Similarity

## 4 TOOLS

### 1. BLAST

The Basic Local Alignment Search Tool (**BLAST**) tool finds regions of local similarity between sequences and calculates the statistical significance of matches

### 2. CDD

Conserved Domain Database Search (**CDD**) finds sequence similarity with genes in conserved orthologous groups (COGs).

### 3. T-Coffee

Tree based Consistency Objective Function For alignment Evaluation (**T-Coffee**) is a multiple sequence alignment program that aligns a set of homologous (similar) sequences

### 4. WebLogo

WebLogo is a program that enables easy creation of sequence logos from the multiple sequence alignments

# Conserved Domain Database Search

**CDD**

click on the CDD search results at the top of the BLAST results page

COG number (*top hit*)

COG name

Score

E-value

Significant COG number (*second hit*)

COG name

Score

E-value

**T-Coffee**

go to <http://www.ebi.ac.uk/Tools/msa/tcoffee/>

Sequences used for alignment

Multiple sequence alignment

**WebLogo**

go to <http://weblogo.berkeley.edu/>

Sequence logo

# Protein Domains

## COG (Clusters of Orthologous Groups ) Conserved Domain Database

- Sequence –based domains in proteins have a particular structure that is related to function.
- Can be seen as building blocks put together in different ways in different proteins.
- Parts of the proteins with similar and vital functions are conserved –clusters of ortholog groups in conserved domain database.
- When a very significant COG hit is observed for the query gene it can be interpreted as a strong likelihood that the protein has the same function.

# BLAST RESULTS PAGE- Swissprot database

Both a Conserved Domain Database  
and BLAST searches are done simultaneously.

RID [631JYCW7016](#) (Expires on 02-12 21:27 pm)

Query ID |d|Query\_44273  
Description KSED\_RS00005- Ksed\_00010-aa sequence  
Molecule type amino acid  
Query Length 506

Database Name swissprot  
Description Non-redundant UniProtKB/SwissProt sequences  
Program BLASTP 2.8.1+ [Citation](#)

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.

**CDD search (conserved domain database)**

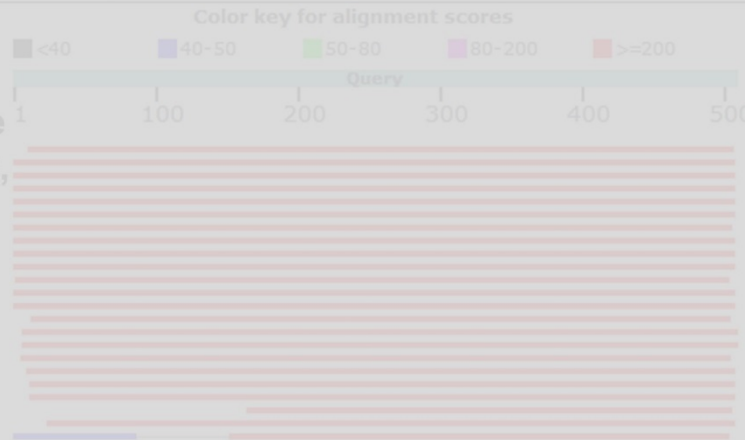
Query seq. Specific hits Superfamilies

dnaA  
dnaA superfamily

Distribution of the top 106 Blast Hits on 100 subject sequences

**Double Click anywhere**

A high score and close to 100% coverage would indicate a high quality alignment, suggesting this sequence is highly conserved in a number of different organisms.



Quick scan  
Visual rep  
of coverage

Done

[Questions/comment](#)

# COG – Clusters of Orthologous Groups

NCBI

HOME SEARCH GUIDE NewSearch Structure Home 3D Macromolecular Structures Conserved Domains Pubchem BioSystems

Conserved domains on [gi|256687299|gb|ACV05101|] View Standard Results

chromosomal replication initiator protein DnaA [Kytococcus sedentarius DSM 20547]

**Protein Classification**

chromosomal replication initiator protein DnaA (domain architecture ID 11478209)  
 chromosomal replication initiator protein DnaA plays a key role in the initiation and regulation of chromosomal replication

**Graphical summary**  Zoom to residue level show extra options

Query seq. 1 75 150 225 300 375 450 506

Specific hits

Superfamilies

Click [+]

Search for similar domain architectures Refine search

**List of domain hits**

	Name	Accession	Description	Interval	E-value
[+]	dnaA	PRK00149	chromosomal replication initiation protein; Reviewed	157-500	0e+00
[+]	DnaA	TIGR00362	chromosomal replication initiator protein DnaA; DnaA is involved in DNA biosynthesis; ...	33-498	0e+00
[+]	DnaA	COG0593	Chromosomal replication initiation ATPase DnaA [Replication, recombination and repair];	147-496	2.48e-153
[+]	Bac_DnaA	pfam00308	Bacterial dnaA protein;	164-381	2.73e-115
[+]	Bac_DnaA_C	cd06571	C-terminal domain of bacterial DnaA proteins. The DNA-binding C-terminal domain of DnaA ...	409-498	1.32e-37
[+]	Bac_DnaA_C	smart00760	Bacterial dnaA protein helix-turn-helix domain; Could be involved in DNA-binding.	408-476	3.34e-30

COG hit

### List of domain hits

Name	Accession	Description	Interval	E-value
[+] dnaA	PRK00149	chromosomal replication initiation protein; Reviewed	157-500	0e+00
[+] DnaA	TIGR00362	chromosomal replication initiator protein DnaA; DnaA is involved in DNA biosynthesis; ...	33-498	0e+00
-] DnaA	<b>COG0593</b>	Chromosomal replication initiation ATPase DnaA [Replication, recombination and repair];	147-496	2.48e-153

Chromosomal replication initiation ATPase DnaA [Replication, recombination and repair];

COG name

Pssm-ID: 223666 Cd Length: 408 Bit Score: 447.50 E-value: 2.48e-153

Length, bit score, and E-value

```

      10      20      30      40      50      60      70      80
...*...|...*...|...*...|...*...|...*...|...*...|...*...|...*...|
gi 256687299 147 SLTATNSSPGVERDYSALNHKTYFDTFVLGSSNRFAHAAATAVAEAPARAYNPLFIYGGSGLGKTHLLHAIGHYARTLDS 226
Cdd:COG0593  63 KVEVRASAPAQPLPSGLNPKYTFDNEFVVGPSNRLAYAAAKVAENPGGAYNPLFIYGGVGLGKTHLLQAIGNEALANGP 142

      90     100     110     120     130     140     150     160
...*...|...*...|...*...|...*...|...*...|...*...|...*...|...*...|
gi 256687299 227 SVRVKYVNSEEFINQFINAVSAGQANAFQRQYrDVDVLLIDDIQFLQKEQTMEEFFHTFNILHNSEKQIVITSDQPPKK 306
Cdd:COG0593 143 NARVVYLISEDFTNDFVKALRDNEMEKFKKEY-SLDLLIDDIQFLAGKERTQEEFFHTFNALLENGKQIVLTSDRPPKE 221

     170     180     190     200     210     220     230     240
...*...|...*...|...*...|...*...|...*...|...*...|...*...|...*...|
gi 256687299 307 LSGFAERMRSRFEWGLLTDVQPPDLETRIAILRRKAAADKLDIPDDVLHLIASKISSNIRELEGALTRVTAFAFASLSGSPL 386
Cdd:COG0593 222 LNGLEDRLRSRLEWGLVVEIEPPDDETRLAAILRKAEDRGIEIPDEVLEFLAKRLDRNVRELEGALNRLDAFALFTKRAI 301

     250     260     270     280     290     300     310     320
...*...|...*...|...*...|...*...|...*...|...*...|...*...|...*...|
```

If there are no hits, write “no significant hits” in notebook

If there **are several** hits, click the **[+]** sign next to the hits and record

# COG – Clusters of Orthologous Groups

NCBI

HOME SEARCH GUIDE NewSearch Structure Home 3D Macromolecular Structures Conserved Domains Pubchem BioSystems

Conserved domains on [gi|256687299|gb|ACV05101|] View Standard Results

chromosomal replication initiator protein DnaA [Kytococcus sedentarius DSM 20547]

**Protein Classification**

chromosomal replication initiator protein DnaA (domain architecture ID 11478209)  
 chromosomal replication initiator protein DnaA plays a key role in the initiation and regulation of chromosomal replication

**Graphical summary**  Zoom to residue level show extra options

Query seq. 1 75 150 225 300 375 450 506

Specific hits

Superfamilies

Click

Search for similar domain architectures Refine search

**List of domain hits**

	Name	Accession	Description	Interval	E-value
[+]	dnaA	PRK00149	chromosomal replication initiation protein; Reviewed	157-500	0e+00
[+]	DnaA	TIGR01462	chromosomal replication initiator protein DnaA; DnaA is involved in DNA biosynthesis; ...	33-498	0e+00
[+]	DnaA	COG0593	Chromosomal replication initiation ATPase DnaA [Replication, recombination and repair];	147-496	2.48e-153
[+]	Bac_DnaA	pfam00308	Bacterial dnaA protein;	164-381	2.73e-115
[+]	Bac_DnaA_C	cd06571	C-terminal domain of bacterial DnaA proteins. The DNA-binding C-terminal domain of DnaA ...	409-498	1.32e-37
[+]	Bac_DnaA_C	smart00760	Bacterial dnaA protein helix-turn-helix domain; Could be involved in DNA-binding.	408-476	3.34e-30



# COG – Clusters of Orthologous Groups

NCBI

## Conserved Protein Domain Family *DnaA*

HOME SEARCH SITE MAP Entrez CDD Structure Protein Help

**COG0593: DnaA** ?

Chromosomal replication initiation ATPase DnaA [Replication, recombination and repair]

**Links** ?

- Source: [COG](#)
- Taxonomy: [Bacteria](#)
- Protein: [Representatives](#)  
[Specific Protein](#)  
[Related Protein](#)  
[Related Structure](#)  
[Architectures](#)
- Superfamily: [d07055](#)

**Statistics** ?

**Structure** ?

COG0593 is a member of the superfamily [d07055](#).


Click

NCBI

## Conserved Protein Domain Family *Bac\_DnaA\_C*

HOME SEARCH SITE MAP Entrez CDD Structure Protein Help

**d07055: Bac\_DnaA\_C Superfamily** ?



C-terminal domain of bacterial DnaA proteins. The DNA-binding C-terminal domain of DnaA contains a helix-turn-helix motif that specifically interacts with the DnaA box, a 9-mer motif that occurs repetitively in the replication origin *oriC*. Multiple copies of DnaA, which is an ATPase, bind to 9-mers at the origin and form an initial complex in which the DNA strands are being separated in an ATP-dependent step.

# Sequence-based Similarity Data Module

## 4 TOOLS

### 1. BLAST

The Basic Local Alignment Search Tool (**BLAST**) finds regions of local similarity between sequences and calculates the statistical significance of matches

### 2. CDD

Conserved Domain Database Search (**CDD**) finds sequence similarity with genes in conserved orthologous groups (COGs).

### 3. T-Coffee

Tree based **C**onsistency **O**bjective **F**unction **F**or alignment **E**valuation (**T-Coffee**) is a multiple sequence alignment program that aligns a set of homologous (similar ) sequences

### 4. WebLogo

WebLogo is a program that enables easy creation of sequence logos from the multiple sequence alignments

# T-Coffee

Tree-based Consistency Objective Function for alignment Evaluation  
**Multiple sequence alignment tool**

- Across evolution amino acids are likely to be conserved because they are important for structure and function.
- One way to measure conservation is to align multiple similar protein sequences from related organisms (orthologs)

*T-Coffee: A novel method for multiple sequence alignments.* Notredame, Higgins, Heringa, **JMB** 302(205-217)2000

T-Coffee is a freeware open source package distributed under the [GNU public license](#)

T-Coffee Server is hosted by the [Centre for Genomic Regulation](#) (CRG) of Barcelona, SPAIN

# In Notebook

## T-Coffee

go to <http://www.ebi.ac.uk/Tools/msa/tcoffee/>

Sequences used for alignment



sequences

Multiple sequence alignment

alignment

# RECALL: What are orthologs?

- Homologs

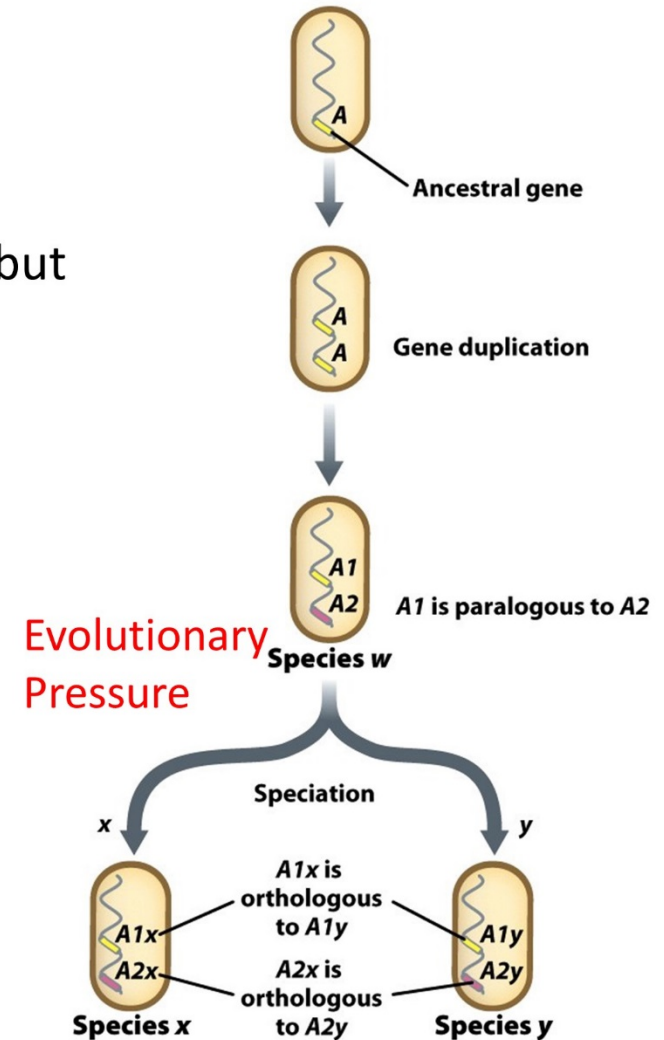
- Orthologs

- Genes that share similarity in function but are found in different organisms

- Paralogs

- Genes duplicated within a species

- Perform slightly different tasks in cell
        - » Can develop new capabilities
        - » Can become pseudogene if functionality lost but sequence similarity retained



# How and where are the orthologs to compare?

## Go back to BLAST Search from nr database

Sequences producing significant alignments:

Select: [All](#) [None](#) Selected: 10

[Alignments](#) [Download](#) [GenPept](#) [Graphics](#) [Distance tree of results](#) [Multiple alignment](#)

	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input checked="" type="checkbox"/>	<a href="#">chromosomal replication initiation protein [Kytococcus sedentarius]</a>	1031	1031	99%	0.0	99%	<a href="#">WP_012801520.1</a>
<input checked="" type="checkbox"/>	<a href="#">chromosomal replication initiation protein [Ornithinimicrobium pekingense]</a>	610	610	99%	0.0	63%	<a href="#">WP_022920049.1</a>
<input checked="" type="checkbox"/>	<a href="#">chromosomal replication initiation protein [Serinicoccus profundii]</a>	589	589	98%	0.0	60%	<a href="#">WP_010147278.1</a>
<input type="checkbox"/>	<a href="#">chromosomal replication initiation protein [Serinicoccus marinus]</a>	565	565	97%	0.0	59%	<a href="#">WP_022923463.1</a>
<input checked="" type="checkbox"/>	<a href="#">chromosomal replication initiation protein [Janibacter sp. HTCC2649]</a>	552	552	96%	0.0	58%	<a href="#">WP_009776970.1</a>
<input checked="" type="checkbox"/>	<a href="#">hypothetical protein [Arsenicococcus bolidensis]</a>	539	539	96%	0.0	56%	<a href="#">WP_029212190.1</a>
<input checked="" type="checkbox"/>	<a href="#">chromosomal replication initiator protein DnaA-DNA-binding transcriptional dual regulator [Tetrasphaera elongata]</a>	537	537	98%	0.0	55%	<a href="#">WP_010851794.1</a>
<input checked="" type="checkbox"/>	<a href="#">chromosomal replication initiation protein [Cellvibrio gilvus]</a>	536	536	98%	0.0	58%	<a href="#">WP_013882065.1</a>
<input checked="" type="checkbox"/>	<a href="#">chromosomal replication initiation protein [Mobilicoccus pelagius]</a>	531	531	97%	0.0	59%	<a href="#">WP_009482734.1</a>
<input checked="" type="checkbox"/>	<a href="#">chromosomal replication initiation protein [Paraoskovia marina]</a>	528	528	97%	2e-180	57%	<a href="#">WP_029253865.1</a>
<input checked="" type="checkbox"/>	<a href="#">chromosomal replication initiation protein [Sanquibacter keddieii]</a>	526	526	98%	7e-180	56%	<a href="#">WP_012865049.1</a>

Click on 10 different significant orthologs from the list

Do NOT use the top 10 without checking out if they are from different organisms if available  
Sometimes the same organism appear multiple times (different strains of the same organism)

Remember: you may need to do an “exclusion blast” if you don’t find enough different organisms in your routine nr blast search!

# BLAST Search from nr database- Why?

The screenshot shows the BLAST search results interface. At the top, it says "Sequences producing significant alignments:" and "Select: All None Selected: 10". Below this is a navigation bar with "Alignments", "Download", "GenPept", "Graphics", "Distance tree of results", and "Multiple alignment". A dropdown menu is open under "Download", showing options: FASTA (complete sequence), FASTA (aligned sequences), GenBank (complete sequence), Hit Table (text), Hit Table (CSV), Text, XML, and ASN.1. A red arrow points from the "Continue" button in the dropdown menu to the "FASTA (complete sequence)" option. Below the menu is a table of search results with columns: Description, Max score, Total score, Query cover, E value, Ident, and Accession. The first row is highlighted in blue and corresponds to the top of the dropdown menu.

Description	Max score	Total score	Query cover	E value	Ident	Accession
...tion protein [Kytococcus sedentarius]	1031	1031	99%	0.0	99%	<a href="#">WP_012801520.1</a>
...tion protein [Ornithinimicrobium pekingense]	610	610	99%	0.0	63%	<a href="#">WP_022920049.1</a>
...tion protein [Serinicoccus profundii]	589	589	98%	0.0	60%	<a href="#">WP_010147278.1</a>
...tion protein [Serinicoccus marinus]	565	565	97%	0.0	59%	<a href="#">WP_022923463.1</a>
...tion protein [Janibacter sp. HTCC2649]	552	552	96%	0.0	58%	<a href="#">WP_009776970.1</a>
...ccus bolidensis]	539	539	96%	0.0	56%	<a href="#">WP_029212190.1</a>
...tor protein DnaA/DNA-binding transcriptional dual regulator [Tetrasphaera elongata]	537	537	98%	0.0	55%	<a href="#">WP_010851794.1</a>
...tion protein I [[Cellvibrio] qilvus]	536	536	98%	0.0	58%	<a href="#">WP_013882065.1</a>
<input checked="" type="checkbox"/> chromosomal replication initiation protein [Mobilicoccus pelagius]	531	531	97%	0.0	59%	<a href="#">WP_009482734.1</a>
<input checked="" type="checkbox"/> chromosomal replication initiation protein [Paraoerskovia marina]	528	528	97%	2e-180	57%	<a href="#">WP_029253865.1</a>
<input checked="" type="checkbox"/> chromosomal replication initiation protein [Sanquibacter keddjeii]	526	526	98%	7e-180	56%	<a href="#">WP_012865049.1</a>
<input type="checkbox"/> chromosomal replication initiation protein [Dermatophilus congolensis]	526	526	97%	2e-179	57%	<a href="#">WP_028327210.1</a>
<input type="checkbox"/> chromosomal replication initiation protein [Kineosphaera limosa]	526	526	97%	5e-179	58%	<a href="#">WP_006591943.1</a>
<input type="checkbox"/> chromosomal replication initiation protein [Cellulomonas flavigena]	526	526	97%	5e-179	56%	<a href="#">WP_013115255.1</a>
<input type="checkbox"/> chromosomal replication initiation protein [Actinopolymorpha alba]	524	524	91%	7e-179	56%	<a href="#">WP_026257010.1</a>
<input type="checkbox"/> chromosomal replication initiation protein [Ruania albidiflava]	523	523	96%	1e-178	57%	<a href="#">WP_022917303.1</a>
<input type="checkbox"/> hypothetical protein [Demetria terraqena]	522	522	96%	1e-178	59%	<a href="#">WP_018157546.1</a>

Click the Download pull down menu at the top of the page and make sure the FASTA (complete sequence) link is clicked –Continue  
**Copy and paste all 10 sequences** including *Kytococcus sedentarius* in a word document.  
Change font to Courier New 10 if needed. Edit *Kytococcus* FASTA header

# Paste all sequences in your notebook and EDIT



```
>WP_012801520.1 chromosomal replication initiator protein DnaA [Kytococcus sedentarius]
MSQTPDDHATAIWQEAMVHLQGAGLAPRDIQVLRRLATLVGLLEGTALLAVKYDHVKDAVEGHLREDVSTALAEVLDRDIR
LAVSVDPDAVSAAQEEAAPPAPSPAEDDDPATGEGPLSTAVDQAVEKHEGSSPARAGESVAPATTASLTATNSSPGVERD
YSALNHKYTFDTFVLGSSNRFAHAATAVAEAPARAYNPLFIYGSGLGKTHLLHAIGHYARTLDSSVRVKYVNSEFTN
QFINAVSAGQANAFQYRDVVDVLLIDDIQFLQGKEQTMEEFFHTFNTLHNSEKQIVITSDQPPKKLSGFAERMRSRFEW
GLLTDVQPPDLETRIAILRRKAAADKLDI...
GGDSGQITPTMILEETAGYFVISVEEIQ...
LGEDRRVYDEVSELTSIIRKKAARGR
>WP_022920049.1 chromosomal replication initiator protein DnaA [Kytococcus sedentarius]
pekingense]
MTSQSPAESAQVWQRVVSQLESQGVTDRAFLRLTQLVGLLDTTALLAVPYQHTKETLETTLRQPIVDALAGELGHDVR
LAITVDEDLRRQVEDEGDPAPGPAVTEQVPSDPDRTPYRSNGAGPGEPRSDGHRTPSGAVQTASAEDARLNPKYTFDTFV
SGSSNRFAHAASLAVAESPARAYNPLFIYGESGLGKTHLLHAIGHYARSLYPGVRVRYVNSEEFNTDFINSIRDDKAGAF
QRRYRNVDVFLVDDIQFLQGKEQTVVEEFFHTFNTLHNSEKQVVITSDQPPKRLSGFAERMRSRFEWGLLTDVQPPDLETR
IAILKKAQEGMQLPDEVLELIGSKI STNIRELEGALIRVTAFAASLSSTPPDAALASHVLKDIIPNSESAITVPTIMG
EVADYFQISNDDLCGTSRSRTL VNARQIAMYLCRELTDLSLPKIGQEFGGRDHTTVMHAERKIRQLIGERRALYDQITEL
TGIIRKASAR
>RIK14929.1 chromosomal replication initiator protein DnaA [Acidobacteria bacterium]
MTHDPSPAASAEVWERVVVAELDQGVTDRAFLRLTQLVGLLDTTALLAVPYQHTKDTLETTLRQPIVSALAEELGHDVR
LAITVDESLRQELKAEAGAVTPPQVAPAGGSTPYPVEVEPTSVPPVAEPTPRRAGATQGTGPDEARLNPKYSFDTFVSGS
SNRFAHAASLAVAESPARAYNPLFIYGESGLGKTHLLHAIGHYARSLYPGVRVRYVNSEEFNTDFINSIRNQEAGAFQRR
YRNVDVFLVDDIQFMQKKEQTVVEEFFHTFNTLHNSEKQVVITSDQPPKRLSGFAERMRSRFEWGLLTDVQPPDLETRIAI
LRKKAQEGMNTPDEVLELIASRITTNIRELEGALIRVTAFAASLSSEPLTAEAAHVLDKDIIPSGEAAAIGVPTIIAEVS
DYFQITRDELTCGTSRSRSLVNARQIAMYLCRELTELSPKIGQEFGGRDHTTVMHAERKIRQLMGERRALYDQITDLTGI
IRKASAR
>WP_010147278.1 chromosomal replication initiator protein DnaA [Serinicoccus profundus]
MSQPSTDSGDTWRRVVSELEDKGLGAREKAFRLRTTMVGVLDSTVLLAVPYPHTKEMLETTLRQPIVDLLSRELDREVRL
AITVDDDVQRQVEDEADDEADEDAQTRESLTPASQPSSSAGAGVPGPSGNGIIPRPATPAGPAVTGAADARLNPKYSFD
TFVSGPSNRFAHAASLAVAESPARAYNPLFIYGESGLGKTHLLHAIGHYARKLYPGVRVRYVNSEEFNTDFINSIRDDKA
GAFQRRYRNVDVFLVDDIQFLQGKEQTVVEEFFHTFNTLHNSEKQVVITSDQPPKRLSGFAERMRSRFEWGLLTDVQPPDL
ETRIAILRRKKAQEGMQLPDEVLELIASRITTNIRELEGALIRVTAFAASLSQPADADLAHVLDKDIVPGSDTAQITVST
IIREVSEYFQISIDELCGTSRSRTL VNARQIAMYLCRELTDLSLPKIGQEFGGRDHTTVMHAERKIRAQIGERRALYDQI
AELTGTIRRASQR
```

Select all letters before the word *Kytococcus* and delete



# The headers should look like this



```
>[Kytococcus sedentarius]  
MSQTPDDHATAIWQEAMVHLQAGLAPRDIGVLRRLATLVGLLEG TALLAVKYDHVKDAVEGHLREDVSTALAEVLD RDIR  
LAVSVDPDAVSAAQEEAAPAPSPADEDDPATGEGPLSTAVDGA VEKHEGSSPARAGESVAPATTASLTATNSSPGVERD  
YSALNHKYTFDFTFVLGSSNRF AHAATAVAEAPARAYNPLFIYGG SGLGKTHLLHAIGHYARTLDSSVRVKYVNSEFTN  
QFINAVSAGQANAFQRQYRDVDVLLIDDIQFLQGKEQTMEEF FHTFNTLHNSEKQIVITSDQPPKKLSGFAERMRSRFEW  
GLLTDVQPPDLETRIAILRRKAAADKLDIPDDVLHLIASKI SSNIRELEGALTRVTAFAFASLSGSPLDEYLARTVLKDVMP  
GGDSGQITPTMILEETAGYFVIVSVEEIQGASRSRNLTRAR QIAMYLCRELTDLSLPGKIGKEFGGRDHTTVMHAERKIKQL  
LGEDRRVYDEVSELSIIRKKAARGR  
>[Qrnithinimicrobium pekingense]  
MISQSPAESA EVWQRVVSQLESQGVITARDRAFLRLTQLVGL LDTTALLAVPYQHTKETLETTLRQPIVDALAGELGHDVR  
LAITVDEDLRRQVEDEGDPAPGPAVTEQVPSDPDRTPYRS NAGAGPGEPRSDGHRTPSGAVQTASAEDARLNPKYTFDFTFV  
SGSSNRF AHAASLAVAESPARAYNPLFIYGESGLGKTHLLH AIGHYARSLYPGVRVRYVNSEEFNTDFINSIRDDKAGAF  
QRRYRNVD FLLVDDIQFLQGKEQTVEEFHTFNTLHNSEKQ VVITSDQPPKRLSGFAERMRSRFEWGLLTDVQPPDLETR  
IAILKKA AQEGMQLPDEVLELIGSKI STNIRELEGALIRV TAFASLSSTPPDAALASHVLKDIIPNSESAAITVPTIMA  
EVADYFQIS NDDLCGTSRSRTL VNARQIAMYLCRELTDLS LPGKIGQEFGGRDHTTVMHAERKIRQLIGERRALYDQITEL  
TGIIRKASAR  
>[Acidobacteria bacterium]  
MTHDPSPAESA EVWERVVAELDQGVITARDRAFLRLTRLVGL LDGTVLLAVPYQHTKDTLETTLRQPIVSALAEELGHDVR  
LAITVDES LRQELKAEAGAVTPPQVAPAGGSTPYPVEVEPT SVPPVAEPTPRRAGATQGTGPDEARLNPKYSFDTFVSGS  
SNRF AHAASLAVAESPARAYNPLFIYGESGLGKTHLLH AIGHYARSLYPGVRVRYVNSEEFNTDFINSIRNQEAGAFQRR  
YRNVD FLLIDDIQFMQGKEQTVEEFHTFNTLHNSEKQ VVITSDQPPKRLSGFAERMRSRFEWGLLTDVQPPDLETRIAI  
LRKKA AQEGMNTPDEVLELIASRITTNIRELEGALIRV TAFASLSSEPLTAE LAHV LKDIIPSGEAAAIGVPTIIAEVS  
DYFQITRDEL CGTSRSRSLVNARQIAMYLCRELTEL SLPKIGQEFGGRDHTTVMHAERKIRQLMGERRALYDQITDLTGI  
IRKASAR  
>[Serinicoccus profundus]  
MSQPSTDSGDTWRRV VSELEDKGLGAREKAFRLRTTMVGVLD STVLLAVPYPHTKEMLETTLRQPIVDLLSRELDREVRL  
AITVDDVQRVVEDEADDEADEDAQTRESLTPASQPSSSAGAGV PGPSPGNGIPRPATPAGPAVTGAAD EARLNPKYSFD  
TFVSGPSNRF AHAASLAVAESPARAYNPLFIYGESGLG KTHLLHAIGHYARKLYPGVRVRYVNSEEFNTDFINSIRDDKA  
GAFQRRYRNVD FLLVDDIQFLQGKEQTVEEFHTFNTL HNSEKQVVITSDQPPKRLSGFAERMRSRFEWGLLTDVQPPDL  
ETRIAILR KKA AQEGMQLPDEVLEH IASRITTNIRELEGALIRV TAFASLS SQPADADLAHV LKDIVPGSDTAQITVST  
IIREVSEYFQIS IDELCGTSRSRTL VNARQIAMYLCREL TDLSLPGKIGQEFGGRDHTTVMHAERKIRAQIGERRALYDQI  
AELTGTIRRASQR
```

# Click on Link: T-Coffee tools

<http://www.ebi.ac.uk/Tools/msa/tcoffee/>

**T-Coffee**

Input form | Web services | Help & Documentation | Bioinformatics Tools FAQ | Feedback | Share

Tools > Multiple Sequence Alignment > T-Coffee

## Multiple Sequence Alignment

T-Coffee is a multiple sequence alignment program. Its main characteristic is that it will allow you to compare

**Important note:** This tool can align up to 500 sequences or a maximum file size of 1 MB.

**STEP 1 - Enter your input sequences**

Enter or paste a set of

PROTEIN

sequences in any supported format:

```
>Kytococcus sedentarius]
MSQTPDDHATAIWQEAMVHLQGAGLAPRDIGVLRLATLVGLLEGTALLAVKYDHYKDAVEGHLREDVSTALAEVLRDRDIR
LAVSVDPDAVSAAQEEAAPPAPSPAEDDDPATGEGPLSTAVDGAVEKHEGSSPARAGESVAPATTASLTATNSSPGVERD
YSALNHKYTFDTFVLGSSNRFHAAATAVAEAPARAYNPLFIYGGSSGLGKTHLLHAIGHYARTLDSSVRVKYVNSSEETIN
QFINAVSAGQANAFQRQYRDVDVLLIDDIQFLQGKEQTMEFFHTFNTLHNSEKQIVITSDQPPKLSGFAERMRSREFEW
GLLTDVQPPDLETRIAILRRKAAADKLDIPDDVLHUASKISSNIRELEGALTRVAFASLGGSPIDEYLARTVLDKVMF
GGDSGQITPTMILEETAGYFVISVEEIQGASRSRNLTRARQIAMYLCRELTDLSLPKIGKEFGGRDHTVMHAERKIKQL
LGEDRRVYDEVSELTSIIRKKAARGR
```

Or upload a file:  No file selected. [Use an example sequence](#) | [Clear sequence](#) | [See more example inputs](#)

**STEP 2 - Set your Parameters**

OUTPUT FORMAT:

ClustalW

The default settings will fulfill the needs of most users.

(Click here, if you want to)

**STEP 3 - Submit your job**

Be notified by email (Tick this box if you want to be notified by email when the results are available)

**Copy / paste all the 10 amino acid sequences in FASTA format into the query window.**

**CLICK SUBMIT**

# Multiple Sequence Alignment (msa)

Copy / paste this alignment into your lab notebook-  
Courier New 9

Results for job tcoffee-I20190215-165604-0351-23150635-p1m

Alignments

Result Summary

Guide Tree

Phylogenetic Tree

Submission Detail

Download Alignment File

Show Colors

View result with Jalview

Send to Simple

CLUSTAL W (1.83) multiple sequence alignment

```
[Acidobacteria      MTHDPSPAASAEVNERVVVELD--QGVTARDRAFLRLRLVGLLDGTVLL
[Arsenicicoccus    MSQ-PST-DTGDTRRRVVSLEED-KGLGAREKAFRLRTTMVGLDSTVLL
[Cellulomonas      MSG-QDD-QLSQVNSAAMAQLEVPDITPRQLAFVRLAKPLGLLDGTMLL
[Janibacter        MSE-PSE-DLDAVWRSILAAVSH-DGVFAPHRAFPLSLARFVGLLDGTALV
[Knoellia          MD-----QIWRITLDALDS-DGIPVQQRAFPLSLAKLVGLLDGTALI
[Kribbia           MGG-TDE-DFAQIWHNILDTLDA-DGVFVTERAFQIGKLVGLLDGTAVI
[Kytococcus        MSQTPDD-HATAINQEAHVHLQG-AGLAPRDIGVRLRLATLVGLLEGSTALL
[Ornithinimicrobium MTSQSPA-ESAENVQRVVSQLES-QGVTARDRAFLRLTQVGLLDITALL
[Serinicoccus      MSQ-PST-DSGDTWRRVVSLEED-KGLGAREKAFRLRTTMVGLDSTVLL
[Tetrasphaera      MD-----QIWRITLDALDN-DGIPVQQRAFPLSLARLVGLLDITALI
*               *       :  .  .  .  .  .  .  *  :*:*: *  ::

[Acidobacteria      AVFYQHTKDTLETTLRQPIVSALAEELGHDVRLAITVDESRLQELKAE-E
[Arsenicicoccus    AVFPYHTKEMLETTLRQPIVDLLSRELDREVRLAITVDDDVQRVEDE-A
[Cellulomonas      AVGNDLTKDYLETVRQEVTDALAAALGRDARFAITVDFSLDGAGDPS-L
[Janibacter        AVFNMYKTYVERALRVFVTQAFSAHYGQDVRLAVTVDFDLDDTEDEL-P
[Knoellia          AVFNDFTKDIVETRLRDRVTETLRSQGLGHDVRLAVTVDFSLGDAPVLVPA
[Kribbia           AVFNDFSKQFVEHRLRQHVTLSALSAQLGSEVRLAVTVDSLSLAEGGDTD-T
[Kytococcus        AVKYDHVKDAVEGHLREDVSTALAEVLDLDRDIRLAVSVDFDAVSAAQEE-A
[Ornithinimicrobium AVFYQHTKETLETTLRQPIVDALAGELGHDVRLAITVDEDLRRQVEDE-G
[Serinicoccus      AVFPYHTKEMLETTLRQPIVDLLSRELDREVRLAITVDDDVQRVEDE-A
[Tetrasphaera      AVFNDFTKDIVETRLRERVTETLSSQGLGHDVRLAVTVDQSLADAPAPE-A
**      *  :*  :*  :  :  .  :  *:*:***  .
```

**Asterix : (\*),** conserved amino acid in all sequences

**Colon (:),** A position of the MSA composed of residues having the same physicochemical properties

**Dot (.),** indicates the column of MSA for which semi-conserved substitutions are observed

Repeating blocks of 50 amino acid stretches

# Return to T-COFFEE Results - show color

CLUSTAL W (1.83) multiple sequence alignment

```

gi|497130903|ref|WP_009482734.1| MSVSGESSTPSEPGRIWGATLRALDQ-AGIPAPQRAFRLQATLVGVLDT
gi|497462772|ref|WP_009776970.1| M-----DQIWRITLDALDS-DGIPVQQRAFSLAKLVGLLDE
gi|497833122|ref|WP_010147278.1| MSQ-----PSTDSGDTWRRVVSLEED-KGLGAREKAFRLRTTMVGVLDSD
gi|499072896|ref|WP_010851794.1| MAD-----ASMTSVVRIIRALDR-EGVSHQERAFLSITRLAGVLDE
gi|502479361|ref|WP_012801520.1| MSQT----PDDHATAIWQEMVHLOG-AGLAPRDIGVLRLATLVGLLEG
gi|502628385|ref|WP_012865049.1| MAT-----TDDNISEIWKQAI AELEASPDITPRQLAFVKLAKPLGLFDG
gi|503647989|ref|WP_013882065.1| MAQ-----DEELSRVWGHVVITLEESPDITQRQLAFVRLAQPLGLLDG
gi|551300082|ref|WP_022920049.1| MTSQ----SPAESA EVWQRVVSQLES-QGVTARDRAFLRLTQLVGLLDT
gi|656266264|ref|WP_029212190.1| MTD-----AQVDVPRVWRDILRALES-GGISAQHRGFLRLSRLVGLLEG
gi|656321871|ref|WP_029253865.1| MPA-----AEVSI DEVWEQTIATLGSNPHMTRRQMGYVVKMAKPRAVFEG
*          *          : *          : . . : : .:::

```

```

gi|497130903|ref|WP_009482734.1| TALI AVPDDFTKEIVESRRARDLVRALTDQVGHVRLAVTVDASLREQFA
gi|497462772|ref|WP_009776970.1| TALI AVPNDFTKDIVETRLRDRVETETLSSQLGHDVRLAVTVDHSLADVPV
gi|497833122|ref|WP_010147278.1| TVLLAVPYPHTKEMLETTLRQPIVDLLSRELDREVRLAITVDDDVQRQVE
gi|499072896|ref|WP_010851794.1| TALI AVPNDFSKDIVETRLRGRISGHILTAELDRPLRLAVTVDPVSLAEAEF
gi|502479361|ref|WP_012801520.1| TALLAVKYDHVKDAVEGHLREDVSTALAEVLDLDRDIRLAVSVDPAVSAAQ
gi|502628385|ref|WP_012865049.1| TVIIAVANDHTRDFLETRVRAEVVQALSALGRDARFAITVDPVPELGFDEE
gi|503647989|ref|WP_013882065.1| TII LAVGNEYTKLEYLTKVRAEVT SALGSALGRDGRFAITVDPVSLVDDAP
gi|551300082|ref|WP_022920049.1| TALLAVPYQHTKETLETTLRQPIVDALAGELGHDVRLAITVDEDLRRQVE
gi|656266264|ref|WP_029212190.1| TALI AVPNDYTRDIVEKRIRTELVAALQEQLGRDVRRLAVTVDSSELSEA
gi|656321871|ref|WP_029253865.1| NVFLAVPADHVRTFIESSLRDDLVEALTSVLGTEVRF AISVEPDMDVQPP
. : : * . : : * * : * : . * : * : : .

```

Amino Acid Residue	3-Letter Code	1-Letter Code
Alanine	Ala	A
Arginine	Arg	R
Asparagine	Asn	N
Aspartic Acid	Asp	D
Cysteine	Cys	C
Glutamine	Gln	Q
Glutamic Acid	Glu	E
Glycine	Gly	G
Histidine	His	H
Isoleucine	Ile	I
Leucine	Leu	L
Lysine	Lys	K
Methionine	Met	M
Phenylalanine	Phe	F
Proline	Pro	P
Serine	Ser	S
Threonine	Thr	T
Tryptophan	Trp	W
Tyrosine	Tyr	Y
Valine	Val	V

AVFPMILW	RED	Small (small+ hydrophobic (incl.aromatic -Y))
DE	BLUE	Acidic
RK	MAGENTA	Basic
STYHCNGQ	GREEN	Hydroxyl + Amine + Basic - Q
Others	Gray	

The colors give information about the amino acid

# Sequence-based Similarity Data Module

## 4 TOOLS

### 1. BLAST

The Basic Local Alignment Search Tool (**BLAST**) finds regions of local similarity between sequences and calculates the statistical significance of matches

### 2. CDD

Conserved Domain Database Search (**CDD**) finds sequence similarity with genes in conserved orthologous groups (COGs).

### 3. T-Coffee

Tree based **C**onsistency **O**bjective **F**unction **F**or alignment **E**valuation (**T-Coffee**) is a multiple sequence alignment program that aligns a set of homologous (similar ) sequences

### 4. WebLogo

WebLogo is a program that enables easy creation of sequence logos from the multiple sequence alignments



<http://weblogo.berkeley.edu/>

- This is a program designed to enable easy creation of sequence logos from multiple sequence alignments.
- One simple graphic is generated.
- At least 10 sequences should be used.
- Save image as .png and attach to notebook.

[Crooks GE](#), [Hon G](#), [Chandonia JM](#), [Brenner SE](#) WebLogo: A sequence logo generator, *Genome Research*, 14:1188-1190, (2004)

Created by [Computational Genomics Research Group](#), Department of Plant and Microbial Biology, University of California, Berkeley

# WebLogo

## WebLogo

go to <http://weblogo.berkeley.edu/>

Sequence logo



logo image

Comments/observations about the Multiple Sequence Alignment WebLogo

comments

Provide an overall summary of your findings from the Sequence Similarity Module in the box below.

**Be sure to save this document after completing the sequence-based information module!**

# Home page for WebLogo



· [about](#) · [create](#) · [examples](#) ·

[Version 2.8.2 \(2005-09-08\)](#)

(= [WebLogo 3](#))

## References

[Crooks GE, Hon G, Chandonia JM, Brenner SE](#) WebLogo: A sequence logo generator. *Genome Research*, 14:1188-1190, (2004) [[Full Text](#)]

Schneider TD, Stephens RM. 1990. [Sequence Logos: A New Way to Display Consensus Sequences](#). *Nucleic Acids Res.* 18:6097-6100

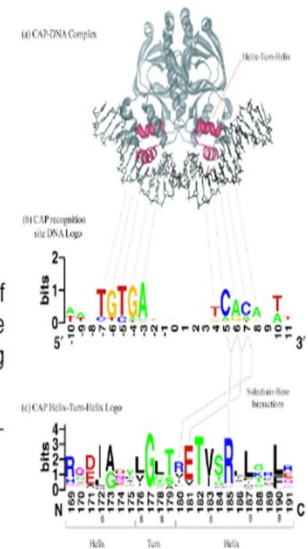
## Introduction

[WebLogo](#) is a web based application designed to make the [generation](#) of sequence logos as easy and painless as possible. Click [here](#) to create your own sequence logos.

[Sequence logos](#) are a graphical representation of an amino acid or nucleic acid multiple sequence alignment developed by [Tom Schneider](#) and [Mike Stephens](#). Each logo consists of stacks of symbols, one stack for each position in the sequence. The overall height of the stack indicates the sequence conservation at that position, while the height of symbols within the stack indicates the relative frequency of each amino or nucleic acid at that position. In general, a sequence logo provides a richer and more precise description of, for example, a binding site, than would a consensus sequence.

## WebLogo Source Code

The WebLogo source code is available for [download](#). See the [README](#) file for installation instructions and release notes, and [LICENSE](#) for the open source license.





# Paste **multiple alignment** with top and bottom line intact

CLICK 32 so that aa are chopped in 32 bits and shown

The screenshot shows the WebLogo interface for creating a sequence logo. The main section is titled "Multiple Sequence Alignment" and displays a CLUSTAL W (1.83) multiple sequence alignment. The alignment shows several protein sequences with gaps represented by dashes. Below the alignment, there are sections for "Image Format & Size" and "Advanced Logo Options".

**Multiple Sequence Alignment**

```
CLUSTAL W (1.83) multiple sequence alignment
Ksed_00010          VSQIFP-----
DHATAIQEAMVH      MVADQ-----
qi|118706|sp|P21173.1|DNAA_MICLU      MVADQ-----
AVLSSWRSVVGS
qi|123144805|sp|Q0SAG7.1|DNAA_RHOSR      MNDDFN-----
RLARSLIDVYAD      MSEGG-----
qi|123774918|sp|Q47U23.1|DNAA_THEFY      MSEGG-----
INLAMVRSRVLDN
qi|166214685|sp|A1T102.1|DNAA_MYCVF      MIIIDP-----
```

**Image Format & Size**

Image Format: PNG (bitmap) | Logo Size per Line: 18 X 5 cm

**Advanced Logo Options**

Sequence Type:  amino acid  DNA / RNA  Automatic Detection

First Position Number: 1

Small Sample Correction:

Multi-line Logo (Symbols per Line):  ( 32 )

**Advanced Image Options**

Bitmap Resolution: 96 pixels/inch (dpi)

Show Y-Axis:

Show X-Axis:

Show Error Bars:

Boxed / Boxed Shrink Factor:  / 0.5

Show fine print:

Antialias Bitmaps:

Y-Axis Height: (bits)

Y-Axis Label: bits

X-Axis Label:

Label Sequence Ends:

Outline Symbols:

Y-Axis Tic Spacing: 1 (bits)

**Colors**

Color Scheme:  Default  Black & White  Custom (See Below.)

# Create → Paste **multiple alignment** with top and bottom line intact

**CLICK 32**

**aa are shown in repeating blocks of 32**

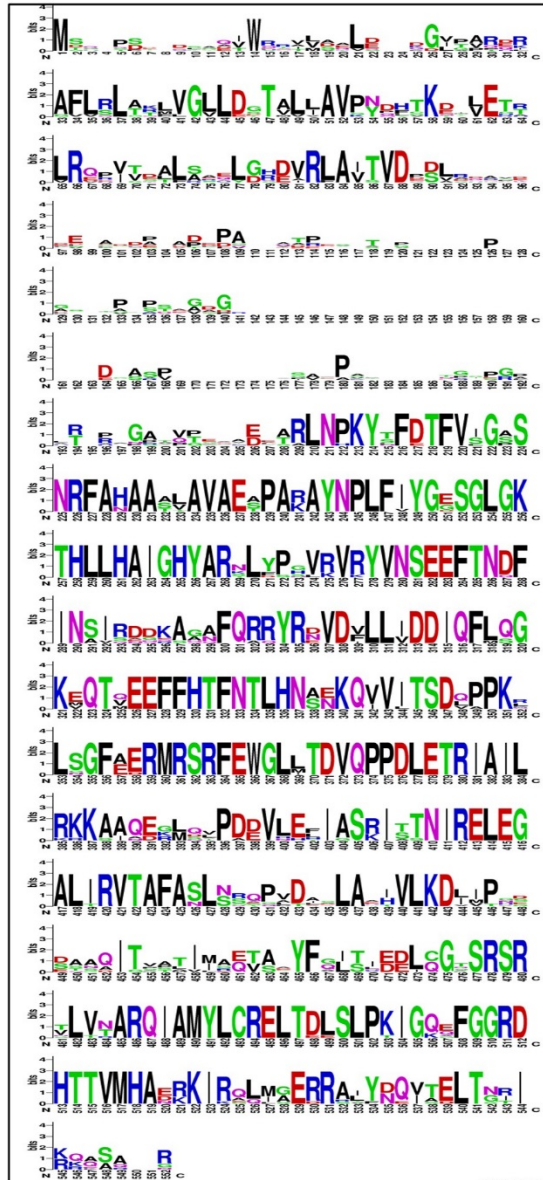
The screenshot shows a web-based logo generation tool. At the top, there is a text area containing a multiple sequence alignment of DNA sequences. Below this is the 'Upload Sequence Data' section with a 'Browse...' button. The 'Image Format & Size' section includes a dropdown for 'Image Format' (set to PNG (bitmap)) and a 'Logo Size per Line' field (set to 18 X 5 cm). The 'Advanced Logo Options' section has radio buttons for 'Sequence Type' (amino acid, DNA / RNA, Automatic Detection) and a 'Multiline Logo (Symbols per Line)' field set to 32. The 'Advanced Image Options' section includes 'Bitmap Resolution' (96 pixels/inch), 'Show Y-Axis', 'Show X-Axis', 'Show Error Bars', 'Boxed / Boxed Shrink Factor', and 'Show fine print'. The 'Colors' section has radio buttons for 'Color Scheme' (Default, Black & White, Custom) and a table for defining colors. The table has columns for 'Symbols', 'Color', and 'RGB'. The 'Create Logo' and 'Reset' buttons are at the bottom right.

Color Scheme	Symbol	Color	RGB	Symbol	Color	RGB
Default	KRH	green			purple	
	DE	blue			orange	
	AVLIPWFM	red			black	
		black		Other	black	

**CLICK**

# Weblogo of the entire protein

## Save as a .png file. Upload into notebok



### COLOR CODE

Polar amino acids  
(G,S,T,Y,C,Q,N) : green

Basic (K,R,H) blue,

Acidic (D,E) red

Hydrophobic  
(A,V,L,I,P,W,F,M) : black.

The most common residue at each position in the alignment will be the largest letter at that position.

The relative height of the stack of letters will be proportional to the % conservation.

The relative widths of the stacks indicate the proportion of valid readings of nucleic bases or amino acids at that position.

The more gaps in the sequence at a specific position means a thinner stack.