

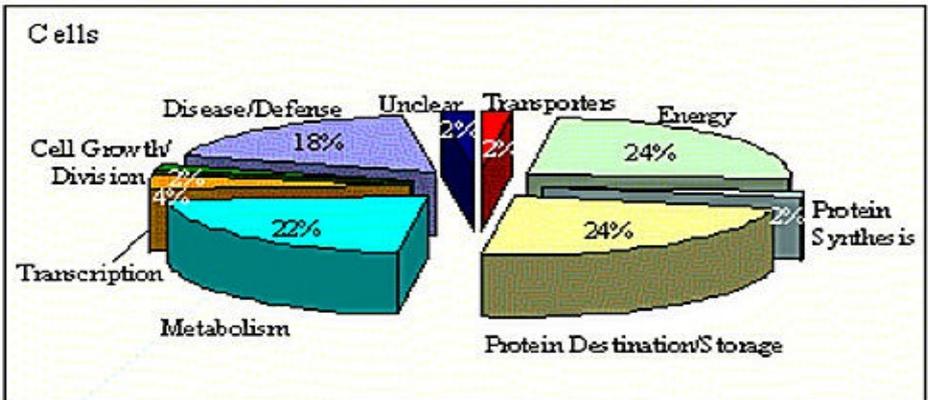
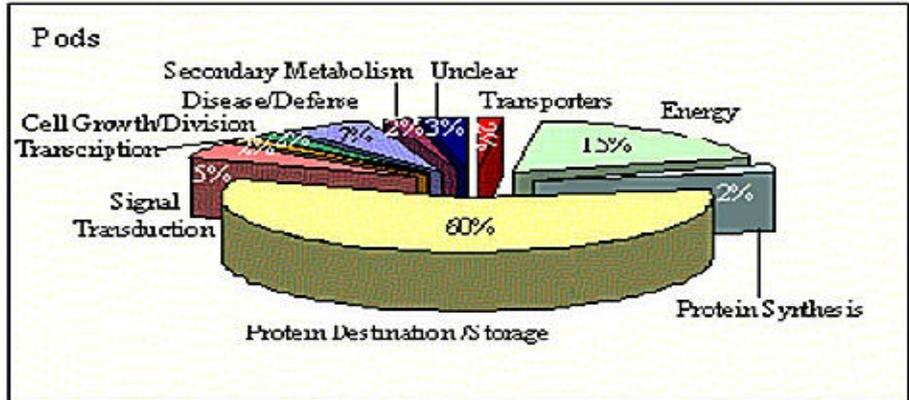
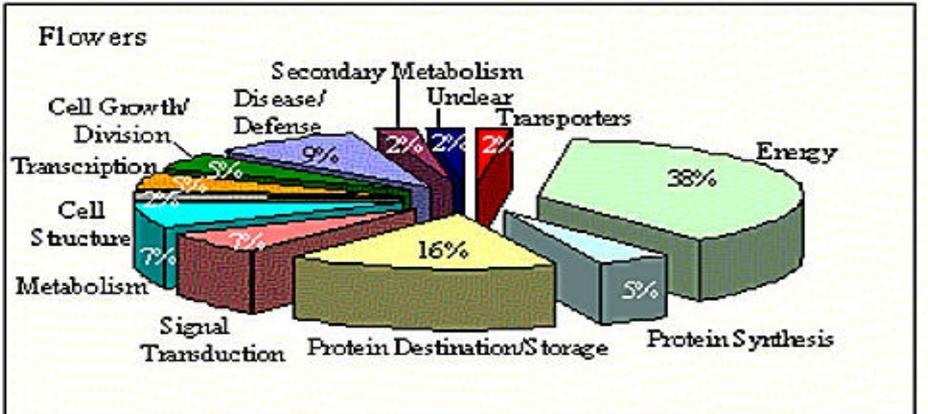
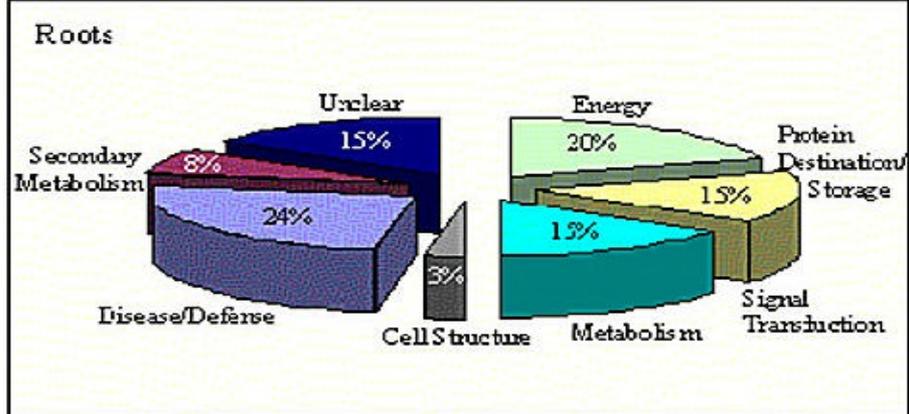
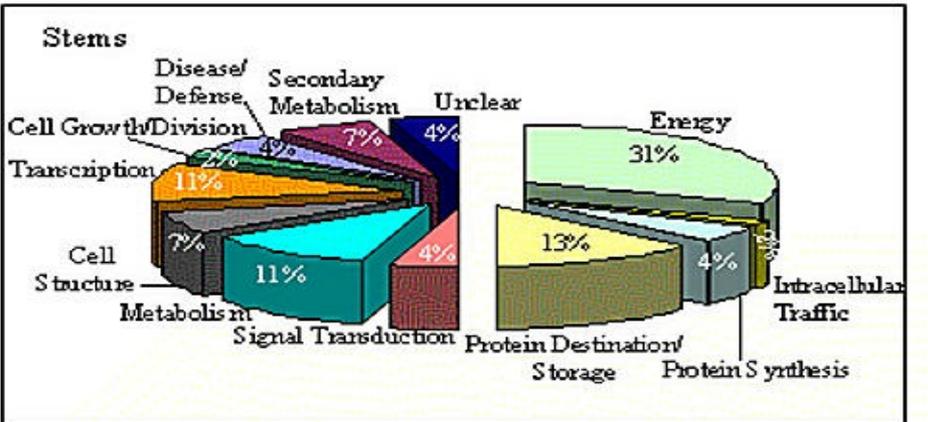
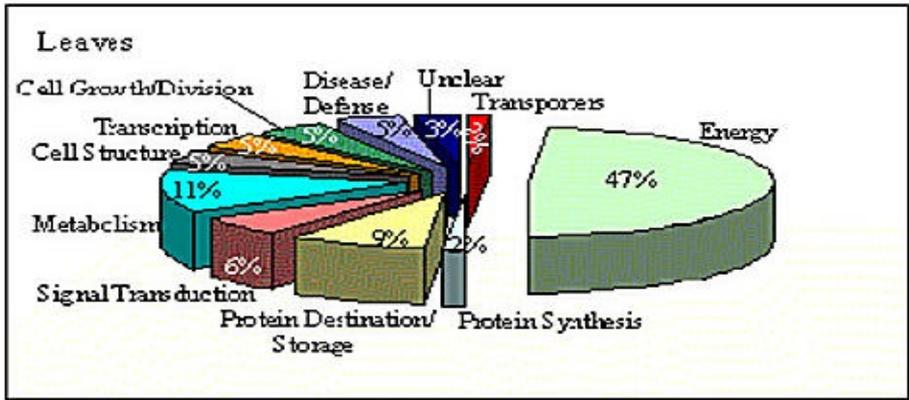
Introduction to **Systems Biology**, **P4 Medicine** and **Bioinformatics**

CH370 - Biochemistry

Spring, 2008

Marvin Hackert

Summary of the functions of various proteins identified in specific tissues of *M. truncatula*.



Genomics

Proteomics

Interactomics

Systems Biology –

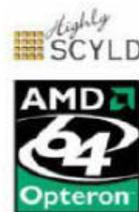
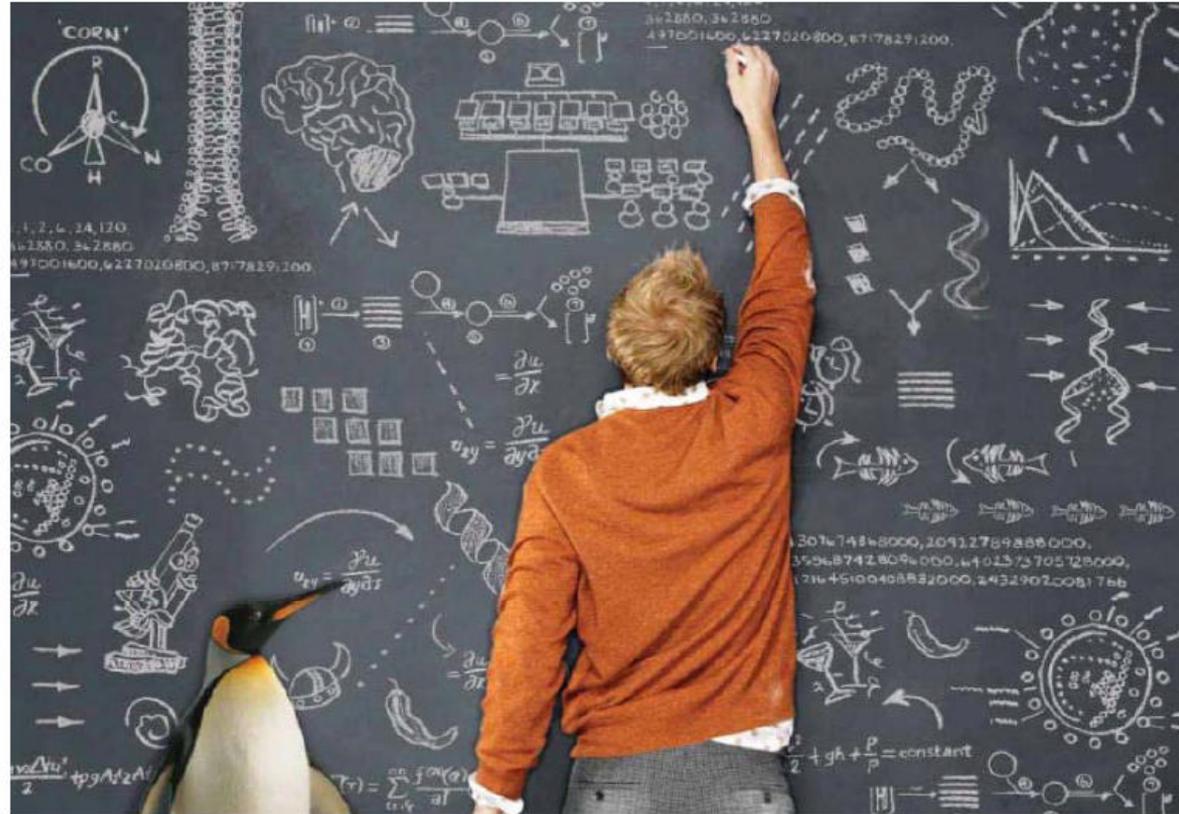
None of these fields of research would be possible without **Bioinformatics**, which would not be possible with lots of **computing power!**

THE PENGUIN CURES

Millions of compounds to go. Database analyses keep one computer clogged; while a microarray analysis chokes the other. The computer hopping begins; so does the throbbing in your brain. Exhale. Penguin Computing® Clusters combine the economy of Linux with the ease of Scyld. Unique,

centrally-managed Scyld ClusterWare™ HPC makes large pools of Linux servers act like a single virtual system. So you get supercomputer power, manageability and scalability, without the supercomputer price. Penguin Computing. So many drugs. So little time.

HEADACHES



PENGUIN HIGH DENSITY CLUSTER. The highest density modular blade server architecture on the market, with powerful Scyld ClusterWare™ HPC for single point command and control, and AMD Dual Core Opteron™ for a highly productive user experience.

To find out more about optimizing your discovery engine, read our whitepaper at www.penguincomputing.com/go/whitepaper



Penguin Computing and the Penguin Computing logo are registered trademarks of Penguin Computing, Inc. Scyld ClusterWare and the Highly Scalable logo are trademarks of Scyld Systems Corporation. AMD Opteron and the AMD logo are trademarks or registered trademarks of Advanced Micro Devices, Inc. Linux is a registered trademark of Linus Torvalds. ©2006 Penguin Computing, Inc. All rights reserved.

Genome – the genome of an organism is its whole hereditary information encoded in its DNA (or, RNA for some viruses) and includes both the coding (genes) and non-coding sequences of the DNA.

Proteome – Proteomics is often considered the next step in the study of biological systems, after **genomics**. It is much more complicated than genomics, mostly because while an organism's genome is rather constant, a **proteome differs from cell to cell and constantly changes** through its biochemical interactions with the genome and the environment. The **Proteome is dynamic** (*Cell type, time, conditions*).

Interactome – whole set of molecular interactions in cells, in the context of proteomics, it refers to protein-protein interaction network(PPI), or protein network (PN).

Systems Biology - seeks to understand how biological systems function. By studying the relationships and interactions between various parts of a biological system (e.g. metabolic pathways, organelles, cells, physiological systems, organisms etc.), it is hoped that eventually a model of the whole system can be developed.

Brief Introduction to Bioinformatics

Terms: NCBI / EMBL

(Sequence Alignments)

Sequence databases

FASTA

Scoring Matrix

PAM

BLOSUM

Smith – Waterman

BLAST

PSI – BLAST

Raw Score

Probability Value

E-value

ClustalW

Acknowledgement: This brief introduction on Sequence Alignments is based on information found at web sites such as that at NCBI and EMBL-EBI, and the slides illustrating the alignment algorithm were taken from a handout provided by Dr. Ed Marcotte (Univ. of Texas at Austin) who teaches a course on Bioinformatics (CH391L) and on-line web notes of Michael Yaffe at MIT.

Ref:

<http://www.ncbi.nlm.nih.gov/>

<http://www.ebi.ac.uk/clustalw/#>



National Center for Biotechnology Information

National Library of Medicine

National Institutes of Health

PubMed All Databases **BLAST** OMIM Books TaxBrowser Structure

Search All Databases for Go

SITE MAP

Alphabetical List
Resource Guide

About NCBI

An introduction to
NCBI

GenBank

Sequence
submission support
and software

Literature databases

PubMed, OMIM,
Books, and PubMed
Central

Molecular databases

Sequences,
structures, and
taxonomy

What does NCBI do?

Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease. [More...](#)

100 Gigabases

GenBank and its collaborating databases, the European Molecular Biology Laboratory and the DHA Databank of Japan, have reached a milestone of 100 billion bases from over 165,000 organisms. See the [press release](#) or find more information on [GenBank](#).

CCDS Database

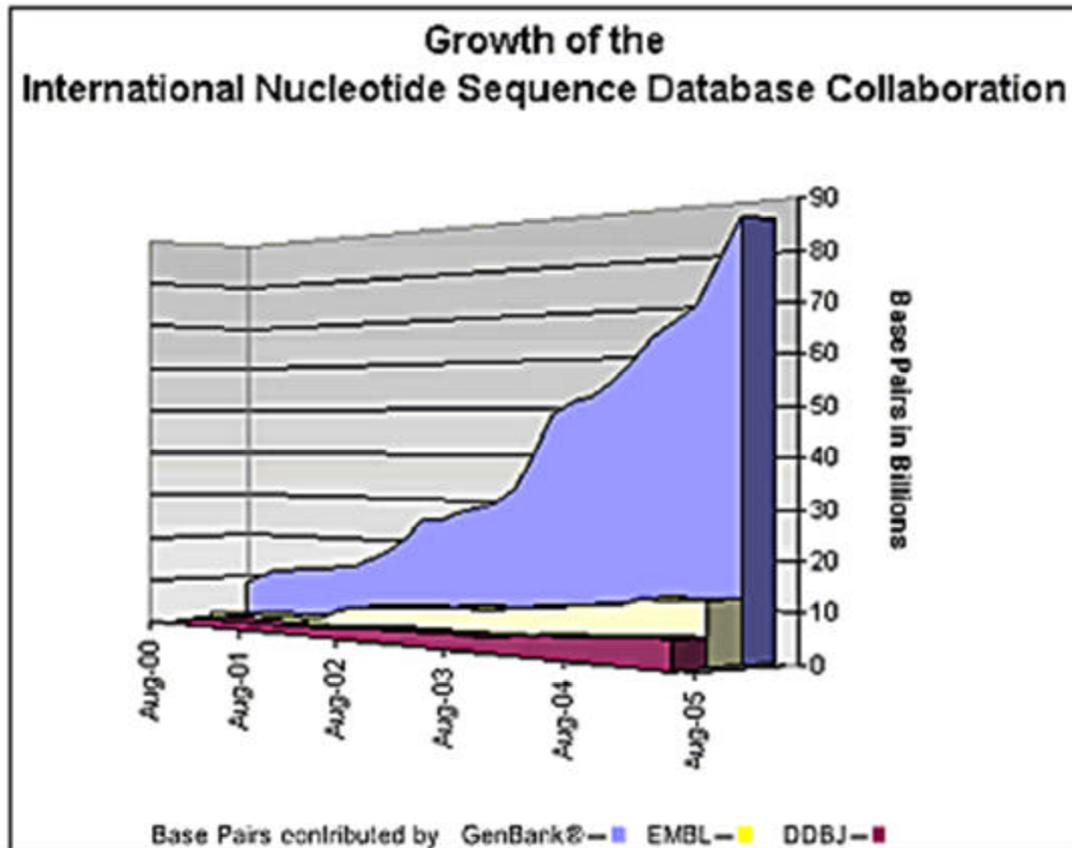
Hot Spots

- ▶ Assembly Archive
- ▶ Clusters of orthologous groups
- ▶ Coffee Break, Genes & Disease, NCBI Handbook
- ▶ Electronic PCR
- ▶ Entrez Home
- ▶ Entrez Tools
- ▶ Gene expression omnibus (GEO)
- ▶ Human genome resources
- ▶ Malaria genetics & genomics

International sequence databases exceed 100 gigabases

In August 2005, the INSDC announced the DNA sequence database exceeded 100 gigabases. GenBank is proud of its contributions toward this milestone. We thank all the scientists who have worked through the submission process at GenBank and made their sequence data available to the world. See the related [press release](#).

>100,000,000,000 bases



> 200,000 organisms!!

Computational biology & Bioinformatics

Computational biology and **bioinformatics** focus on the computational/ theoretical study of biological processes, and much of the disciplines involve constructing models like those above, then testing/validating/proving/applying these models using computers, hence the nickname “*in silico* biology”. The fields are closely related: computational biology is the more inclusive name, and **bioinformatics often refers more specifically to the use of “informatics” tools like databases and data mining.**

Big problems tackled by these fields include:

Assembling complete genomes from pieces of sequenced DNA

Finding genes in genomes

Modeling networks & interactions of proteins

Predicting protein/RNA folding, structure, and function

Sequence alignments (BLAST)

BLAST – Basic Local Alignment Search Tool

Many “flavors” of BLAST

<u>Program</u>	<u>Query</u>	<u>Database</u>
BLASTP	aa	aa
BLASTN	nt	nt
BLASTX	nt (\Rightarrow aa)	aa
TBLASTN	aa	nt (\Rightarrow aa)
TBLASTX	nt (\Rightarrow aa)	nt (\Rightarrow aa)
PsiBLAST	aa (aa msa)	aa

(Position-Specific Iterative)

Why Align Sequences?

Identify Protein or Gene from Partial Information

Infer Functional Information

Infer Structural Information

Infer Evolutionary Relationships

Assumes:

conservation of
sequence



conservation of
function

BUT: Function carried out at level of proteins, i.e.
3-D structure

Sequence conservation carried out at level of DNA
1-D sequence

Sequence Alignment

The *Smith-Waterman algorithm* considers a simple model for protein sequence evolution that allows us to align amino acid sequences of proteins to see if the proteins are related. **BLAST** is designed to **mimic this algorithm**, but BLAST is much faster due to some shortcuts and approximations and clever programming tricks.

This **process of gene evolution** can be modeled as **a stochastic process of gene mutation** followed by a “**selection**” **process for those sequences still capable of performing their given roles** in the cell. Over enough time, as new species evolve & diverge from related species, this has the result of producing families of related gene sequences, more similar in regions where that particular sequence is critical for the function of the molecule, and less similar in regions less critical for the molecule's function. Frequently, **we observe only the products of millions of years of this process**. Given a set of molecules (DNA, RNA or protein sequences) - **??** How can we **decide if they are similar enough to be considered part of the same family** or if the **observed similarity is just present by random chance**.

Database Searching

The Assumptions:

The sequences being sought have **an evolutionary ancestral sequence** in common with the query sequence (our newly determined sequence).

Our best guess at the **actual path of evolution is the path that requires the fewest evolutionary events.**

All substitutions are not equally likely and should be weighted to account for this.

Gaps: **Insertions and deletions are less likely than substitutions** and should be weighted to account for this. In most alignment and search programs, the gap penalty consists of two terms, the cost to **open** the gap and the cost to **extend** the gap.

FASTA Format

- This format contains a one line header followed by lines of sequence data.
- Sequences in fasta formatted files are preceded by a line starting with a ">" symbol.
- The first word on this line is the name of the sequence. The rest of the line is a description of the sequence.

Term	Entry Name	Molecule Type	Gene Name	Sequence Length
e.g.	FOSB_MOUSE	Protein	fosB	338 bp

- The remaining lines contain the sequence itself.
- Blank lines in a FASTA file are ignored, and so are spaces or other gap symbols (dashes, underscores, periods) in a sequence.
- Fasta files containing multiple sequences are just the same, with one sequence listed right after another. This format is accepted for many multiple sequence alignment programs.

```
>FOSB_MOUSE Protein fosB. 338 bp
MFQAFPGDYDSGSRCSSSPSAESQYLSSVDSFGSPPTAAASQECAGLGEMPGSFVPTVTA
ITTSQDLQWLVOPTLISSMAQSQGQPLASQPPAVDPYDMPGTSYSTPGLSAYSTGGASGS
GGPSTSTTTSGPVSARPARARPRRPREETLTPEEEEEKRRVRRRERNKLAAAKCRNRRRELT
DRLQAETDQLEEEKAELESEIAELQKEKERLEFVLVAHKPGCKIPYEEGPGGPLAEVRD
LPGSTSAKEDGFGWLLPPPPPPPLPFQSSRDAPPNLTA SLFTHSEVQVLGDPFPVSPSY
TSSFVLTCEVSAFAGAQR TSGSEQPSDPLNSP L LAL
```

Examples of aligned protein sequences:

Shown are 3 pairs of sequences, showing aligned sequences of proteins named FlgA1, FlgA2, FlgA3, and HvcPP. Between each pair the perfect matches and close matches (shown by + symbols, indicating chemically similar amino acids) are written.

Two biologically related proteins with similar sequences:

FlgA1 EAGNVKLRGRLDTLPPRTVLDINQLVDAISLRDLSPDQPIQLTQFRQAWRVKAGQQRVNVIASGD
++K+K+GRLDTLPP +L+ N A+SLR ++ QP+ R+ W +KAGQ V V+A G+ (28/65)
FlgA2 TLQDIKMKQGRDLTLPPGALLEPNFAQGAVSLRQINAGQPLTRNMLRRLWI KAGQDVQVLALGE (186)

Also biologically related (& fold up into the same 3D protein structure):

FlgA1 EAGNVKLRGRLDTLPPRTVLDINQLVDAISLRDLSPDQPIQLTQFRQAWRVKAGQQRVNVIASGD
A + P +L I+ R L P + I R+AW V+ G V V (15/65)
FlgA3 LAALKQVTTLIAGKHKPDAMATHAEELQGKIAKRTLLPGRYIPTAAIREAWLVEQGA AVQVFFIAG (50)

But these are biologically unrelated (& fold up into unrelated structures):

FlgA1 AGNVKLRGRLDTLPPRTVLDINQLVDAISLRDLSPDQPIQLTQFRQA-WRVKAGQQRVNVIASGD
AG+V K G + + PRT ++ I+ P PI +++A WRV A + V V+ GD (21/65)
HvcPP AGHV--KNGTMRIVGPRTCSNVWNGTFPINATTTGPSIPI PAPNYKKALWRV SATEYVEVVRVGD (128)

The problem we face is how to distinguish the biologically meaningless match (FlgA1-HvcPP) from the biologically meaningful ones (FlgA1-FlgA2 and FlgA1-FlgA3)?

To align two sequences, we need:

1. Some way to decide which alignments are better than others. For this, we'll **invent a way to give the alignments** a “**score**” indicating their quality.

→ **“Scoring Matrix”**

2. Some **way to align the proteins** so that they **get the best possible score**.

→ ***Smith-Waterman algorithm***

dynamic programming, recursive manner

3. Then finally, some way to **decide when a score is “good enough”** for us to believe the alignment is biologically significant.

→ **“Scramblings - Expect Values”**

extreme value distribution

What is a **scoring matrix**?

The aim of a **sequence alignment**, is to match "the most similar elements" of two sequences. This similarity must be evaluated somehow.

For example, consider the following two alignments:

AIWQH	AIWQH
AL--QH	A--LQH

They seem quite similar: both contain one "gap" and one "substitution," just at different positions. However, the first alignment is the better one because isoleucine (I) and leucine (L) are similar sidechains, while tryptophan (W) has a very different structure. This is a **physico-chemical** measure; we might prefer these days to say that leucine simply substitutes for isoleucine more frequently---without giving an underlying "reason" for this observation.

However we explain it, **it is much more likely that a mutation changed I into L and that W was lost, than W was changed into L and I was lost.** We would expect that a change from I to L would not affect the function as much as a mutation from W to L--but this deserves its own topic.

To **quantify the similarity achieved by an alignment**, **scoring matrices** are used: they contain a value for each possible substitution, and the **alignment score** is the sum of the matrix's entries for each aligned amino acid pair. For gaps a special **gap score** is necessary---just add a constant penalty score for each new gap. The **optimal alignment** is the one which **maximizes the alignment score**.

Importance of scoring matrices

Scoring matrices appear in all analysis involving sequence comparison. The choice of matrix can strongly influence the outcome of the analysis. Scoring matrices implicitly represent a particular theory of evolution. Understanding theories underlying a given scoring matrix can aid in making proper choice.

Types of matrices

Ref: <http://www.ebi.ac.uk/clustalw/#>

[PAM](#)

[BLOSSUM](#)

[GONNET](#)

[DNA Identity Matrix](#)

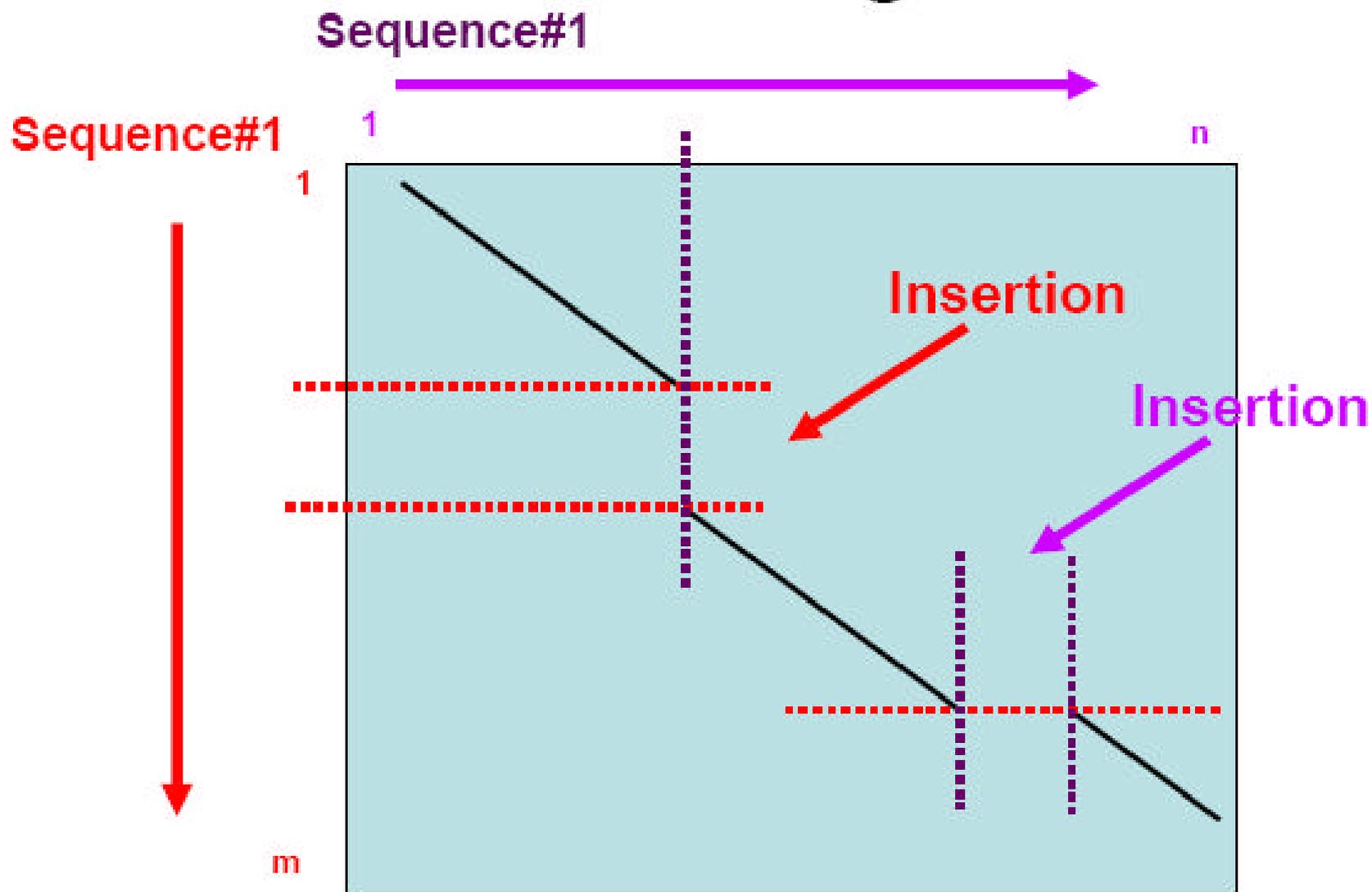
Unitary Scoring Matrices

Early sequence alignment programs used **unitary scoring matrix**. A unitary matrix **scores all matches the same and penalizes all mismatches the same**. Although this scoring is sometimes appropriate for DNA and RNA comparisons, **for protein alignments using a unitary matrix amounts to proclaiming ignorance about protein evolution and structure**. Thirty years of research in aligning protein sequences have shown that different matches and mismatches among the 400 amino acid pairs that are found in alignments require different scores.

	A	T	G	C
A	<u>1</u>			
T	-10000	<u>1</u>		
G	-10000	-10000	<u>1</u>	
C	-10000	-10000	-10000	<u>1</u>

Many alternatives to the unitary scoring matrix have been suggested. One of the earliest suggestions was **scoring matrix based on the minimum number of bases that must be changed to convert a codon for one amino acid into a codon for a second amino acid**. This matrix, known as the **minimum mutation distance matrix**, has succeeded in identifying more distant relationships among protein sequences than the unitary matrix approach.

Dot Matrix Alignments



Evolutionary Distances

The best improvement achieved over the unitary matrix was based on **evolutionary distances**. **Margaret Dayhoff** pioneered this approach in the 1970's. She made an extensive study of the frequencies in which amino acids substituted for each other during evolution. The **studies involved carefully aligning all of the proteins in several families of proteins and then constructing phylogenetic trees for each family**. Each phylogenetic tree was examined for the substitutions found on each branch. This led to a **table of the relative frequencies with which amino acids replace each other over a short evolutionary period**.

This table and the relative frequency of occurrence of the amino acids in the proteins studied were combined in computing the **PAM (Point Accepted Mutations) family of scoring matrices**.

From a biological point of view **PAM matrices are based on observed mutations**. Thus **they contain information about the processes that generate mutations** as well as the criteria that are important in selection and in fixing a mutation within a population. From a **statistical point of view PAM matrices, and other log-odds matrices**, are the most accurate description of the changes in amino acid **composition** that are expected after a given number of mutations that can be derived from the data used in creating the matrices. Thus the highest scoring alignment is statistically the most likely to have been generated by evolution rather than by chance.

Log-odds scoring

Log-odds matrices: Each score in the matrix is the **logarithm of an odds ratio**. The odds ratio used is the **ratio of the number of times residue "A" is observed to replace residue "B" divided by the number of times residue "A" would be expected to replace residue "B" if replacements occurred at random.**

Deriving realistic substitution matrices:

First need to know frequency of one amino acid substituting for another in related proteins [=P(ab)] c/w the chance that substituting one for the other occurred by chance, based on the relative frequencies of each amino acid in proteins, q(a) and q(b). Call this the "odds ratio": $P(ab)/q(a)q(b)$

If we do this for all positions in an alignment, then the total probability will be the product of the odds ratios at each position....but multiplication is computationally expensive....so....take the **log (odds ratio)** and add them instead.

The **BLOSUM family of matrices** developed by Steven and Jorja Henikoff are one of these newly developed log-odds scoring matrices. The **improved performance of the BLOSUM** matrices can be attributed to **many more protein sequences known now**, thus they incorporate many more observed amino acid substitutions, and because the **substitutions used in constructing the BLOSUM matrices are restricted to those substitutions found within well conserved blocks in a multiple sequence alignment.**

PAM (Percent Accepted Mutation)

A **unit** introduced by M.O. Dayhoff et al. to quantify the amount of **evolutionary change** in a protein sequence. **1.0 PAM unit, is the amount of evolution which will change, on average, 1% of amino acids in a protein sequence.** A PAM(x) substitution matrix is a look-up table in which scores for each amino acid substitution have been calculated based on the frequency of that substitution in closely related proteins that have experienced a certain amount (x) of evolutionary divergence.

PAM matrices are based on global alignments of closely related proteins.

71 groups of protein sequences, 85% similar
1572 amino acid changes.

Functional proteins → "Accepted" mutations by natural selection

PAM1 matrix means 1% divergence between proteins - i.e. 1 amino acid change per 100 residues. Some texts re-state this as the probability of each amino acid changing into another is ~ 1% and probability of not changing is ~99%

The optimal alignment of two very similar sequences with PAM 500 may be less useful than that with PAM 50.

Construction of a Dayhoff Matrix: PAM1

Step 1: Measure pairwise substitution frequencies for each amino acid within families of related proteins

↓

... . GDS**F**H**Y**FV**S**HG... . .
... . GDS**F**H**Y**YV**S**F**G**... . .
... . GDS**Y**H**Y**FV**S**F**G**... . .
... . GDS**F**H**Y**FV**S**F**G**... . .
... . GDS**F**H**F**FV**S**F**G**... . .

900 Phe (F)....+ another 100 probable Phe but...

100 Phe (F) → 80 Tyr (Y), 3 Trp (W), 2 His (H)....

Gives f_{ab} , i.e. $f_{FY}=80$
 $f_{FW}=3$

....by evolution!

<u>Amino Acid Change</u>	<u>PAM 1 Score</u>	<u>PAM 250 Score</u>
F→A	0.0002	0.04
F→R	0.0001	0.01
F→N	0.0001	0.02
F→D	0.0000	0.01
F→C	0.0000	0.01
F→Q	0.0000	0.01
F→E	0.0000	0.01
F→G	0.0001	0.03
F→H	0.0002	0.02
F→I	0.0007	0.05
F→L	0.0013	0.13
F→K	0.0000	0.02
F→M	0.0001	0.02
F→F	0.9946	0.32
F→P	0.0001	0.02
F→S	0.0003	0.03
F→T	0.0001	0.03
F→W	0.0001	0.01
F→Y	0.0021	0.15
F→V	0.0001	0.05

These are the M_{ab} values!
i.e. the chance that one amino acid will replace another at 250 PAMs in two proteins that are evolutionarily related to each other!

SUM = 1.0

But we have to use the right matrix!!!

PAM 250 matrix – 250% expected change

Sequences still ~ 15-30 % similar, i.e. Phe will match Phe ~ 32% of the time
Ala will match Ala ~ 13% of the time

Expected % similarity

Other PAM matrices: PAM 120 – 40%
PAM 80 – 50%
PAM 60 – 60% } Use for similar sequences

PAM250 – 15-30% similarity.

Use the correct PAM matrix for alignments based on how similar the sequences to be aligned are! But wait.....how do we know that in the first place? Usually don't!!!!.

So..... try PAM200, PAM120, PAM60, PAM80, and PAM30 matrix and use the one that gives the highest ungapped alignment score

Alternative amino acid matrices

Problems with Dayhoff:

- Based on amino acids, not nucleotides.
- Assumes evolutionary model with explicit phylogenetic relationships, and circular arguments: alignment → matrices; matrices → new alignments.
- Based on a small set of closely related molecules.
- **Gonnett, Cohen & Benner**
 - All against All database matching using DARWIN
 - 1,700,000 matches
 - Compile mutation matrices at different PAMs DIRECTLY*
- **BLOSUM = Blocks Amino Acid Substitution Matrices-Henikoff&Henikoff 1992**
 - based on a much larger dataset from ~500 Prosite families identified by Bairoch using conserved amino acid patterns “blocks” that define each family.

Typically used for multiple sequence alignment.
AA substitutions noted, log odds ratios derived.

for example...Block patterns 60% identical give rise to Blosum60 matrix, etc....i.e. conservation of functional blocks based on un-gapped alignments.
Blosum62 - best match between information content and amount of data
Not based on explicit evolutionary model

BLOSUM matrices are based on local alignments.

BLOSUM (BLOcks SUBstitution Matrix): BLOSUM 62 is a matrix calculated from comparisons of sequences with no less than 62% divergence.

BLOSUM 62 is the default matrix in BLAST 2.0. Though it is tailored for comparisons of moderately distant proteins, it performs well in detecting closer relationships. A search for distant relatives may be more sensitive with a different matrix.

Differences between PAM and BLOSUM

PAM matrices are based on an **explicit evolutionary model** (that is, replacements are counted on the branches of a phylogenetic tree), whereas the **BLOSUM** matrices are based on an **implicit** rather than explicit **model of evolution**.

The sequence variability in the alignments used to count replacements. The **PAM** matrices are based on mutations observed throughout a **global alignment**, this includes both highly conserved and highly mutable regions. The **BLOSUM** matrices are **based only on highly conserved regions** in series of alignments forbidden to contain gaps.

BLOSUM62 Substitution Scoring Matrix. The BLOSUM 62 matrix is a 20 x 20 matrix in which every possible identity and substitution is **assigned a score based on the observed frequencies of such occurrences in alignments of related proteins.** Identities are assigned the most positive scores. **Frequently observed substitutions also receive positive scores** and **seldom observed substitutions are given negative scores.**

Blosum 45 Amino Acid Similarity Matrix

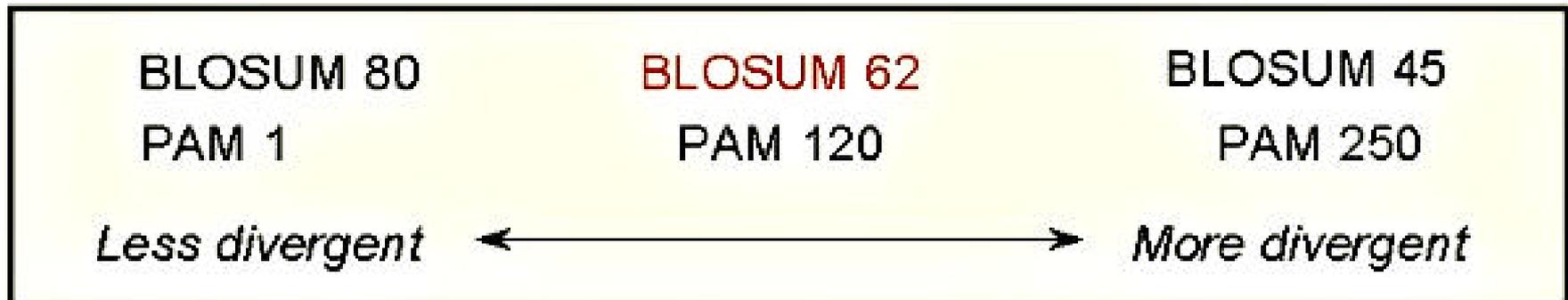
<u>Gly</u>	7																			
<u>Pro</u>	-2	9																		
<u>Asp</u>	-1	-1	7																	
<u>Glu</u>	-2	0	2	6																
<u>Asn</u>	0	-2	2	0	6															
<u>His</u>	-2	-2	0	0	1	10														
<u>Gln</u>	-2	-1	0	2	0	1	6													
<u>Lys</u>	-2	-1	0	1	0	-1	1	5												
<u>Arg</u>	-2	-2	-1	0	0	0	1	3	7											
<u>Ser</u>	0	-1	0	0	1	-1	0	-1	-1	4										
<u>Thr</u>	-2	-1	-1	-1	0	-2	-1	-1	-1	2	5									
<u>Ala</u>	0	-1	-2	-1	-1	-2	-1	-1	-2	1	0	5								
<u>Met</u>	-2	-2	-3	-2	-2	0	0	-1	-1	-2	-1	-1	6							
<u>Val</u>	-3	-3	-3	-3	-3	-3	-3	-2	-2	-1	0	0	1	5						
<u>Ile</u>	-4	-2	-4	-3	-2	-3	-2	-3	-3	-2	-1	-1	2	3	5					
<u>Leu</u>	-3	-3	-3	-2	-3	-2	-2	-3	-2	-3	-1	-1	2	1	2	5				
<u>Phe</u>	-3	-3	-4	-3	-2	-2	-4	-3	-2	-2	-1	-2	0	0	0	1	8			
<u>Tyr</u>	-3	-3	-2	-2	-2	2	-1	-1	-1	-2	-1	-2	0	-1	0	0	3	8		
<u>Trp</u>	-2	-3	-4	-3	-4	-3	-2	-2	-2	-4	-3	-2	-2	-3	-2	-2	1	3	15	
<u>Cys</u>	-3	-4	-3	-3	-2	-3	-3	-3	-3	-1	-1	-1	-2	-1	-3	-2	-2	-3	-5	12
<u>Gly</u>		<u>Pro</u>	<u>Asp</u>	<u>Glu</u>	<u>Asn</u>	<u>His</u>	<u>Gln</u>	<u>Lys</u>	<u>Arg</u>	<u>Ser</u>	<u>Thr</u>	<u>Ala</u>	<u>Met</u>	<u>Val</u>	<u>Ile</u>	<u>Leu</u>	<u>Phe</u>	<u>Tyr</u>	<u>Trp</u>	<u>Cys</u>

The PAM family

- PAM matrices are based on global alignments of closely related proteins.
- The PAM1 is the matrix calculated from comparisons of sequences with no more than 1% divergence.
- Other PAM matrices are extrapolated from PAM1.

The BLOSUM family

- BLOSUM matrices are based on local alignments.
- BLOSUM 62 is a matrix calculated from comparisons of sequences with no less than 62% divergence.
- All BLOSUM matrices are based on observed alignments; they are not extrapolated from comparisons of closely related proteins.
- BLOSUM 62 is the default matrix in BLAST 2.0. Though it is tailored for comparisons of moderately distant proteins, it performs well in detecting closer relationships. A search for distant relatives may be more sensitive with a different matrix.



The relationship between BLOSUM and PAM substitution matrices. BLOSUM matrices with higher numbers and PAM matrices with low numbers are both designed for comparisons of closely related sequences. BLOSUM matrices with low numbers and PAM matrices with high numbers are designed for comparisons of distantly related proteins. If distant relatives of the query sequence are specifically being sought, the matrix can be tailored to that type of search.

Sequence Analysis: Which scoring method should I use?

Comparable Blosum and PAM Tables

<u>Blosum</u> Tables	(Entropy)	<u>PAM</u> Tables	(Entropy)
<u>Blosum</u> 90	(1.18)	<u>PAM</u> 100	(1.18)
<u>Blosum</u> 80	(0.99)	<u>PAM</u> 120	(0.98)
<u>Blosum</u> 60	(0.66)	<u>PAM</u> 160	(0.70)
<u>Blosum</u> 52	(0.52)	<u>PAM</u> 200	(0.51)
<u>Blosum</u> 45	(0.38)	<u>PAM</u> 250	(0.36)

Percent
Sequence
Identity
PAM Tables

43
38
30
25
20

The **entropy** as defined by **information theory** is the **average amount of information per position in a sequence alignment that is available to determine whether or not the sequences are homologous**. This amount of entropy is available only if the similarity scores used in the database search or alignment are matched for the appropriate degree of sequence divergence.

An Alignment Algorithm

If we had all the time in the world, we could just make all possible alignments, score them all, & choose the best. But realistically, that won't work, since even for **two 100 amino acid sequences**, there are **10^{59} possible alignments**. So, the following approach was developed.

The particular class of algorithm we'll use is called ***dynamic programming***, which refers to a set of algorithms that allow the optimal solutions to be found for problems that can be defined in a ***recursive manner***. That is, the **problems are broken into subproblems**, which are **in turn broken into subproblems**, etc, until the simplest subproblems can be solved. For sequence alignments, this sequential dependency takes a form where the choice of optimal alignment of a sequence of length n is found from the solution to the optimal alignment of a sequence of length $n-1$ plus the alignment of the n th symbol, and the optimal alignment of the $n-1$ case is a function of the $n-2$ case, and so on. **Dynamic programming was developed by Richard Bellman 40-50 years ago, but then “rediscovered” by biologists aligning sequences in the 1970's.**

There are 2 types of alignments that we could make: *global* and *local*

Global alignments will require a forced match between every symbol of one string with some symbol (or gap) of the second string, e.g.

ACGTTATGCATGACGTA

-C---ATGCAT----T-

Local alignments will correspond to the best matching subsequences (including gaps). For the above example, this corresponds to:

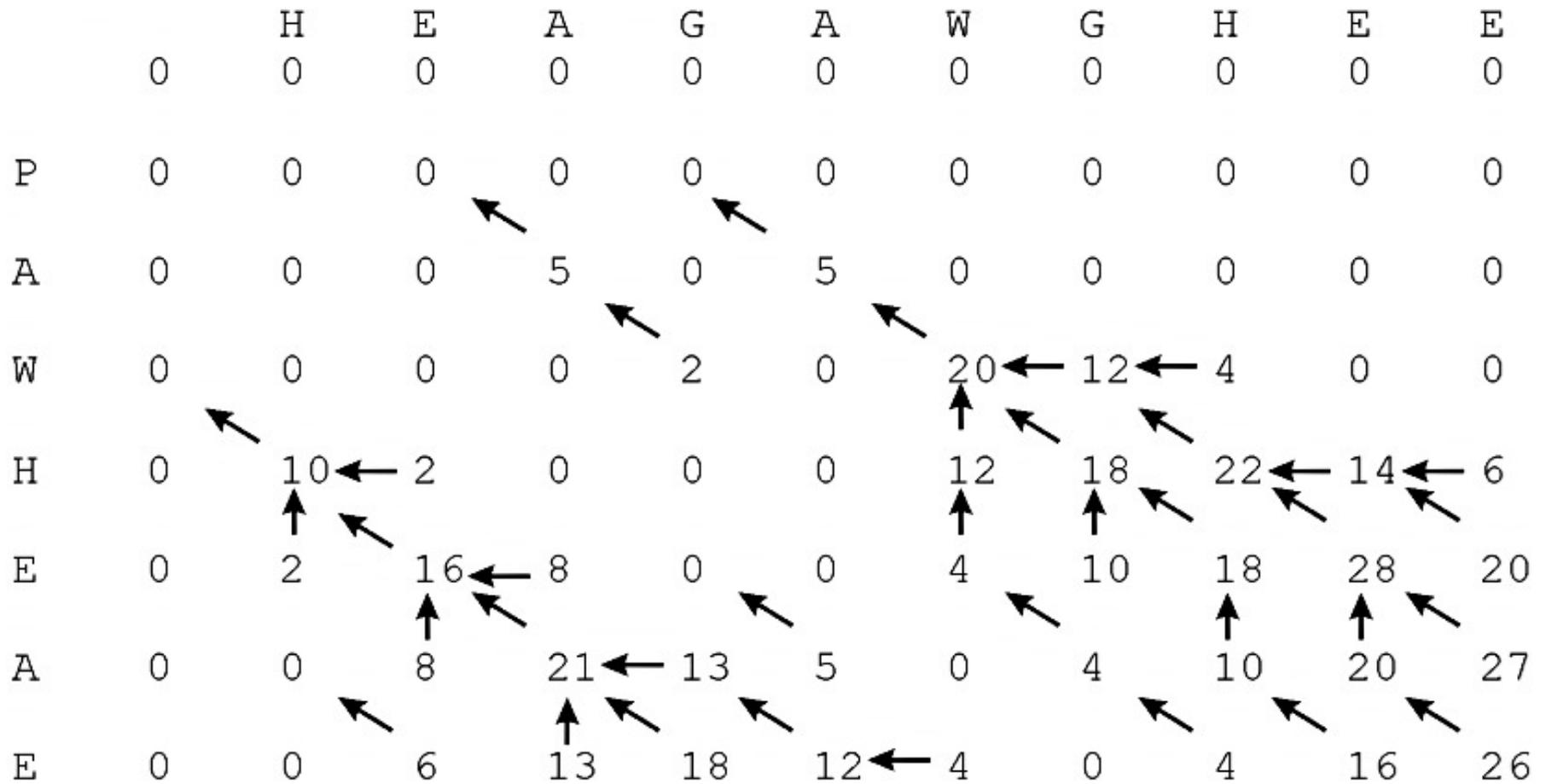
ATGCAT

ATGCAT

We'll look at **local alignments**, since these are what are used in almost any sequence alignment algorithm you might choose. This approach (in biology) is named the ***Smith-Waterman algorithm*** after Temple Smith & Mike Waterman, Journal of Molecular Biology vol. 147, 195-197 (1981).

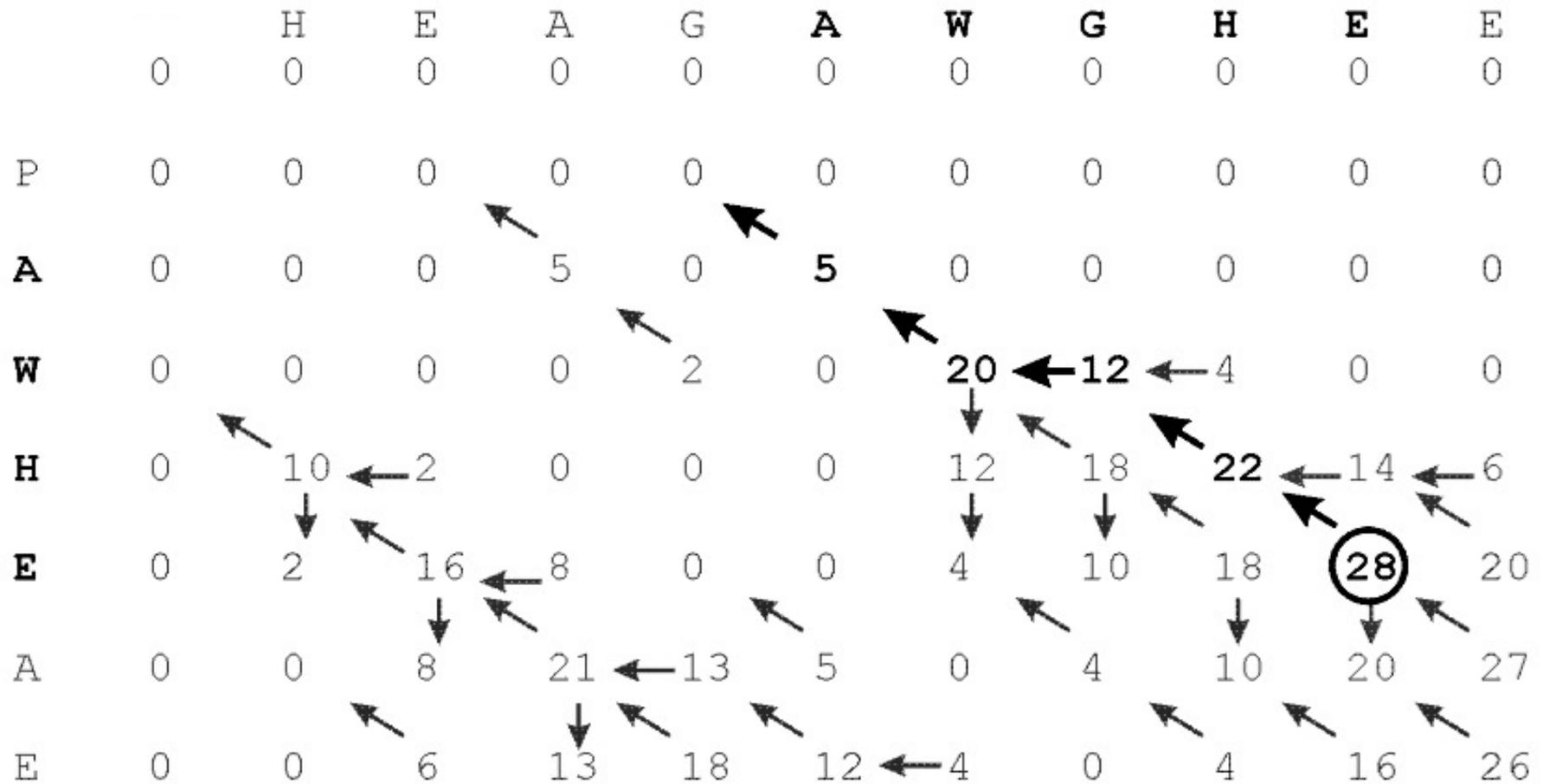
1) H E A G A W G H E E

2) P A W H E A E



1) H E A G A W G H E E

2) A W - H E



GAPS / Gap penalties

In most alignment and search programs, the **gap penalty** consists of **two terms**, the **cost to open the gap** and the **cost to extend the gap**.

Utility	Details
<u>FASTA3</u> , <u>BLAST2</u> , <u>CLUSTALW</u> , <u>ScanPS</u> and <u>MPsrch</u> .	GAOPEN or OPENGAP or OPEN GAP PENALTY : Penalty for the first residue in a gap (e.g. fasta defaults: -12 by with proteins, -16 for DNA). GAPEXT or EXTENDGAP or EXTEND GAP PENALTY : Penalty for additional residues in a gap (e.g. fasta defaults: -2 with proteins, -4 for DNA).

Ref: <http://www.ebi.ac.uk/clustalw/#>

Examples of aligned protein sequences:

Shown are 3 pairs of sequences, showing aligned sequences of proteins named FlgA1, FlgA2, FlgA3, and HvcPP. Between each pair the perfect matches and close matches (shown by + symbols, indicating chemically similar amino acids) are written.

Two biologically related proteins with similar sequences:

FlgA1 EAGNVKLRGRLDTLPPRTVLDINQLVDAISLRDLSPDQPIQLTQFRQAWRVKAGQQRVNVIASGD
++K+K+GRLDTLPP +L+ N A+SLR ++ QP+ R+ W +KAGQ V V+A G+

FlgA2 TLQDIKMKQGRDLTLPPGALLEPNFAQGAVSLRQINAGQPLTRNMLRRLWI KAGQDVQVLALGE (186)

Also biologically related (& fold up into the same 3D protein structure):

FlgA1 EAGNVKLRGRLDTLPPRTVLDINQLVDAISLRDLSPDQPIQLTQFRQAWRVKAGQQRVNVIASGD
A + P +L I+ R L P + I R+AW V+ G V V

FlgA3 LAALKQVTTLIAGKHKPDAMATHAEELQGKIAKRTLLPGRYIPTAAIREAWLVEQGA AVQVFFIAG (50)

But these are biologically unrelated (& fold up into unrelated structures):

FlgA1 AGNVKLRGRLDTLPPRTVLDINQLVDAISLRDLSPDQPIQLTQFRQA-WRVKAGQQRVNVIASGD
AG+V K G + + PRT ++ I+ P PI +++A WRV A + V V+ GD

HvcPP AGHV--KNGTMRIVGPRTCSNVWNGTFPINATTTGPSIPI PAPNYKKALWRV SATEYVEVVRVGD (128)

The problem we face is how to distinguish the biologically meaningless match (FlgA1-HvcPP) from the biologically meaningful ones (FlgA1-FlgA2 and FlgA1-FlgA3)?

1) H E A G A W G H E E

2) A W - H E

How do we know when a score is “good enough”?

Two elements of aligning sequences:

scoring the alignments (by generating substitution matrices)

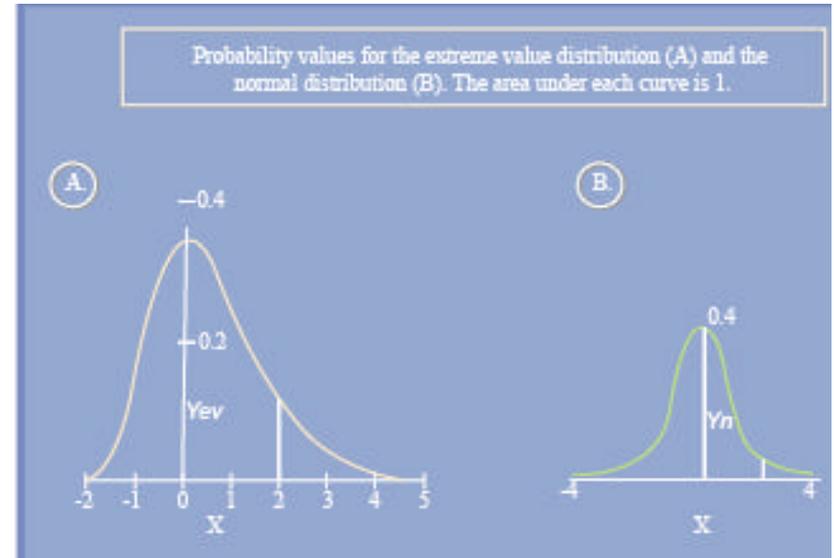
constructing the optimal scoring alignments by dynamic programming.

After we get an alignment, we have to decide if score is “good enough” to be significant. One way to this is to ask how hard it is to get that score from random alignments. Suppose we **“scrambled” one of the sequences, and found the best alignment with the other sequence.** The algorithm will always give us an alignment, even though the score is not very good. Still, let’s **do the scrambling and alignment process 1000 times.** If we look at those scores, and **never see a score as good as the real one, we can say that the real one has a 1 in a 1000 chance of happening just by luck.** If we did this 1,000,000 times and still didn’t see a score that good, we would begin to feel pretty confident in our alignment being significant.

Could do a **million random tests** after an alignment, and that should give a correct feeling for how good the alignment was. However, in practice, we can get away with just doing a **few random trials**, then mathematically modeling the scores we get out to save having to do a million such trials. The histogram of scores turns out to have a particular, predictable shape known as the **extreme value distribution** (also called the Gumbel distribution). Visually, the extreme value distribution looks this:

This distribution can be described by an equation of the form:

$$p(\text{max score} \leq X) \approx e^{-kNe^{\lambda(X-\mu)}}$$



where N is the number of scrambled y 's tested, μ is the mean value of the high scores from the scrambling experiment, and k and λ are numbers that characterize the shape of the particular **extreme value distribution** that comes from aligning x to y . In practice, k and λ can be fit from the scores from a few random alignments,

Multiple Sequence Alignments

For 3 sequences....

```
ARDFSHGLENKLLGCDSMRWE
GRDYKMALLEQWILGCD-MRWD
SRDW--ALIEDCMV-CNFRWD
```

An $O(mnj)$ problem !

Consider sequences each 300 amino acids

2 sequences – $(300)^2$
3 sequences – $(300)^3$
but for v sequences – $(300)^v$

Uh Oh !!!

Our polynomial problem
Just became exponential!

ClustalW

Higgins and Sharp 1988

- 1- Do pairwise analysis of all the sequences (you choose similarity matrix).
- 2- Use the alignment scores to make a phylogenetic tree.
- 3- Align the sequences to each other guided by the phylogenetic relationships in the tree.

New features: Clustal \boxtimes ClustalW (allows weights) \boxtimes ClustalX (GUI-based)

Weighting is important to avoid biasing an alignment by many sequence members that are closely related to each other evolutionarily!

Steps in doing a Multiple Sequence Alignment:

- 1) Get desired sequence in FASTA format.
- 2) NCBI web site – **BLAST** run
- 3) Select best matches to use in alignment
- 4) EMBL web site – **ClustalW** run

>CgX SEQUENCE

MPTYTCWSQRIRISREAKQRIAEAITDAHHELAHAPKYLQVIFNEVEPDSYFIAAQS
ASENHIWVQATIRSGRTEKQKEELLRLTQEIALILGIPNEEVVYITEIPGSNMTEY
GRLLMEPGEEEKWFNSLPEGLRERLLEGSSE

ClustalW Submission Form

Clustal W is a general purpose multiple sequence alignment program for DNA or proteins. It produces biologically meaningful multiple sequence alignments of divergent sequences. It calculates the best match for the selected sequences, and lines them up so that the identities, similarities and differences can be seen. Evolutionary relationships can be seen via viewing Cladograms or Phylograms. **New users, please read the FAQ.**

 [Download Software](#)

YOUR EMAIL	ALIGNMENT TITLE	RESULTS	ALIGNMENT	CPU MODE
<input type="text"/>	<input type="text" value="Sequence"/>	<input type="text" value="interactive"/>	<input type="text" value="full"/>	<input type="text" value="single"/>
KTUP (WORD SIZE)	WINDOW LENGTH	SCORE TYPE	TOPDIAG	PAIRGAP
<input type="text" value="def"/>	<input type="text" value="def"/>	<input type="text" value="percent"/>	<input type="text" value="def"/>	<input type="text" value="def"/>
MATRIX	GAP OPEN	END GAPS	GAP EXTENSION	GAP DISTANCES
<input type="text" value="def"/>	<input type="text" value="def"/>	<input type="text" value="def"/>	<input type="text" value="def"/>	<input type="text" value="def"/>

OUTPUT		PHYLOGENETIC TREE		
<input type="text"/>	<input type="text"/>	TREE TYPE	CORRECT DIST.	IGNORE GAPS

- Help Index
- General Help
- Formats
- Gaps
- Matrix
- References
- ClustalW Help
- ClustalW FAQ
- Jalview Help
- Scores Table
- Alignment
- Guide Tree
- Colours

4-OT– (Tautomerase/MIF Superfamily)

- with Professor Chris Whitman (Pharmacy)



Christian P.
Whitman

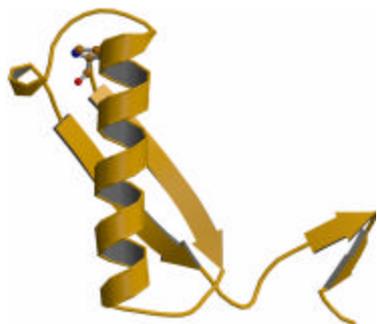
EEEEEEETT HHHHHHHHHHHHHHHHHHT GGG EEEEEEE GGG EETTEETTTT

4OT 1 PIAQIHILEG_RSDEQKETLIREVSEAISRSLDAPLTSVRVIITEMAKGHFGIGGELASKVRR 62

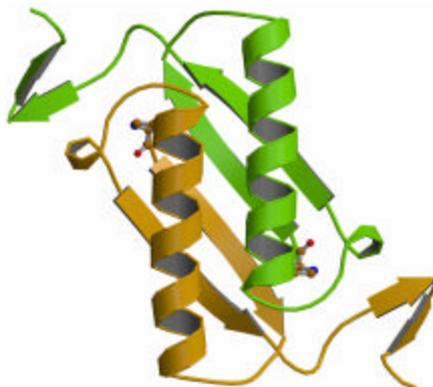
CHMI 1 PHFIVECSDNIREEADLPGLFAKVNPPTLAATGIFPLAGIRSRVHWVDTWQMADGQHDYAFVHM..-125

MIF 1 PMFIVNTNVP_RASVPEGFLSELTQQLAQATGK_PAQYIAVHVVPDQLMTFSGTNDPCALCSL..-114

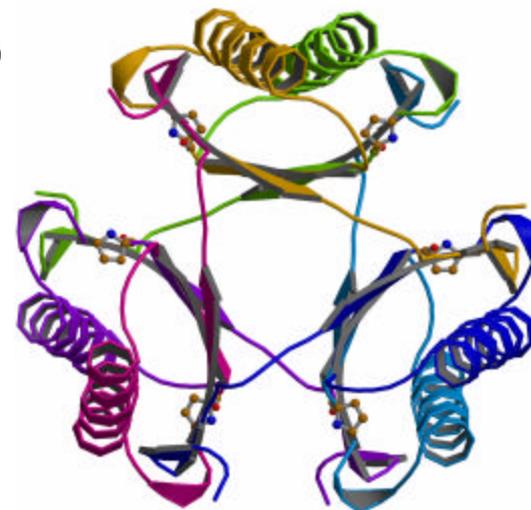
A)



B)



C)



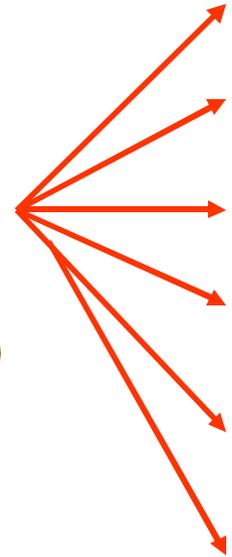
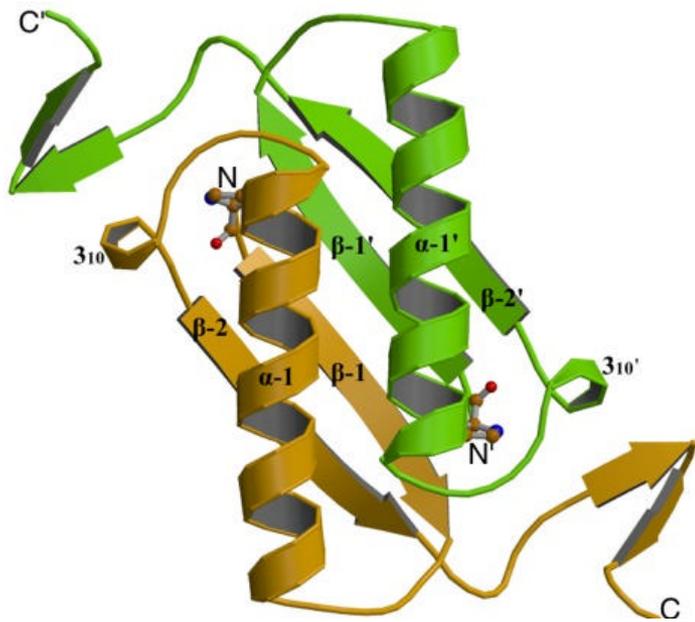
β - α - β Motif



New Activities



New Structures



4-OT - Tautomerase

4OT Homologues

CHMI - Isomerase

MIF - Cytokine / Hormone

Dehalogenase

Decarboxylase

a_6 a_3 a_2 $(ab)_3$

Sample Psi-BLAST Output

Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", Nucleic Acids Res. 25:3389-3402.
RID: 1012187428-16844-19639

Query= Pseudomonas putida - 4-OT (62 letters)

1 piaqihileg rsdeqketli revseaisrs ldapltsvrv iitemakghf giggelaskv rr

Database: All non-redundant GenBank CDS

translations+PDB+SwissProt+PIR+PRF

857,413 sequences; 270,034,499 total letters

Sequences with E-value BETTER than threshold

Round 1 - 30 Hits / Round 2 57 hits / Round 3 - 66 Hits

Sequences with E-value BETTER than threshold

	Score	E
Sequences producing significant alignments:	(bits)	Value
gi 6624277 dbj BAA88507.1 (AB029044) 4-oxalocrotonate isomerase...	<u>81</u>	2e-15
gi 16124116 ref NP_407429.1 (NC_003143) putative tautomerase [Y...	<u>78</u>	2e-14
gi 14715457 dbj BAB62059.1 (D85415) 4-oxalocrotonate tautomerases...	<u>78</u>	2e-14

gi 15642664 ref NP_232297.1 (NC_002505) 5-carboxymethyl-2-hydro...	<u>44</u>	3e-04
gi 15801678 ref NP_287696.1 (NC_002655) ydcE gene product [Esch...	<u>44</u>	3e-04
gi 16079011 ref NP_389834.1 (NC_000964) similar to hypothetical...	<u>43</u>	8e-04

Sequences with E-value WORSE than threshold

gi 15894207 ref NP_347556.1 (NC_003030) Protein related to MIFH...	<u>38</u>	0.014
gi 14600626 ref NP_147143.1 (NC_000854) MRSA protein [Aeropyrum...	<u>37</u>	0.047
gi 17562710 ref NP_506003.1 (NM_073602) macrophage migration in...	<u>35</u>	0.16

gi 5051891 gb AAD38354.1 (AF119571) macrophage migration inhibi...	<u>30</u>	4.4
---	-----------	-----

gi 14600626 ref NP_147143.1 (NC_000854) MRSA protein [Aeropyrum...	<u>30</u>	4.6
---	-----------	-----

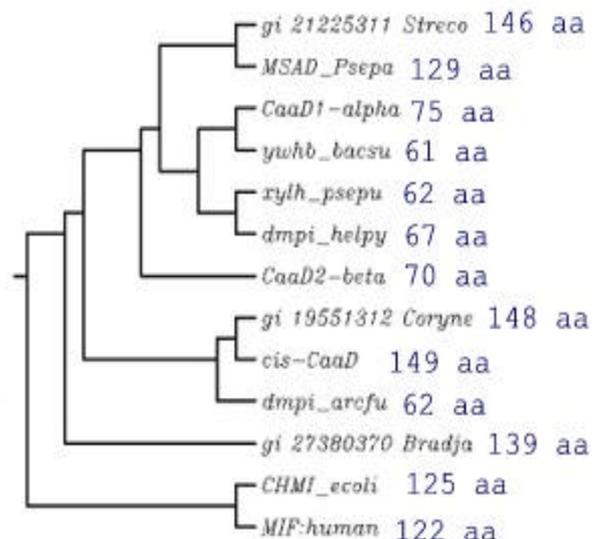
gi 5327268 emb CAB46354.1 (AJ012740) macrophage migration inhib...	<u>30</u>	8.1
---	-----------	-----

clustalw.aln

CLUSTAL W (1.83) multiple sequence alignment

```
gi|19551312|Coryne      PTYTCWSQRIRISREAKQRIAEAITDAHHELAHAPKYLQVIFNEVEPDSYFIAAQ--S-
cis-CaaD                PVYMVVYSQDRLTPSAKHAVAKAITDAHRGLTGTQHFLAQVNFQEQPAGNVFLGGV--Q-
CaaD1-alpha            PMISCDMRY-----
ywhb_bacsu             PYVTVKMLE-----
xylh_psepu             PIAQIHILE-----
dmpi_helpy            PFINIKLVPE-----
dmpi_arcfu            PVLIVYGPK-----
CaaD2-beta            PFIECHIAT-----
gi|21225311|Streco     PLITVSLRQGTTPQYRRLVSEALHKSMDVVKIPQDDQFHFVHEVTDNFMQPVV--FG
MSAD_Psepa             PLLKFDIFYGRDQAIKSLDAAHGAMVDAFGVPANDRYQTVSQHRPGEHVLEDTG--LG
gi|27380370|Bradja     PLITVSYTTSRQSPSLKADIASAVSELTAKILHKDPKVTAIIVKSWDAGDWFAGGRSLAE
MIF_human              PMFIVNTNVPRASVPDGFLELTQQLAQATGKPPQYIAVHVVPDQLMAFGG-----
CHMI_ecoli             PHFIVECSDNIREEADLPGLFAKVNPTLAATGIFPLAG-----
```

*



```
gi|19551312|Coryne      -ASENHIVVQATIRSGRTEKQKEELLLRLTQEIALILGIPNEEVWVYITEIPGSNMTEYG RLLMEPGEEEEKWFNSLPEGLRERLTELEGSSE-
cis-CaaD                -QGGDTIFVHGLHREGRSADLKGQLAQRIVDDVSVAEIDRKHIWVYFGEMPAQOMVEYGRFLPQPGHEGEWFDNLSSDERAFMETNVDVSRTEHL
CaaD1-alpha            -----GRTEQKRALSAGLLRVISEATGEPRENIFFVIREGSGINFVEHGEHLPDYVPGNANDKALIAKLK-----
ywhb_bacsu             -----GRTEQKRMLVEKVTEAVKETTGAEEKIVVFIEEMRKDHAVAVAGKRLSDME-----
xylh_psepu             -----GRSDEQKETLIREVSEAI SRSLDAPLTSVRVIITEMAKGHFGIGGELASKVRR-----
dmpi_helpy            -----NGGPTNEKQQLIEGVSDLMVKVLMKNKASIVVIIDEVDSNNYGLGGESVHHLRQKN-----
dmpi_arcfu            -----LDVGKKREFVERLTSVA AEIYGMDRSAITILIH EPPAENVGVGGKLIADRERE-----
CaaD2-beta            -----GLSVARKQQLIRDVIDVTNKSIGSDPKIINVLLVEHAEANMSISGRINGEAASTERTPAVS-----
gi|21225311|Streco     LRRTSRTLFIQLSFMRRGAEQKARLFRAIVANLRLYADVPEEDVMLVAFETARENWAAARRVVDPATGYDERMTDVPALPDVSEEPGR-----
MSAD_Psepa             YGRSSAVVLLTVISRPRSEEQKVC FYKLLTGALERDCGISPDVIVALVENS DADWSFGRGRAEFLTGDLV-----
gi|27380370|Bradja     QKLASYWIDIHVSEGTNTKDEKAA YLAAMFKRMAEILGPLHPETYLHVDEVKGDAYGFGGLTQERRYIAGKLEVALKAA-----
MIF_human              -SSEPCALCSLHSIGKIGGAQNR SYSKLLCGLLAERLRISPDRVYINYYDMNAANVGVNNSTFALEHHHHHH-----
CHMI_ecoli             IRRVHVVDVTQMQADGQHDYASVHMTL KIGAGRSLESRQQAGEMLFELIKTHFAALMESRLLALSFEIEELHPTLNFKQNNVHALFK-----
```