

BCH441-BIOINFORMATICS

SEQUENCE ALIGNMENT



---

BORIS STEIPE

DEPARTMENT OF BIOCHEMISTRY – DEPARTMENT OF MOLECULAR GENETICS  
UNIVERSITY OF TORONTO

section

FAST

ALIGNMENT:

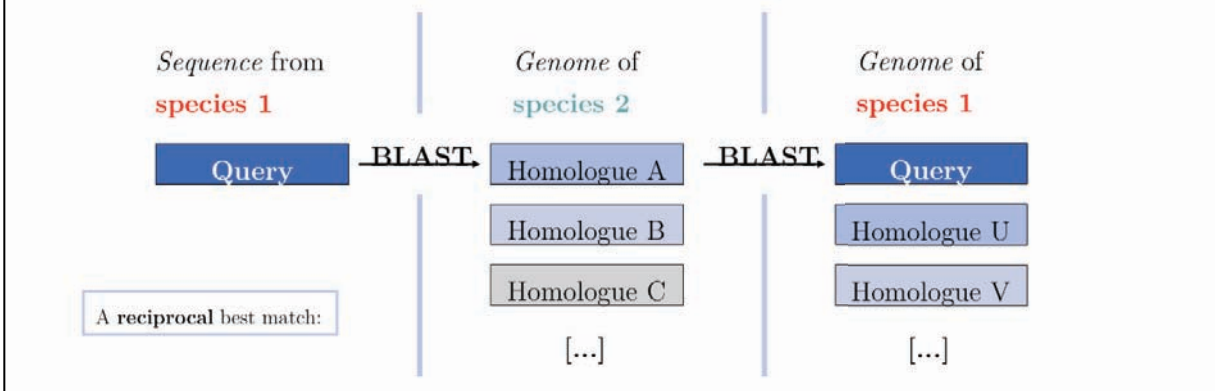
BLAST

why large search spaces?

Some procedures require **genome-wide** or **database-wide** similarity searches. Such searches are not feasible with *optimal sequence alignment* algorithms.

### Example: discovering orthologues

Computationally, an orthologue has been defined when the best match in another specie's genome has the original sequence as it's own best match. This is called the "**reciprocal best-match**" criterion. (*This is a useful procedure, but not the correct definition of an orthologue.*)



An example for searches in very large search spaces is the computational definition of orthologues.

The tight integration of search capacity with database holdings is the key to the utility of the data. Public investments in sequencing **only** pay off when the sequences are easily accessible!

Some computational definitions absolutely require genome-wide searches: e.g. the computational definition of orthologues, or compiling evolutionary conservation patterns.

BLAST

## Basic Local Alignment Search Tool

BLAST encompasses many different implementations and enhancements to a search algorithm that finds "**High Scoring Pairs**" of sequence alignments in databases.

It is a **Fast** way to find similar sequences.

It is **heuristic**, not exact, not optimal.

It is **not** the most **sensitive** way to search.

It is by a wide margin the **most commonly used tool** in bioinformatics.

BLAST was developed as a heuristic alternative to exact alignment, looking for a way to compute much faster, repeated searches in large search spaces than what one can do with pairwise alignments. The strategy is to pre-compute similarities and then piece a match together from quickly retrieved partial matches.

## BLAST principle

### 1: Query Preprocessing

Break query into words

D**TLV**RAIP -> DTL, **TLV**, LVR, VRA ...

Make a table of similar words

TLV -> TLI, **TIV**, SLV

### 2: Search query words in indexed table of database words

Find exact match between table word & db.

**TIV**

SDTDGDKNADGWIE**TIV**RALPTSD

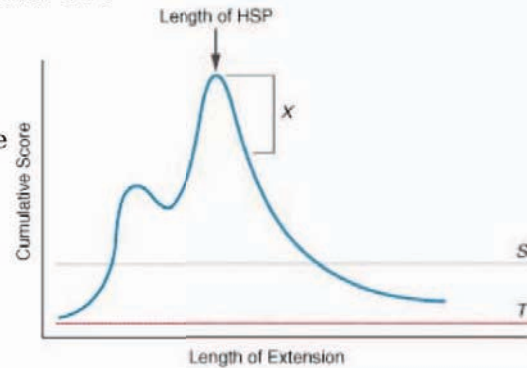
### 3: Extend query match to HSP

- keep HSPs of significant quality.

**DTLVRAIP**

SDTDGDKNADGWIE**DTLVRAIP**RALPTSD

### 4: Assemble HSPs into gapped local alignment



**HSP (High Scoring Segment Pair):**  
An ungapped, high-scoring, local alignment

The enormous speed-up of BLAST is due to its use of an **indexed table** of database "words". The index is a list of positions at which each word occurs in the database. Using an index, it is very easy to examine every occurrence of a word in the database and try to extend the word match on both sides with additional similar sequence. The extension does not introduce gaps, because this is faster, but also because the statistics of ungapped alignments are tractable! The final step is the assembly of significant hits into longer alignments.

Note that BLAST is **heuristic**, not **optimal** and that it is a **local**, not **global** alignment algorithm.

See also: Altschul *et al.* (1990): <http://www.ncbi.nlm.nih.gov/pubmed/2231712>

BLAST

BLAST: Basic Local Alignment Search Tool

blast.ncbi.nlm.nih.gov/blast.cgi

BLAST\* Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

NCBI BLAST Home

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

**DELTA-BLAST, a more sensitive protein-protein search**

BLAST Assembled Genomes

Find Genomic BLAST pages:

[GO](#)

Human Rabbit Zebrafish  
 Mouse Chimp Clawed frog  
 Rat Guinea pig Arabidopsis  
 Cow Fruit fly Rice  
 Pig Honey bee Yeast  
 Dog Chicken Microbes

Basic BLAST

Choose a BLAST program to run.

nucleotide blast Search a nucleotide database using a nucleotide query  
 Algorithms: blastn, megablast, discontiguous megablast

protein blast Search protein database using a protein query  
 Algorithms: blastp, psi-blast, phi-blast, delta-blast

blastx Search protein database using a translated nucleotide query

tblastn Search translated nucleotide database using a protein query

tblastx Search translated nucleotide database using a translated nucleotide query

Specialized BLAST

Choose a type of specialized search (or database name in parentheses.)

- Make specific primers with [Primer-BLAST](#)
- Cluster multiple sequences together with their database neighbors using [MOLE-BLAST](#)
- Find conserved domains in your sequence (cds)
- Find sequences with similar [conserved domain architecture](#) (cdart)
- Search sequences that have [gene expression profiles](#) (GEO)
- Search [immunoglobulins and T cell receptor sequences](#) (IgBLAST)
- Screen sequence for [vector contamination](#) (vecscreen)
- [Align](#) two (or more) sequences using BLAST (bi2seq)
- Search [protein or nucleotide targets](#) in PubChem BioAssay
- Search [SRA by experiment](#)
- [Constraint Based Protein Multiple Alignment Tool](#)
- [Needleman-Wunsch Global Sequence Alignment Tool](#)
- Search [RefSeqGene](#)
- Search [trace archives](#)

Your Recent Results [New!](#)

[All Recent results...](#)

News

[Find Genomic BLAST pages](#)

You can now find Genomic BLAST pages using the search box from the BLAST homepage.

Thu, 02 Oct 2014 11:00:00 EST

[View BLAST news...](#)

Tip of the Day

[View tips...](#)

blastn

blastp

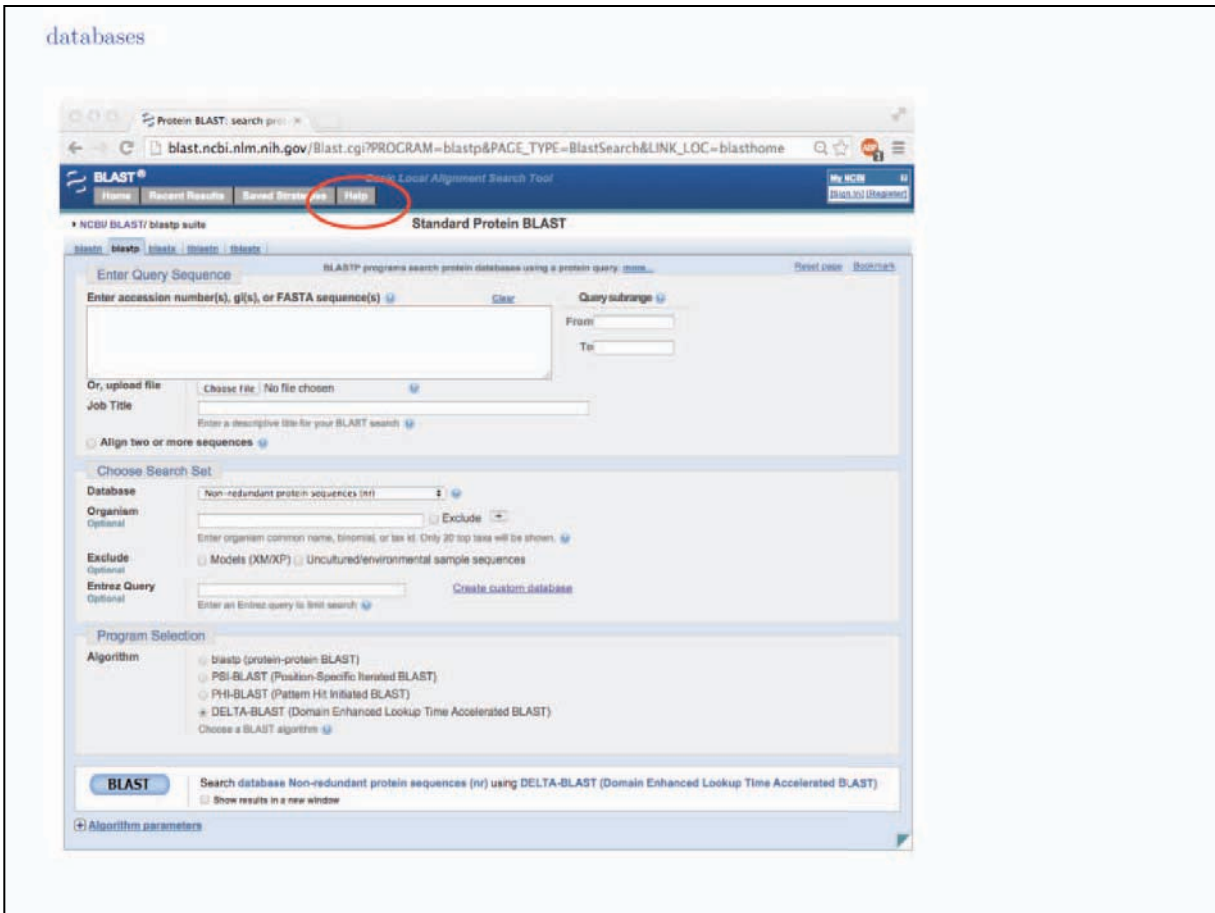
blastx

tblastn

tblastx

PSI-BLAST  
PHI-BLAST  
DELTA-BLAST

The **BLAST** home page offers a number of different BLAST *flavours*.

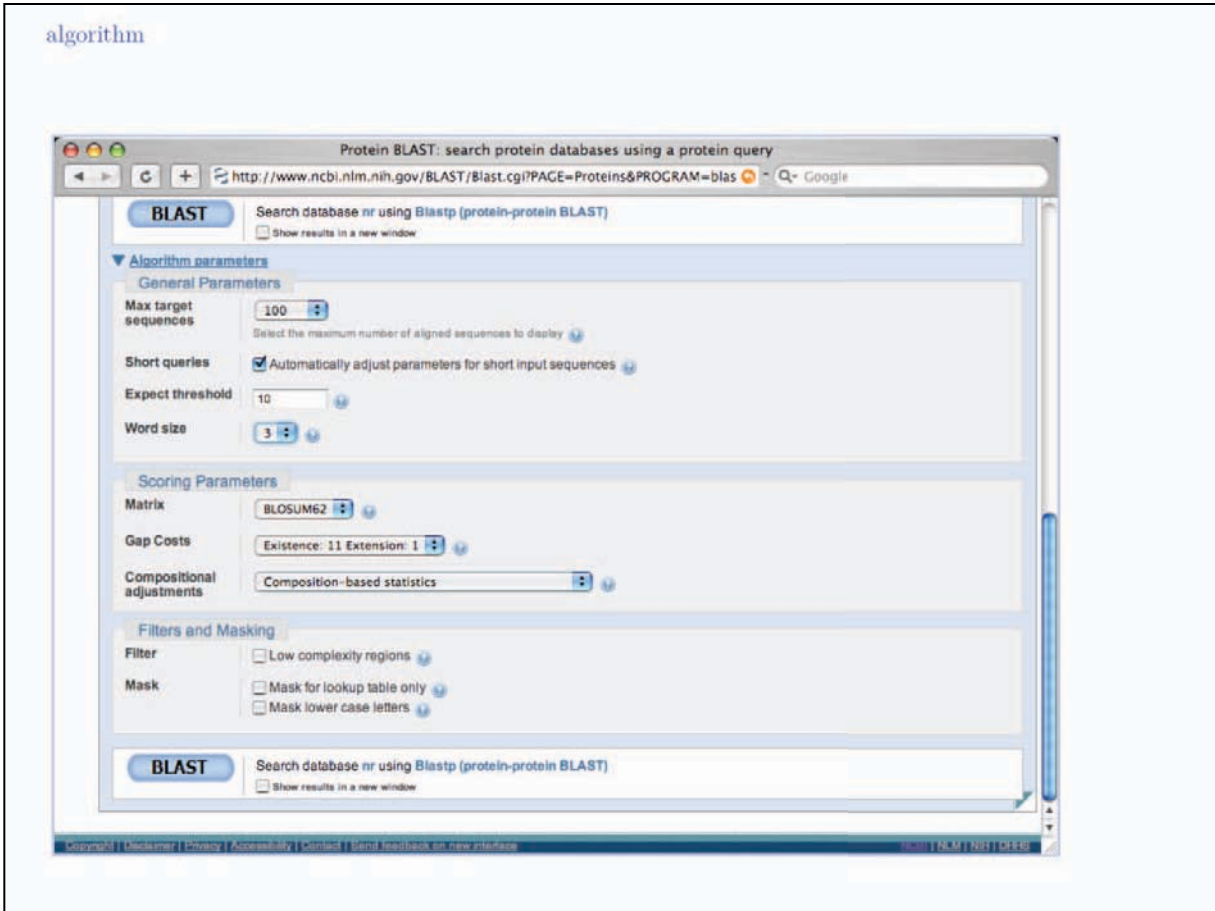


## BLAST parameters: databases

Extensive help is available (and should be read!) for each of the options. Take the time to read the **Web BLAST options document**<sup>1</sup> and be sure to understand how to format input, what databases are available and how the choice of database influences the results. If you are not confident with the document, ask on the course list.

For example, the Help page contains guides to the **search interface** and the **report output**!

<sup>1</sup> <http://www.ncbi.nlm.nih.gov/BLAST/blastcghelp.shtml>



## BLAST parameters: algorithm

Be sure to understand the choices and their consequences for **Composition-based statistics**<sup>1</sup> and for **Filtering and Masking** segments of low complexity in your query. Filtering is an important option to consider especially for PSI-BLAST searches!

<sup>1</sup> [http://www.ncbi.nlm.nih.gov/BLAST/blastcghelp.shtml#compositional\\_adjustments](http://www.ncbi.nlm.nih.gov/BLAST/blastcghelp.shtml#compositional_adjustments)

<sup>2</sup> <http://www.ncbi.nlm.nih.gov/BLAST/blastcghelp.shtml#filter>



results

All identifiers reported for the same hit represent different database versions of the same sequence

Sequence numbers are given relative to the query-string resp. the database entry, excluding gaps

A "hit"

```
>gi|39998047|ref|NP_953998.1| redox-active disulfide protein 2 [Geobacter sulfurreducens PCA]
gi|39984992|gb|AAR36348.1| redox-active disulfide protein 2 [Geobacter sulfurreducens PCA]
Length = 78

Score = 58.5 bits (140), Expect = 2e-08
Identities = 31/76 (40%), Positives = 54/76 (71%)

Query: 1 MMKIQIYGTGCANCQMLEKNAREAVKELGIDAEFEKIKEMDQILEAGLTALPGLAVDGEL 60
+MKI++ GTGCA C+ L +N ++AV+ G +AE K++E+ +I++ G+ + P L +DG +
Sbjct: 3 IMKIEVLGTGCAKCKTLYENVQKAVEMSGKEAEVVRVVEIQKIMKYGVMSPTALVIDGVV 62

Query: 61 KIMGRVASKEEIKKIL 76
K G+V + +EIK +L
Sbjct: 63 KFSGKVPAADEIKGML 78
```

"Query" is the sequence you searched with

"Sbjct" is the sequence that BLAST found

Each Blast **hit** represents an alignment that can contain one or more HSPs (High Scoring Segment Pairs). Note: If a hit is followed by a second hit and no new GI number, it identifies a second region of similarity in the **same sequence**.

## E-value

The **quality** of the alignment is represented by the Score ( $s$ ).

The **significance** of the alignment is computed as an E- value.

### **E-value (E)** *Expectation value:*

The number of alignments with scores equivalent to, or better than a score  $s$  that are expected to occur in a database **of the same size** that does not contain a homologous sequence.

The smaller the E-value (the larger its negative exponent), the more significant the score.

The E-value is a statistically well founded metric that allows us to conclude the likelihood of a spurious alignment. Computing E-values is possible for HSPs since the statistics of gap-less alignments are analytically tractable, whereas gapped alignments have no theoretical description of the distribution of expected scores.

Note that E-values do not represent an assertion about the retrieved sequence, but an assertion about the score and its relation to the expected distribution of scores. Or, to rephrase this, a large E-value does not mean that your hit is not a homologue, but it means that an irrelevant sequence has a high chance of having just as high a score due to chance similarities. To repeat: a large E-value does not mean your hit is not a homologue. However a small E-value does indeed mean that a chance alignment is unlikely.

It is important to realize that the E-value depends on the database size. Obviously, you would expect randomly high-scoring hits more often in a large database than in a small one. Thus an alignment with the **same score** will have a **smaller E-value** when searched against a particular genome than if you search it against the entire "nr" dataset of GenBank.

More detail in the NCBI tutorial "The Statistics of Sequence Similarity Scores" (<http://www.ncbi.nlm.nih.gov/BLAST/tutorial/Altschul-1.html>)

interpretation

E-values are very convenient because they have an obvious interpretation of significance, **but they do not absolve you from using biological common sense!**

**Example: searching pea defensin (1JKZ) against nr: are these sequences homologous?**

```
gi|15226880|ref|NP_178322.1|   plant defensin protein, putative (PDF2.6)
gi|11387216|sp|Q9ZUL8|THG4_ARATH   Gamma-thionin homolog At2g02140 precursor
gi|25330850|pir|D84433   proteinase inhibitor II [imported] - Arabidopsis thaliana
gi|4038038|gb|AAC97220.1|   protease inhibitor II [Arabidopsis thaliana]
gi|21592674|gb|AAM64623.1|   protease inhibitor II [Arabidopsis thaliana]
      Length = 73

Score = 30.8 bits (68), Expect = 6.7
Identities = 14/46 (30%), Positives = 27/46 (58%), Gaps = 1/46 (2%)

Query: 1  KTCEHLADTYRGVCFNASCDDHCKNKAHLISGTCHNWKCFCTQNC 46
      ++ ++ ++GVC + SC  C ++      G C + +C+C++ C
Sbjct: 29 RTCESPSNKFQGVCLNSQSCAKACPSEG-FSGGRCSSLRCYCSKAC 73
```

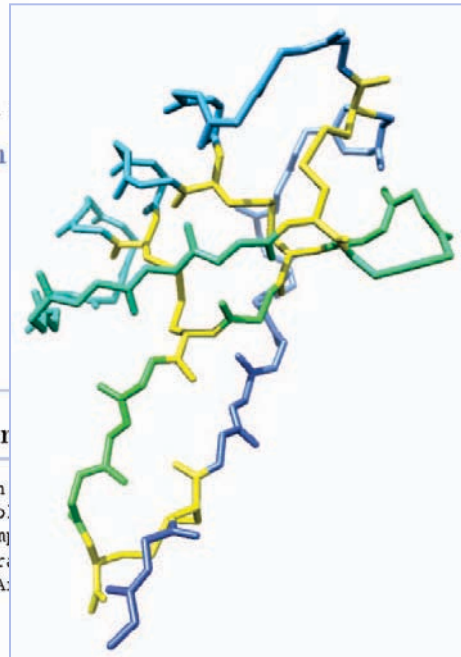
In the example above, the BLAST search of a Pea defensin - PDB structure 1JKZ - achieved an E-value of only 6.7.

interpretation

E-values are very convenient because they have an significance, **but they do not absolve you from sense!**

Always:

- Check the annotation;
- Check the alignment.



Example: searching pea defensin (1JKZ) against r

```
gi|15226880|ref|NP_178322.1| plant defensin protein
gi|11387216|sp|Q9ZUL8|THG4_ARATH Gamma-thionin homo
gi|25330850|pir|D84433 protease inhibitor II [im
gi|4038038|gb|AAC97220.1| protease inhibitor II [Ar
gi|21592674|gb|AAM64623.1| protease inhibitor II [A
Length = 73
```

Score = 30.8 bits (68), Expect = 6.7

Identities = 14/46 (30%), Positives = 27/46 (58%), Gaps = 1/46 (2%)

```
Query: 1 KTCEHLADTYRGVCFNASCDHCKNKAHLISGTCHNWKCFCTQNC 46
      ++ ++ ++GVC + SC C ++ G C + +C+C++ C
Sbjct: 29 RTCESPSNKFQGVCLNSQSCAKACPSEG-FSGGRCSSLRCYCSKAC 73
```

However the hit that was retrieved:

- is annotated as an *arabidopsis* **defensin**;
- has **30% identity** over the entire domain, albeit the domain is small;
- requires only **one single gap** for alignment; and
- **has each and every single cysteine conserved, when compared to the query!**

Each of these additional observations alone could have led you to conclude homology. It should be obvious for example that the aligned cysteines are extremely unlikely to be due to a random similarity of unrelated sequences! The large E-value is primarily due to the fact that the protein sequences are quite short.

## Too many ?

- **Restrict the database to RefSeq** (best representative, non-redundant sequence), or restrict the search to particular organism(s).
- **Search on a substring** of the sequence (e.g. search for a domain if you can define one through RPS-BLAST or SMART). This will suppress, smaller, non-specific results.
- **Increase the number of hits** that are being reported. If you don't do that, relevant hits may be dropping off the end of the results page.
- **decrease E-value** (smaller value means more stringent threshold for reporting hits) but be aware of the potential to loose interesting hits from the "twilight zone".

How can there be too many hits, when *lots-of-hits* is what you are looking for? Either you find redundant sequences or trivially similar sequences that are obscuring the rare, interesting similarities you are looking for (GFP or other fusion proteins and ankyrin domains come to mind, for example), or you are searching in a database section that contains redundant sequences.

Note that restricting by organism does not restrict the search, but only the list of results that are being reported. The search takes just as long. Only the specialized genome search pages and some non-NCBI databases of model-organism genome projects offer BLAST searches on reduced datasets. These searches are faster.



too few hits

## Too few?

- Search with **domains**, rather than full-length proteins (more sensitive)
- remove database restrictions (e.g. search **nr** instead of RefSeq)
- raise the **E-value**: 10, 100, ... (but expect irrelevant hits that may be difficult to verify)
- change the scoring matrix to BLOSUM45 instead of BLOSUM62 (always a good idea when you are looking for distant relationships)
- search additional databases (e.g. use **tblastn** to search EST or genomic data, this may identify frameshifts or inadequate annotations)
- **use a more sensitive algorithm: PSI-BLAST**

But don't expect miracles. Many genes/proteins simply do not have significant non-trivial database matches.

How many genes have no homologues? That depends. Unknown genes (or "ORFans") may comprise a significant (albeit diminishing) fraction of genomes.

In general, between 10 and 30% of sequences may fall into this category and it is likely that even the most closely related species have sequences that are unique.

See Siew&Fischer (2003)<sup>1</sup> and a discussion of the role of viral horizontal gene transfer in ORFans by Yin and Fischer (2006)<sup>2</sup>

<sup>1</sup> <http://www.ncbi.nlm.nih.gov/pubmed/12517334>

<sup>2</sup> <http://www.ncbi.nlm.nih.gov/pubmed/16914045>

Searching with profiles instead of individual sequences is more sensitive!

It reduces the effect of *random drift* in individual sequences and focuses the search on the *commonalities* and the specific *tolerance to variation* that is particular to a *family* of sequences (in which all sequences have diverged from a common ancestor).

## **PSI-BLAST proceeds in five steps:**

1. Select a query and BLAST it against a protein database
2. PSI-BLAST constructs a multiple sequence alignment from the BLAST hits, then creates a "profile" (or position-specific scoring matrix (PSSM))
3. The PSSM is used as a query against the database
4. PSI-BLAST estimates statistical significance (E values) and proposes significant hits for inclusion into the next iteration
5. Repeat steps [3] and [4] iteratively, typically 5 times. At each new search, a new profile is constructed from the previous hits and used as the new query.



PSI-BLAST is useful to detect weak but biologically meaningful relationships between proteins.

Not all sequences can be used for a search. For example, a query with a coiled-coil motif may detect thousands of proteins that have this motif but that are not homologous.

False positives in the search can arise from the inclusion of irrelevant sequences into the profile.

Once even a single irrelevant protein sequence is included in a PSI-BLAST profile, it will not go away.

If the number of related irrelevant proteins is large, they will *take over* the profile: **profile corruption!**

profile corruption

**Apply *filtering* of regions with low-complexity and biased composition**

**Adjust *E-value* from 0.001 (default) to a more stringent value (e.g. 0.0001).**

**Visually inspect the output from each iteration:**

Suspicious hits have irrelevant locations or functions, similarities only to parts of domains, fail to conserve important motifs or disulfide bridges or have poor-quality alignments.

**Remove suspicious hits from inclusion by unchecking the box.**

Be conservative regarding the sequences you include, **true positives will gradually improve their *E-values*** with subsequent iterations, even if they are not included in the profile. You can simply include them in later iterations (or not at all, they will still be reported, even if they don't contribute to the profile). False positives will not improve significantly.

In the end, how many false positives can we expect? Unfortunately, more than we'd think. Jones & Swindells (2002)<sup>1</sup> have run an analysis against decoy sequences that picked up false positives in 5% of all cases, after the fifth iteration, although the *E-value* threshold was set to 0.001.

Even though their methodology was a bit *ad hoc* and finding false positives about 50 times more frequently than expected is not catastrophic, we must realize that protein sequences are not random strings and that significance is often hard to evaluate, because it is hard to get the *null* hypothesis right. Use caution, use common sense and in questionable cases try to use independent confirmation of homology, such as conserved binding sites or functional motifs, if possible.

<sup>1</sup> <http://www.ncbi.nlm.nih.gov/pubmed/11893514>

## filtering

Filtering is used to mask parts of sequences with **low complexity** for database searches, by replacing sequence fragments with "X" (unknown residue) symbols.

Filtering is applied to the QUERY not the database

Filtering can eliminate statistically significant but biologically uninteresting reports from BLAST output (e.g. Basic, acidic, or proline rich regions, but these **may** actually be interesting, depending on your biological question.)

Personally, I turn filtering OFF by default for standard BLAST searches, only turn it ON if low-complexity regions appear to cause problems with specificity. Informed judgement is required.

However: always use filtering for PSI-BLAST (profile corruption)!

importance

To ...

... identify regions of the polypeptide chain that fold independently,  
that are stable on their own

*(folding units; initiation sites for folding)*

... identify gene fusion or gene insertion events  
from analysis of the 3D structure

*(understand evolutionary history)*

... understand protein mechanism as an additive/cooperative result  
of domain function

*(CDART, SMART - domain architecture)*

... allow for meaningful structural classification of proteins

*(SCOP, CATH classifications)*

section

# DOMAINS

The discovery and catalogization of the universe of domains in protein sequences is the greatest achievement of **profile- and model based sequence alignment algorithms.**

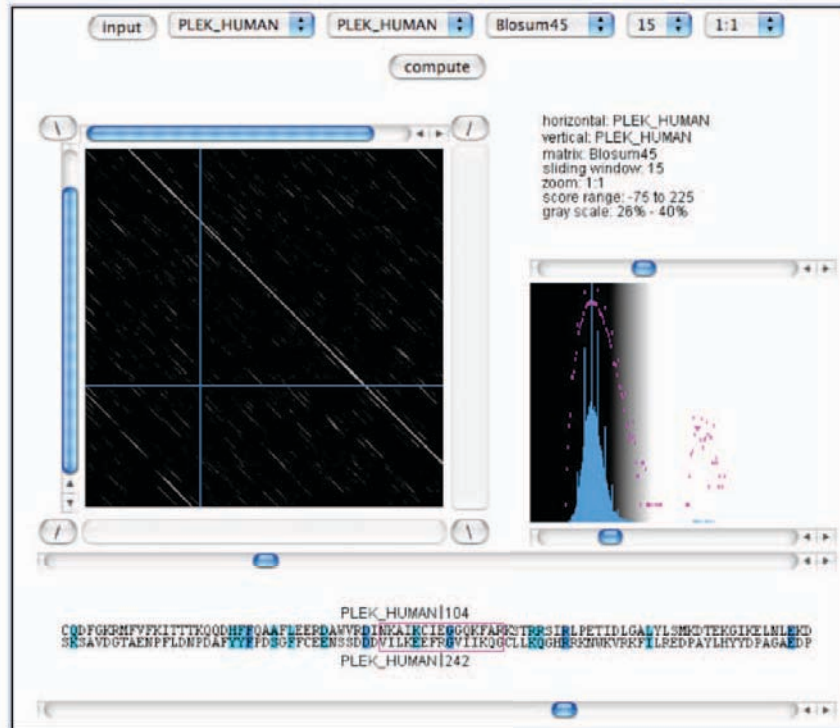
Domains are modules for ...

- inheritance
- distinct function
- separate folding

The *domain* is the natural unit of analysis of protein structure.

## Dotlet -

A dotplot of Pleckstrin (p47) reveals similarity between N-and C terminus !



Here is an example of how a domain was discovered from sequence alignments.



units of inheritance

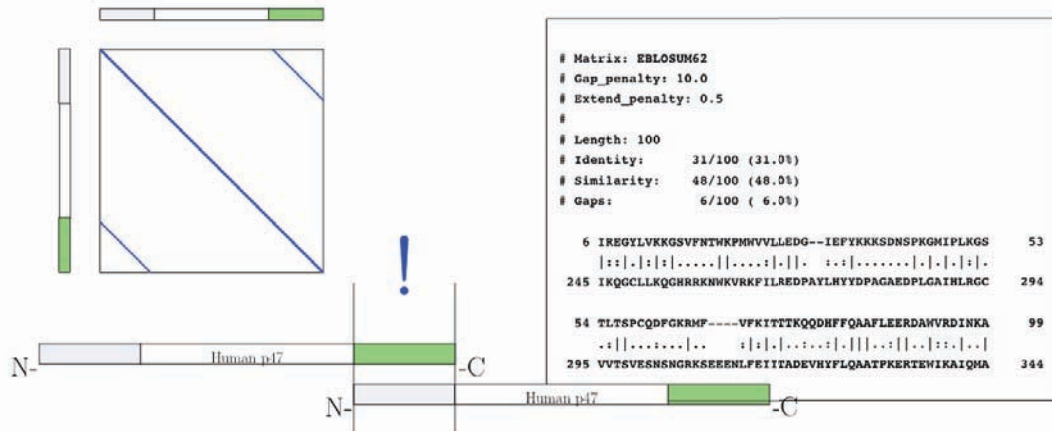
```
# Matrix: EBLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
#
# Length: 100
# Identity:      31/100 (31.0%)
# Similarity:    48/100 (48.0%)
# Gaps:          6/100 ( 6.0%)

  6 IREGYLVKKGSVFNTWKPMWVVLLEDG--IEFYKKKSDNSPKGMIPKGS      53
   |:|. |:|:|. . . . . | | . . . . :|. | | .   :.:|. . . . . | . | . | :|.
245 IKQGCLLKQGHRRKNWKVRKFI LREDPAYLHYYPDAGAEDPLGAIHLRGC    294

 54 TLTSPCQDFGKRMF-----VFKITTTKQDHFQAAFLEERDAWVRDINKA    99
   :|:|. |. . . . . :|. | | | . . . | | . | :.:|. | . |
295 VVTSVESNSNGRKSEENLFEIITADEVHYFLOAATPKERTEWIKAIQMA    344
Optimal sequence alignment: 31% identity over ~100 amino acids.
```

The alignment shows high similarity between N- and C-terminus.

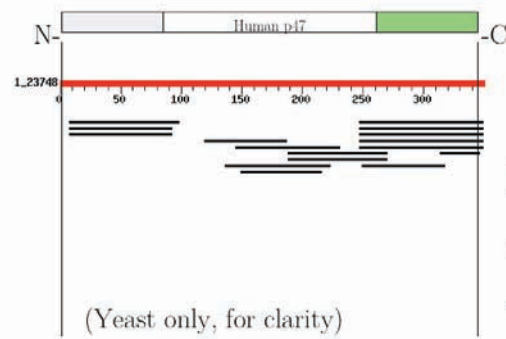
units of inheritance



Overlapping alignments may define domain boundaries ! We can search a database with this knowledge ...

improved specificity

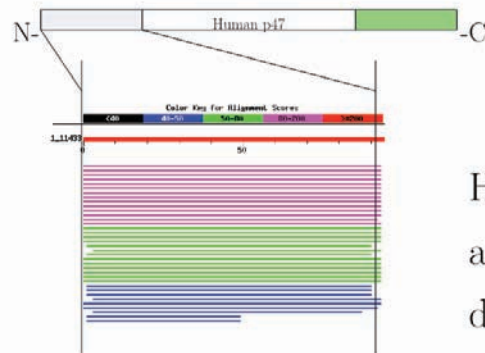
## BLAST search with *full-length* Pleckstrin sequence



Hits extend over the entire domain. **PSI-BLAST** would be difficult ...

improved specificity

## BLAST search with Pleckstrin Domain *only* ...



Hits are smoothly bounded and extend over the entire domain. PSI-BLAST is straightforward.

486 hits ... etc.

## Domain discovery:

- HMMER profiles
- Pfam and SMART databases

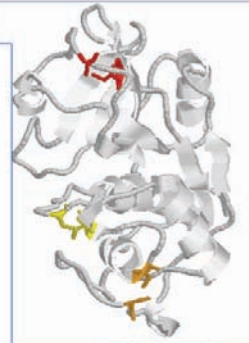
Domain discovery on a large scale has been made possible through Hidden Markov Model alignments, implemented in Sean Eddy's HMMER program. This has been used to compile large databases like Pfam that curate domain profiles. These profiles can be scanned against an unknown sequence, thus allowing the annotation of the sequence with the domains it contains. In many cases this allows to assign at least a coarse description of function and mechanism.

section

# MULTIPLE SEQUENCE ALIGNMENT

added value of MSAs

- Phylogenetic relationships
- Conservation patterns
  - Mutable regions
  - Conserved residues
  - Conserved properties
  - Conserved sequence patterns
  - Domain boundaries
  - [...]



	56	63	96
NLVD	CYSE...	ND.GGGGY...	SCM...
ELVDCER...	RSH.GGGGY...	TQR...	
ELVECTNG.QNS.G	NGGL...	KCD...	
CYSP_HEMSP...	ELVDCDR.S.YNE.G	CDGGL...	VCD...
CYSP_HEMSP...	ELVDCDKEE..NQ.G	CNGGL...	TCD...
CATL_DROME...	NLVDCT.KYGNN.G	CNGGL...	SCH...
CATJ_RAT...	NLLDTRSE...GI.G	LPWGT...	PCR...
ALEU_HORVU...	QLVDCAG.GFNNF.G	CNGGL...	VCH...
BROM			
EUM1			
CPR5			
CYS1			

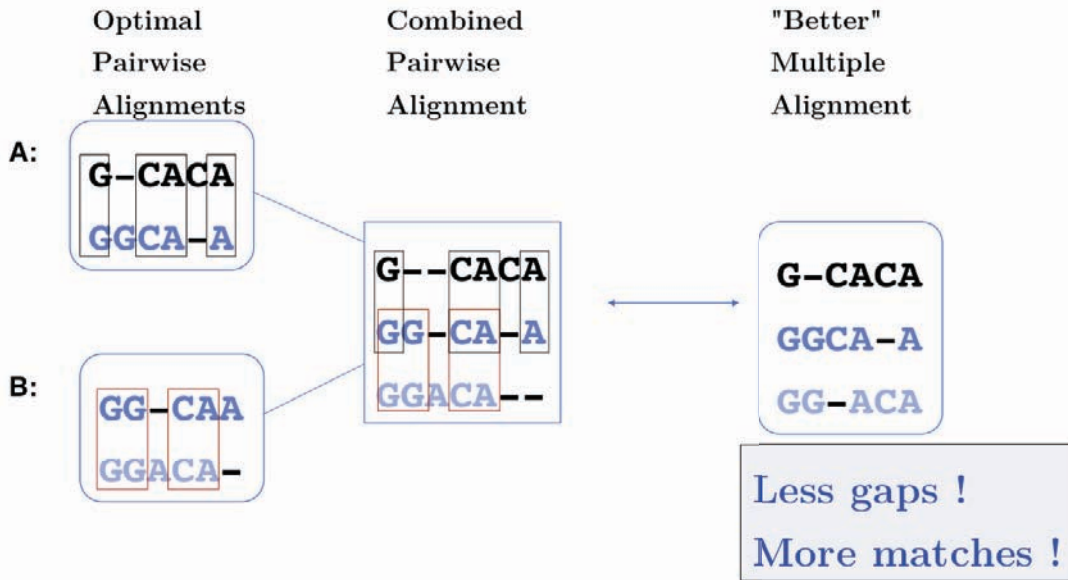
- Homology modelling
- Secondary structure prediction
- Phylogenetic reconstruction
- Sensitive homology searches
- [...]

MSAs show conservation patterns.

Multiple sequence alignments don't just match residues. They also give information on how strongly a residue is conserved, what it can be replaced with, which species share particular sequence patterns, and where in the sequence indels can be tolerated. An analysis of conservation even allows to distinguish between structurally and functionally conserved residues! This makes multiple sequence alignments the method of choice for many applications.

- Multiple sequence alignments are more accurate than pairwise alignments, thus they are the method of choice for starting **homology modeling** projects;
- Combined information from numerous sequences is invaluable for **secondary structure prediction** and **sensitive sequence database searches**;
- They contain the information needed for inferences about **evolutionary relationships**, i.e. the order in which particular sequence changes occurred.

## Optimal pairwise alignment $\neq$ best multiple alignment



Multiple alignments cannot necessarily be **constructed** by merging pairwise alignments. Moreover, it may be actually be impossible to merge three mutually pairwise alignments into a non-contradicting multiple alignment. However the inverse is always possible: a multiple alignment can be **decomposed** into pairwise alignments.



## Optimal multiple alignment is intractable for more than about ten sequences

Pairwise alignment:  $O(L^2)$  → 2D-Pathmatrix

$n$ -wise alignment:  $O(L^n)$  →  $n$ D-Path-hypercube

E.g. 20 APSES domains:  $L = 100$ ,  $n = 20$ :  $10^{21}$  operations

Despite more than thirty years of efforts, multiple sequence alignment is still a topic of current research in bioinformatics.

Besides being intractable, it is questionable how meaningful the objective function of *optimal sequence alignments* is for *multiple alignments*. This pair score maximizes the score derived from a mutation data matrix, for pairs of aligned residues. But – for example – the pair score does not optimize the pattern of indel placements, or whether a particular motif is well-conserved.

**"Objective function" of MSAs:**  
**the score that MSA algorithms try to maximize.**

---

- Many alternatives have been proposed
  - Most common: sum of scores of all pairwise alignments
  - Scores are not comparable across different alignments and across different algorithms
- 

How to define an objective function that identifies the biologically most meaningful MSA?

If we want an algorithm to optimize anything at all, we first must define how we can measure the quality of the result. This metric defines the **target function** or **objective function**.

(Note that "objective" here is not used in the sense of "unbiased" but in the sense of being a "target", or "goal".)

biologically meaningful?

## Biologically motivated objectives for multiple alignments:

- Minimize number of indels, not length
- Minimize number of sites at which indels are tolerated
- Maximize sequence similarity
- Retain conserved motifs and patterns
- Recapitulate phylogeny
- Concentrate on alignable regions not on gapped regions
- [...]

---

Alignment objectives are based on the **biological models** we apply to multiple alignments - they attempt to capture constraints on sequence similarity that go beyond optimizing a pairwise alignment score (which is based on a **context-independent** mutation data matrix and on an **empirical model** for the probability of indels).

*Reasonable* alignment metrics are based on models of how evolution has shaped a family of related sequences.

Each of the reasonable biological objectives suggests a different alignment strategy! The most modern algorithms currently available attempt to satisfy these heuristics simultaneously. Note that these are "heuristics", they are not the result of some rigorously applied theory, but reflect the complex relationship between protein sequence, structure, evolution and selection.

Computational strategies for multiple alignment algorithms derive from the nature of the objective function. The objective function can reflect biological heuristics.

<b>Objective</b>	<b>Algorithm</b>	<b>Type of alignment algorithm</b>
<b>Maximize similarity, Minimize gaps</b>	Bounded optimal solution search	<b>Exact</b>
<b>Align according to phylogeny</b>	Align most similar first, then merge together	<b>Progressive</b>
<b>Retain conserved regions</b>	Conserved regions guide alignment	<b>Consistency based</b>
<b>Maximize similarity to model</b>	Create a model, align each sequence to that	<b>Probabilistic</b>
<b>"Learn" about important regions and extend the alignment from secure seeds</b>	Improve alignment from draft alignments	<b>Iterated</b>

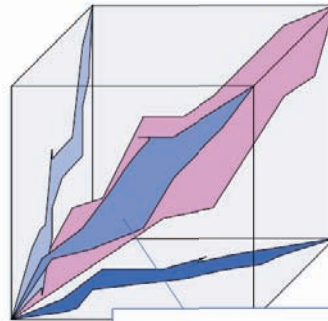
Alignment algorithms that operate according to one or more of these principles are easily accessible online *via* the EBI:

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

Exact MSA (... maximize alignment scores)

**MSA:**

Multidimensional dynamic programming; search space in n-dimensional path matrix can be bounded by optimal pairwise alignments. Optimizes sum-of-pairs score which may not be the most biologically meaningful score anyway. More accurate than e.g. progressive methods but compute intensive. Practical limit ~ 10 seq. of  $L = 200$ .



Pairwise alignments bound the search volume in n-space!

**Novel developments:**

**DCA** (divide and conquer: split into subsequences, then recombine),  
**OMA** (iterated improvement of splits).

Exact methods certainly have their place where it comes to analyzing and improving algorithms; they are especially of interest to computer science because high-dimensional optimal alignment is a difficult problem. However they cannot compete in terms of result-quality with modern heuristic methods. This is not only because they really don't scale to current genome-scale questions or even modest sized protein families, but also because optimizing the score derived from a pair-score mutation data matrix plus an empirical affine gap model is not a really a very good objective for MSAs that inform about biology in the first place.

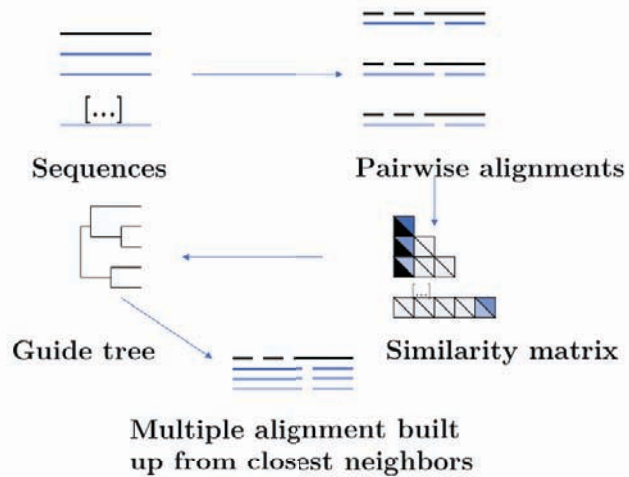
(... align the most similar sequences first)

**Clustal W, X:**

Considered by many to be *The Standard* method.

1. Build pairwise similarity matrix
2. **Build guide tree.**
3. join neighbors into profiles according to guide tree.
4. Align according to tree.

Key limitation: persistence of early errors. Best performance is for globally alignable, gap-poor sequence sets. Performance progressively worse for multidomain proteins and distant similarities.



Bottom line: Don't use CLUSTAL W. There are much better and equally convenient algorithms available.

**Progressive** alignment is one of the fundamental algorithmic approaches to MSA. Pure progressive alignment algorithms are only of historical interest today, since they suffer from unacceptable degradation of accuracy for sequences below ~30% ID due to the fact that early alignment errors cannot be corrected.

## Consistency based multiple alignment: alignment based on motifs and patterns

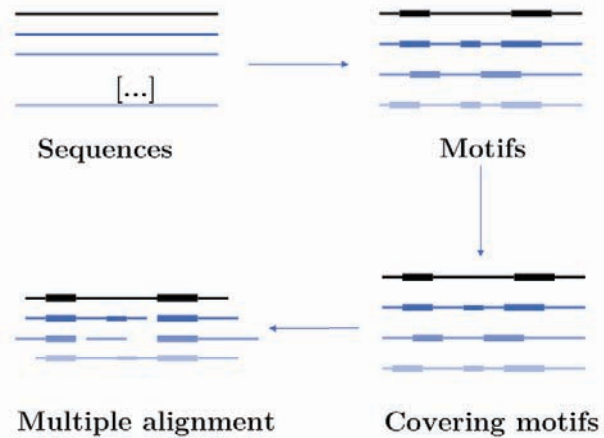
( ... emphasize  
conserved  
-i.e. important-  
regions)

### MUSCA:

- (I) TEIRESIAS motif discovery.
- (II) Use set covering algorithm to select motifs that are common to sequence set.

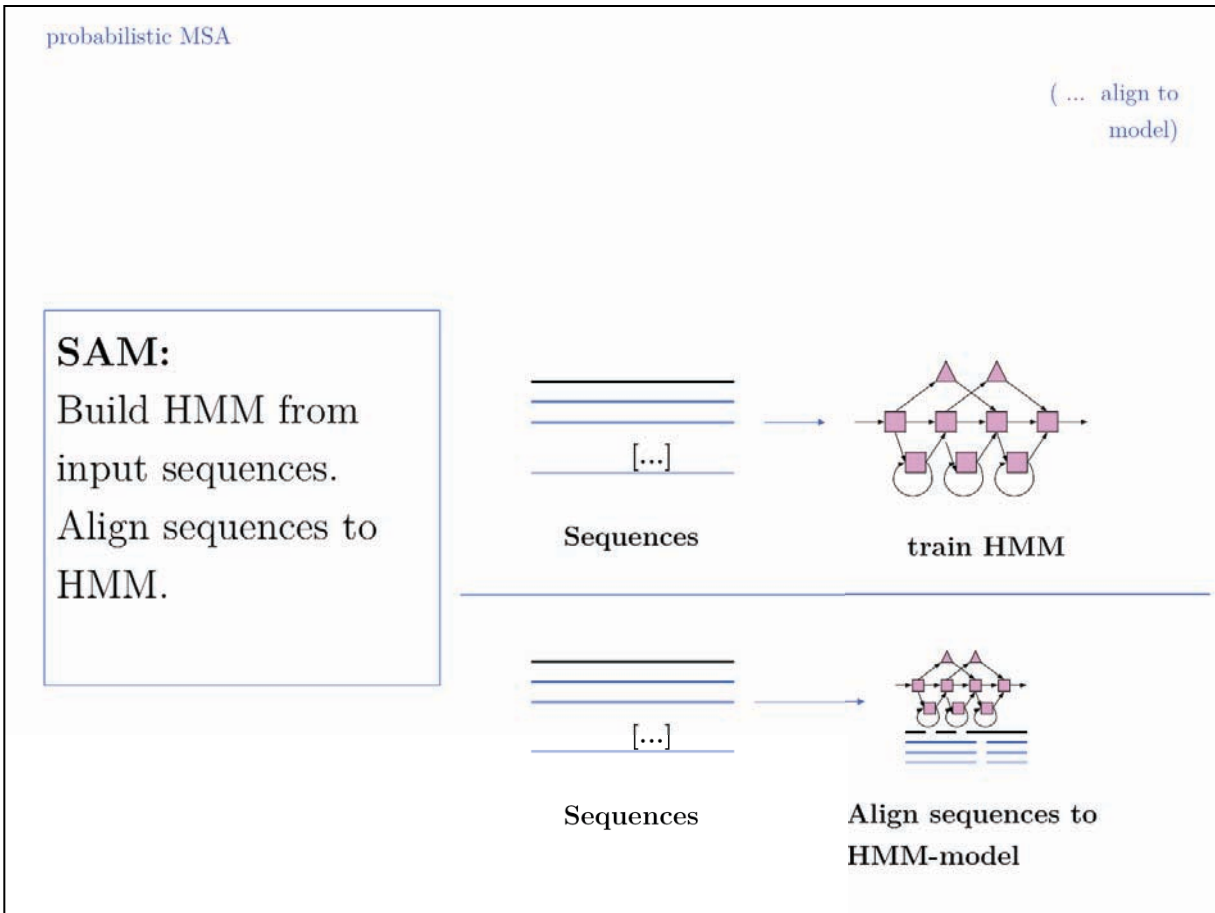
### MEME:

Postulate motif, compare residue composition of motif with background, choose motifs to maximize composition difference, output.



**Consistency** based multiple alignment is one of the fundamental algorithmic approaches to MSA. Many modern algorithms have a consistency based step included, however none of them relies solely on consistency, since problems from spurious local similarity can corrupt the alignment.





**Probabilistic** multiple alignment is one of three fundamental algorithmic approaches to MSA.

A statistical model of the sequences is built, then the alignment can be generated by aligning the sequences to the model. Of course, aligning sequences to a profile is a special case of this procedure: PSI BLAST can thus be used as an alignment algorithm. The most widely used algorithm is Sean Eddy's **HMMER**<sup>1</sup>, a profile hidden Markov model tool, which is also used in the generation of the **Pfam** domain database<sup>2</sup>.

<sup>1</sup> <http://hmmer.janelia.org/>

<sup>2</sup> <http://pfam.sanger.ac.uk/Pfam>



profile based

The MSA derived from aligning sequences to a profile in a **PSI-BLAST** search is also a model based alignment.

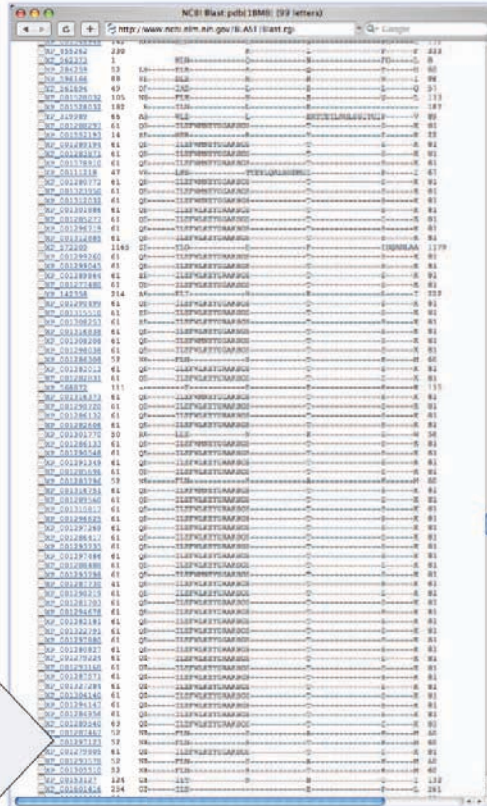
## PSI-BLAST

1. Begin with BLAST search
2. Identify significant hits
3. Align to query
4. Compile into position specific scoring matrix (PSSM or 'sequence profile')
5. **Repeat search with profile**
6. Add new aligned hits to PSSM
7. Iterate until no new sequences can be added

Results can be displayed as an MSA.

Choose formatting option:

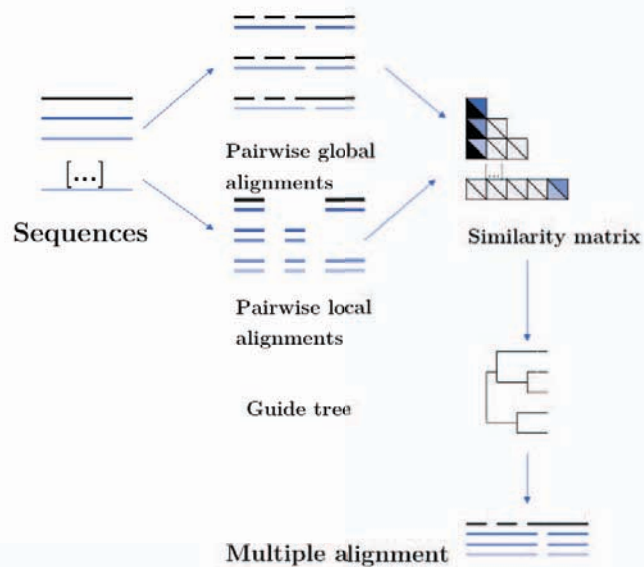
"Flat query-anchored with letters for identities"



Altschul *et al.* (1998) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Research* **25**:3389-3402 (<http://www.ncbi.nlm.nih.gov/pubmed/9254694>)

**TCoffee - a hybrid algorithm significantly improves performance:**

1. Compute **global** pairwise similarity matrix.
2. Compute top 10 non-intersecting **local** alignments.
3. Combine by looking at triplets of sequences.
4. Build guide tree.
5. Align according to tree.



Very good results.

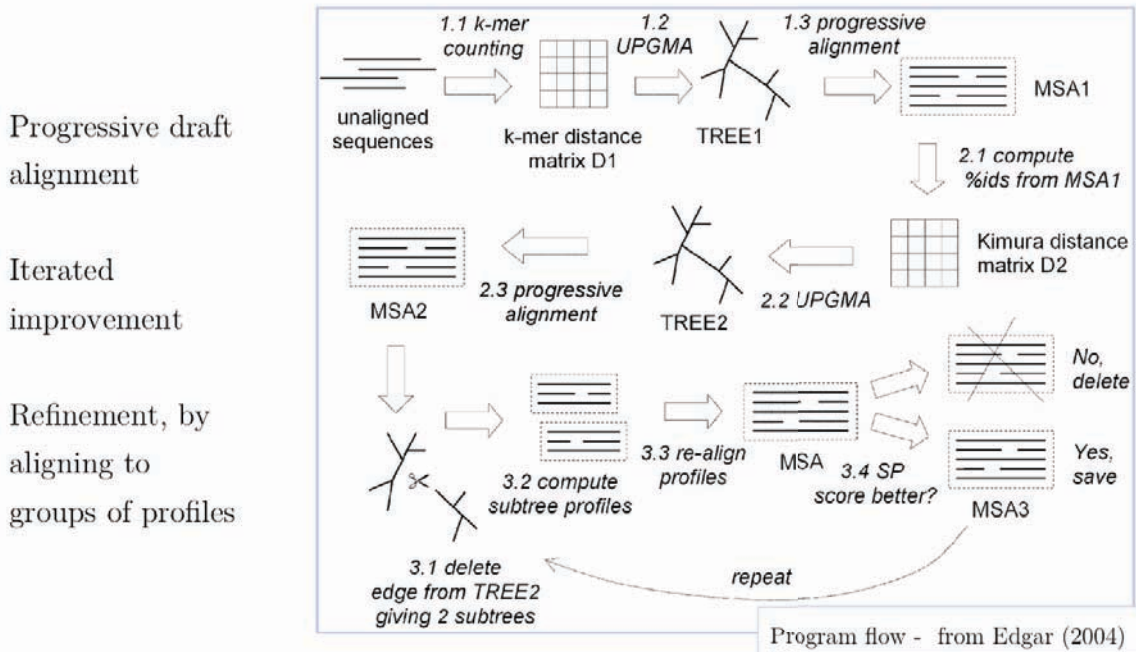


I personally rate Toffee as one of the most useful and useable tools that are currently available. It is robust, fast, and gives reasonable results for many cases. Usually it is **very** noticeably better than CLUSTAL and I would reject any result based on CLUSTAL for that reason.

Run Toffee via the EBI **TCoffee** server which is very easy to use (although alignment size is limited; Source code can be obtained and a local installation on UNIX machines is straightforward. The Toffee Web page<sup>1</sup> links to another Web server and also offers 3DCoffee, a variant that automatically fetches related structures and incorporates structural alignments for increased accuracy.

The inset image shows one of the useful features of Toffee: an alignment output in which sequence is coloured according to the local quality of the alignment. This makes reliable and unreliable regions easy to spot, and immediately highlights outliers that could for example be due to sequence errors, such as frameshifts in exons. (MSA taken from the Mbp1 full-length protein alignment).

## MUSCLE



Better than CLUSTAL, faster than ProbCons - one of the general purpose algorithms of first choice with the capacity to align thousands of sequences in one run.

Run MUSCLE MSAs via the EBI MSA server, which is very easy to use, or via the Berkeley **MUSCLE server**<sup>1</sup>, courtesy of Kimmen Sjolander's lab. Source code and compiled code can be obtained from the MUSCLE homepage<sup>2</sup> and a local installation on UNIX and Windows machines is straightforward. That site also hosts the PREFAB multiple alignment benchmark.

MUSCLE is one of the algorithms provided in the R package **msa**.

<sup>1</sup> [http://phylogenomics.berkeley.edu/cgi-bin/muscle/input\\_muscle.py](http://phylogenomics.berkeley.edu/cgi-bin/muscle/input_muscle.py)

<sup>2</sup> <http://www.drive5.com/muscle/> Muscle

## Practical considerations for MSAs:

Principle: try to use sequences that are well distributed on the evolutionary tree. If a group of sequences biases the alignment, then not all information contributes equally to the result.

Include sequences with known structure wherever possible.

You may include more sequences than the ones that you are actually interested in comparing.

Ensure that your selection of sequences is appropriate to your question, e.g. clearly distinguish between orthologues and paralogues when trying to study function.

Spend some time and thought **before** you run the MSA to review the sequences that you are planning to align. Including un-alignable sequence **will** lead the algorithms astray and has the potential to degrade the entire alignment.

The requirement not to align **non-homologous** sequence should really be extended not to align (or at least: not to evaluate) sequence segments that have evolved in different context, such as in different local structural environments after insertions or deletions have occurred. The reason is: if the structural environment is not conserved, the mutation data matrix scores are irrelevant for the residues that are paired up. They may be "aligned" by the algorithm, but they are really not equivalent in structure or function, thus whether they have a good or poor similarity score is meaningless.

Rule #1 of pairwise alignment applies:

## **Don't align non-homologous sequence!**

### **Here are some heuristics for preparing the sequences for MSAs:**

- Remove unalignable (non-homologous) domains and N- and C-terminal extensions.
- Consider what the alignment is supposed to show! Include no more sequence than that which is needed.
- Be cautious when including widely divergent outlier sequences - better to align those to a profile derived from the result, than to create errors in the sensitive early stages of progressive alignments.
- If individual sequences do not align well, reconsider the evidence for their homology: the alignment may try to tell you that they are not homologous after all and should be removed from the set.
- If sequence segments behave as outliers, they may indicate sequencing frameshifts, skipped exons or erroneously translated introns: review the evidence for the gene-model that was used to define the sequence.
- If domain annotations exist: align domains separately.



## Data formats for MSA: multi FASTA .mfa

```

>0918_CANAL_XP_710918_256..352
-VIWDYETGWVHLTGIWKASLTIDGSNVSPSHLKADIVKLESTPKE----Y--QQYIKR
IRGGFLKIQ---GTWLPYKLCILARRFCYLRYSLIP-IFGTDFFPDS
>9773_DEBHA_XP_459773_187..274
-IIWDYETGFVHLTGIWKASIND--EVNTHRNKADIVKLESTPKQ----Y--HQHIKR
IRGGFLKIQ---GTWLPFDLCKMLAKRFCYHIRFQLIP-IF-----
>MBP1_SACCE_NP_010227_024..107
SIMKRKDDWVNATHILKAANF-----AKAKRTRILEKEV-----L--KETHEK
VQGGFGKYQ---GTWVPLNIAKQLAEKFSVY--DQLKP-LFDFTQTDG
>2599_ASPTE_XP_001212599_130..218
-IMWDYNIGLVRTTPLFRS-----QNSKTTPAKVLDPANPGL--REISHS
ITGGAIVAQDKPGYWIPFEAAKAVAATFCWRIRYALTP-IFGLDFPSQ
>3510_ASPFU_XP_753510_089..163
-LMRRSKDGYVSATGMFKIAFPW--AKLEEEKAEREYLKTRGTSEDEIAG-----
-----NIWVSPLLALELAKEY-----QMYDWRALLD---
>7766_ASPNI_XP_657766_089..163
-LMRRSKDGYVSATGMFKIAFPW--AKLEEERSEREYLKTRPETSEDEIAG-----
-----NVWISPVLALELAEEY-----KMYDWRALLD---
>2267_NEUCR_XP_962267_085..162
-LMRRSQDGYISATGMFKATFPY--ASQEEEEAEKRYIKSIPTTSSEETAG-----
-----NVWIPPEQALILAEY-----QITPWIRALLDPSD
>3762_MAGGR_XP_363762_084..161
-LMRRSSDGYVSATGMFKATFPY--ADADEEAERNYIKSLPATSKEETAG-----
-----NVWISPDQALALAEY-----SIATWIRALLDPTD
>5459_GIBZE_XP_385459_077..154
-LMRRSYDGFVSATGMFKASFPY--AEASDEDAERKYIKSLPTTSHEETAG-----
-----NVWIPPEQALILAEY-----KISPWIRALLDPTP
>3412_CANAL_XP_723412_087..178
-VLRRVQDSFVNVTQLFQILIKL--EVLPTSQVDNYFDNEILSNLKYFGSSSNTPYQLDL
RKHQNIYLO---GIWIPYDKAVNLALKFD-----IYEITKKLF----
>9901_DEBHA_XP_459901_067..158
-ILRRVQDSYINISQLFSILLKI--GHLSEAQLTNFLNNEILTNTQYLSSGGSNPQFNDL
RNHEVRDLR---GLWIPYDRAVSLALKFD-----IYELAKSLF----

```

Three common formats exist for MSA results. An **aligned** multi FASTA file contains FASTA formatted sequences into which gap characters have been inserted. Of course, multi FASTA files can also be unaligned and they are the most common way of formatting **input files** for MSAs.

## CLUSTAL

# Data formats for MSA: CLUSTAL .aln

```
CLUSTAL FORMAT for T-COFFEE Version_5.05 SCORE=36, Nseq=11, Len=108

0918 CANAL      -VIWDYETGWVHLTGIWKASLTIDGSNVSPSHLKADIVKLESTPK-----Y--QQYIKR
9773 DEBHA      -IIWDYETGFVHLTGIWKASIND--EVNTHRNKADIVKLESTPKQ-----Y--HQHIKR
MBP1 SACCE      SIMKRKKDDWVNATHILKAANF-----AKAKRTRILEKEV-----L--KETHEK
2599 ASPTE      -IMWDYNIGLVRTTPLFRS-----QNYSKTTPAKVLDANPGL--REISHS
3510 ASPFU      -LMRRSKDGYVSATGMFKIAFPW--AKLEEEKAEREYLKTRREGTSEDEIAG-----
7766 ASPNI      -LMRRSKDGYVSATGMFKIAFPW--AKLEEEEREREYLKTRPETSEDEIAG-----
2267 NEUCR      -LMRRSQDGYISATGMFKATFPY--ASQEEEAERKYIKSIPPTSSEETAG-----
3762 MAGGR      -LMRRSSDGYVSATGMFKATFPY--ADADEEAERNYIKSLPATSKEETAG-----
5459 GIBZE      -LMRRSYDGFVSATGMFKASFPY--AEASDEDAERKYIKSLPTTSHEETAG-----
3412 CANAL      -VLRVQDSFVNVTQLFQILIKL--EVLPTSQVDNYFDNEILSNLKYFGSSSNTPOYLDL
9901 DEBHA      -ILRRVQDSYINISQLFSILLKI--GHLSEAQLTNFLNNEILTNTQYLSSGGSNPQFNDL
      ::      . : : :

0918 CANAL      IRGGFLKIQ---GTWLPYKLCILARRFCYLYRYSLIP-IFGTDFFDS
9773 DEBHA      IRGGFLKIQ---GTWLPPDLCKMLAKRFCYHIRFQLIP-IF-----
MBP1 SACCE      VQGGFGKYQ---GTWVPLNIAKQLAEKFSVY--DQLKP-LFDFDTQTDG
2599 ASPTE      ITGGAIVAQDKPGYWIPPEAAKAVAATFCWRIRYALTP-IFGLDFPSQ
3510 ASPFU      -----NIWVSPLLALELAKEY-----QMYDWVRALLD---
7766 ASPNI      -----NVWISPVLALELAEEY-----KMYDWVRALLD---
2267 NEUCR      -----NVWIPPEQALILAEY-----QITPWIRALLDPSD
3762 MAGGR      -----NVWISPDQALALAEY-----SIATWIRALLDPTD
5459 GIBZE      -----NVWIPPEQALILAEY-----KISPWIRALLDPTP
3412 CANAL      RKHQNIYLO---GIWIPYDKAVNLALKFD-----IYEITKKLF---
9901 DEBHA      RNHEVRDLR---GLWIPYDRAVSLALKFD-----IYELAKSLF---
      . * : . . : * : :
```

Three common formats exist for MSA results.

The CLUSTAL format is not the same as the CLUSTAL algorithm. A CLUSTAL formatted alignment is the format in most common use.

Take care when formatting input FASTA files to ensure the **first 10 characters in your input file are unique** and contain **no special characters!** These are the characters that are usually used for the sequence names of the .aln files. I have seen programs break if they contain blanks, hyphens and | (the pipe character). The latter is especially annoying, since the | character is used in NCBI FASTA files to separate the database identifier from the accession number.



MSF

## Data formats for MSA: MSF .msa

```
MSF: 108 Type: P Check: 3302 ..
Name: 0918 CANAL oo Len: 108 Check: 8295 Weight: 1.000
Name: 9773 DEBHA oo Len: 108 Check: 3488 Weight: 1.000
Name: MBP1 SACCE oo Len: 108 Check: 808 Weight: 1.000
Name: 2599 ASPTE oo Len: 108 Check: 241 Weight: 1.000
Name: 3510 ASPFU oo Len: 108 Check: 9082 Weight: 1.000
Name: 7766 ASPNI oo Len: 108 Check: 9952 Weight: 1.000
Name: 2267 NEUCR oo Len: 108 Check: 5383 Weight: 1.000
Name: 3762 MAGGR oo Len: 108 Check: 3063 Weight: 1.000
Name: 5459 GIBZE oo Len: 108 Check: 4901 Weight: 1.000
Name: 3412 CANAL oo Len: 108 Check: 5134 Weight: 1.000
Name: 9901 DEBHA oo Len: 108 Check: 2955 Weight: 1.000

//

0918 CANAL .VIWDVETGW VHLTGIWKAS LTIDGSNVSP SHLKADIVKL LESTPKE... .Y..OOYIKR
9773 DEBHA .ITWDVETGF VHLTGIWKAS IND..EVMTH RNLKADIVRL LESTPKQ... .Y..HOHIKR
MBP1 SACCE SIMKRKDDW VNATHILKAA NF..... .AKAKRTRI LEKEV... .L..KETHEK
2599 ASPTE .IMWDYNIGL VRTTPLFRS. .... .ONYSKT TPAKVLANP GL..REISHS
3510 ASPFU .LMRRSKDGY VSATGMFKIA FPW..AKLEE EKAEREYLKT REGTSEDEIA G.....
7766 ASPNI .LMRRSKDGY VSATGMFKIA FPW..AKLEE ERSEREYLKT RPTSEDEIA G.....
2267 NEUCR .LMRRSODGY ISATGMFKAT FPY..ASOEE EEAERKYIKS JPTSSSEETA G.....
3762 MAGGR .LMRRSDGY VSATGMFKAT FPY..ADAD EEAERNYIKS LPATSSSEETA G.....
5459 GIBZE .LMRRSYDGF VSATGMFKAS FPY..AEASD EDAERKYIKS LPTTSSEETA G.....
3412 CANAL .VLRVQDSF VNVTLFOIL IKL..EVLPT SQVDNYFDNE ILSNLKYFGS SSNTPOYLDL
9901 DEBHA .ILRRVQDSY INISQLFSIL LKI..GHLSE AQLTNFLNNE ILTNTQYLSS GGSNPQFDL

0918 CANAL IRGGFLKIO. .GTWLPYKL CKILARRFCY YLRYSLIP. I FGDFPDS
9773 DEBHA IRGGFLRIO. .GTWLEFDL CKMLAKRFCY HIRFOLIP. I F.....
MBP1 SACCE VQGGFGKYQ. .GTWVPLNI AKQLAEEKFSV Y..DOLKP.L FDFQTQDQ
2599 ASPTE ITGGAIVAQD KPGYWIPFEA AKAVAATFCW RIRYALTP. I FGLDFPSQ
3510 ASPFU ..... .NIWVSPLL ALELAKEY. . . . .OMYDWW RALLD...
7766 ASPNI ..... .NVWISPVL ALELAEEY. . . . .KMYDWW RALLD...
2267 NEUCR ..... .NVWIPPEO ALILAEEY. . . . .OITPWI RALLDPSD
3762 MAGGR ..... .NVWISPDO ALILAEEY. . . . .SIATWI RALLDPTD
5459 GIBZE ..... .NVWIPPEO ALILAEEY. . . . .KISPWI RALLDPTP
3412 CANAL RKHQNIYLO. .GIWIPYDK AVNLALKFD. . . . .IYEIT KKLK...
9901 DEBHA RNHEVRDLR. .GLWIPYDR AVSLALKFD. . . . .IYELA KSLF....
```

Three common formats exist for MSA results.

MSF is a legacy format from the GCG package of sequence alignments, also produced by the EMBOSS tool EMMA, and supported as a valid input format for many programs. Gaps are denoted by periods and checksums are calculated for the sequences and for the alignment.

## manual editing

Manual editing can often improve MSAs because you can take context into account whereas algorithms typically calculate scores based on objective functions that are computed over columns only.

( Context: e.g. where are indels placed, which residues are part of a functional site ... )

- Move sequences of different length (thus having obviously different structure) for dissimilar families into separate columns. Try not to include non-alignable residues in the same column and create the impression they are alignable.
- Move indels into regions adjacent to secondary structure elements, even if this creates a lower sequence alignment score. This is where they are commonly **accommodated** in structure, even if they have been generated elsewhere.
- Move two/four/six-residue insertions equally to both sides of a conserved beta-turn, rather than postulating only a single insertion on one side.
- Adjust gap positions to minimize the number of indel **sites**.
- Consider the evolutionary tree to minimize number of indel **events**.
- Conserve hydrophobic patterns in regions of secondary structure, e.g. alternating hydrophobic residues in beta-strands and  $i \rightarrow i+4$  patterns in alpha-helices.
- Preserve conservation especially in binding sites and functional motifs, even at the cost of larger indels.

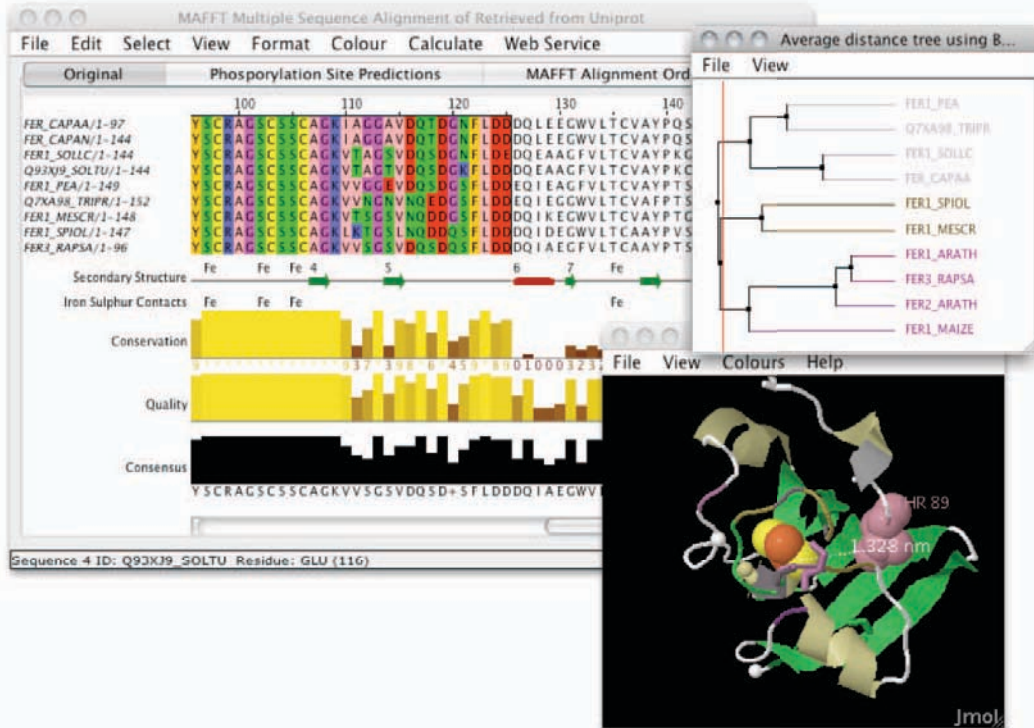
It is common and perfectly permissible to manually edit a MSA with some biologically motivated heuristic in mind **as long as you document what you have done!** In the early days of MSAs, editing was simply required since the results were often obviously inadequate. In all cases in which the algorithm uses only the input sequences for the alignment, this still holds true. However, regarding the more modern template-based procedures (e.g. SPEM, PROMALS or PRALINE) I would be more reluctant to edit, since we may be actively ignoring/discarding the additional information the algorithm has used.

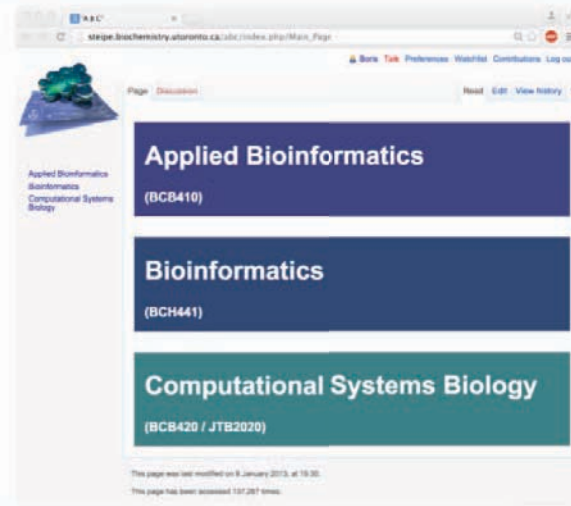
You **can** use a text-editor to edit MSA's, (e.g. MS Word) ...

1. Format your sequences in a **fixed-width** font (e.g. **Courier New**).
2. Display the entire sequence in one line: small font, landscape -page setup, legal- or custom wide paper size, small margins, magnified view.
3. Press the **<ALT>** key
4. You can now select **columns** and **blocks** of text.
5. You can **format** your selection; you can **delete, copy, cut...**
6. You can also **paste** entire columns or column ranges into the alignment.

... but you should avoid it, because **much better tools are freely available.**

JALVIEW





<http://steipe.biochemistry.utoronto.ca/abc>

B O R I S . S T E I P E @ U T O R O N T O . C A

DEPARTMENT OF BIOCHEMISTRY & DEPARTMENT OF MOLECULAR GENETICS  
UNIVERSITY OF TORONTO, CANADA