

# DIALIGN

A segment-segment multiple sequences alignment program



Kai Ye & Qilan Li

## Outline of method section

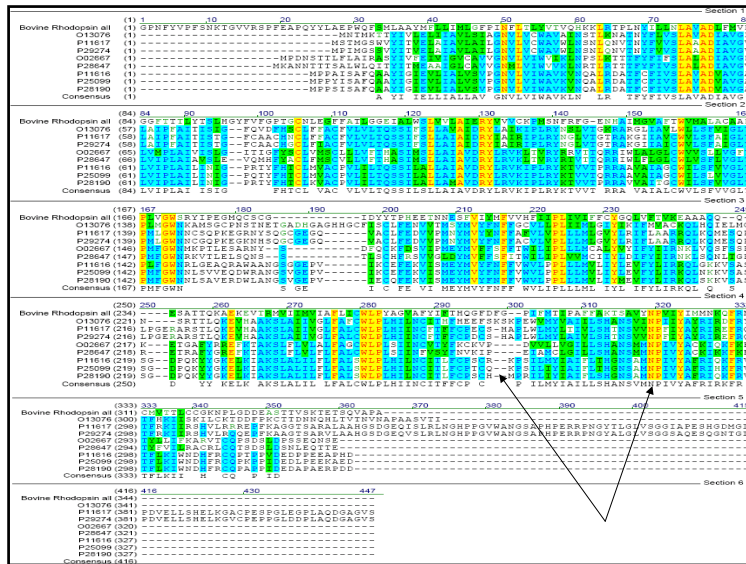
- 1. What is protein sequence?
- 2. What is sequence alignment?
- 3. Why we need sequence alignment?
- 4. Methods to do sequence alignment.
- 5. Advantage and disadvantage of traditional method: ClustalW
- 6. DIALIGN
- 7. Advantage and disadvantage of DIALIGN
- 8. Improvement of speed

## What is protein sequences?

- A protein sequence is represented by a string a of letters coding for the 20 different types of amino acid residues.
- >P04201
- MDGSNVTSFVVEEPTNISTGRNASVGNHRQIPVHW  
VIMSISPVGFVENGILLWFLCFMRMRNPFTVYITHLSIA  
DISLLFCIFILSIDYALDYELSSGHYYTIVTL SVTFLFGYN  
TGLYLLTAISVERCLSVLYPIWYRCHRPKYQSALVCAL  
LWALSCLVTTMEYVMCIDREEESHNRNDCRAVIFAIL  
SFLVFTPLMLVSSSTILVVKIRKNTWASHSSKLYIVIMVTII  
IFLIFAMPMRLLYLLYYEYWSTFGNLHHISLLFSTINSSA  
NPFYFFVVGSSKKKRFKESLKVVLTRAFKDEMQPRRQ  
KDNCNTVTVETVV

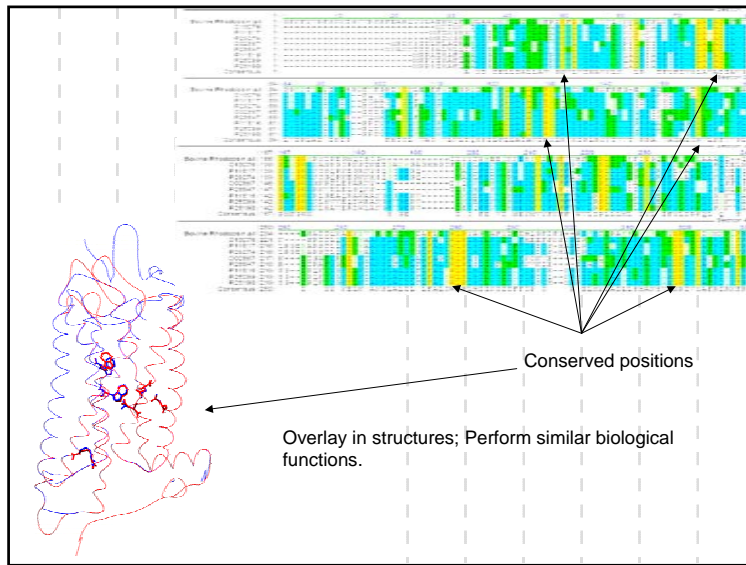
## What is sequence alignment?

- A protein sequence alignment is created when the residues in one sequence are lined up with those in at least one other sequence.
- Optimal alignment of the two sequences will usually require the insertion of gaps in one or both sequences in order to find the best alignment.



## Why we need sequence alignment?

- Alignment of two residues implies that those residues are performing similar roles in the two different proteins.
- This allows for information known about specific residues in one sequence to be potentially transferred to the residues aligned in the other.



## Methods to make sequence alignment

- Global progressive alignment algorithm:
  - MULTALIGN; MULTAL; PILEUP; CLUSTALW
  - differ mainly in the method used to determine the order of alignment of the sequences.
- Local dynamic programming algorithm
  - PIMA;

Traditional methods

## Methods to make sequence alignment

- Segment-segment comparisons rather than residue-to-residue comparisons: DIALIGN
- Genetic algorithm: SAGA
- Hidden Markov models (HMMs): HMMT  
Speech recognition

"New" methods

## Advantage and disadvantage of traditional method: ClustalW

- ClustalW is superior on globally related protein families ( same ancestor ).
- However, many protein families share only isolated regions of local similarity.
- The resulting alignments depend sensitively on a set of user-defined parameters, especially the gap penalty.

## DIALIGN

- Segment-to-segment alignments with variable segment length
- No gap penalties
- Regions of low similarity are excluded from DIALIGN alignment.

## DIALIGN consistent collections of segment pairs

```

I A V L F A E D
| | | | | |
L A V I F G S /
W D D V T F D A E
    
```

A

```

I A V L F A E D
| | | | | |
L A V I F G S /
W D D V T F D A E
    
```

B

```

I A V L F A E D
| | | | | |
L A V I F G S /
W D D V T F D A E
    
```

C

```

I A - V L F - A E d
L A - V I F - G s -
w d d V T F d A E -
    
```

D

## Advantage and disadvantage of DIALIGN

- Alignments produced by DIALIGN seem to be comparable to the results of standard global alignment methods as CLUSTALW. However, if sequences are only locally related, DIALIGN seems to be clearly superior to other methods.
- Too slow to align sequences in the order of hundreds of kilobases or more.

## Improvement of speed

- Solution 1:  
CHAOS, a novel algorithm for rapid identification of chains of local pair-wise sequence similarities. Alignments calculated by CHAOS are used as anchor points to improve the running time of DIALIGN. 95%

## Improvement of speed

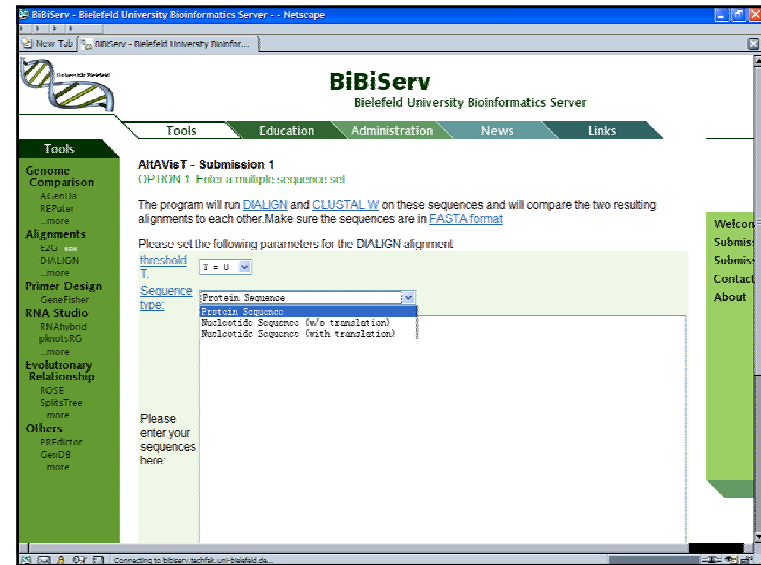
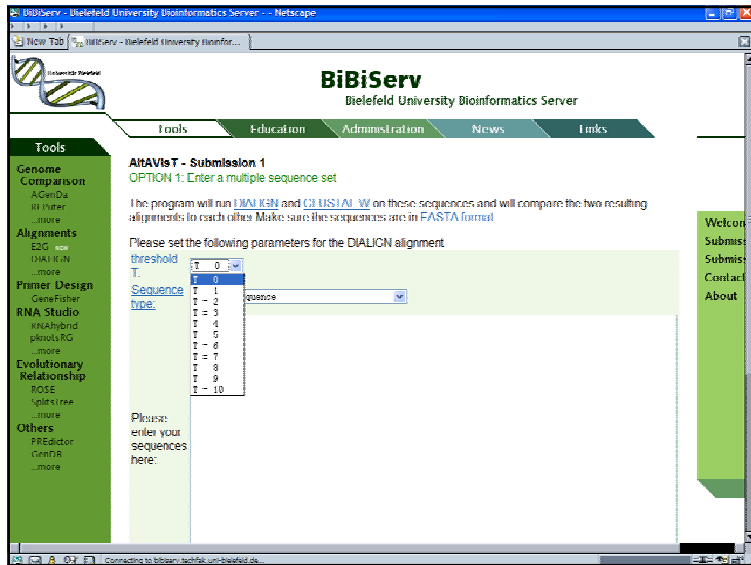
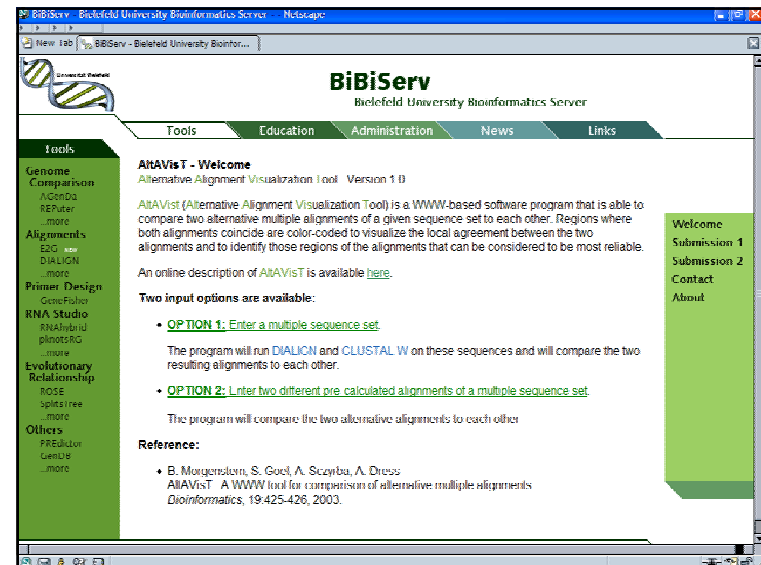
- **Solution 2: DIALIGN P, a parallel version of the multi-alignment program DIALIGN.**
- **A lot of sequences**
  - (a) pair-wise sequence alignments that are used as a first step to multiple alignment account for most of the CPU time in DIALIGN. Since alignments of different sequence pairs are completely independent of each other, they can be distributed to multiple processors without any effect on the resulting output alignments.
- **Longe sequences**
  - (b) For alignments of large genomic sequences, DIALIGN use a heuristics by splitting up sequences into sub-sequences based on a previously introduced *anchored alignment* procedure. For test sequences, this combined approach reduces the program running time of DIALIGN by up to 97%.

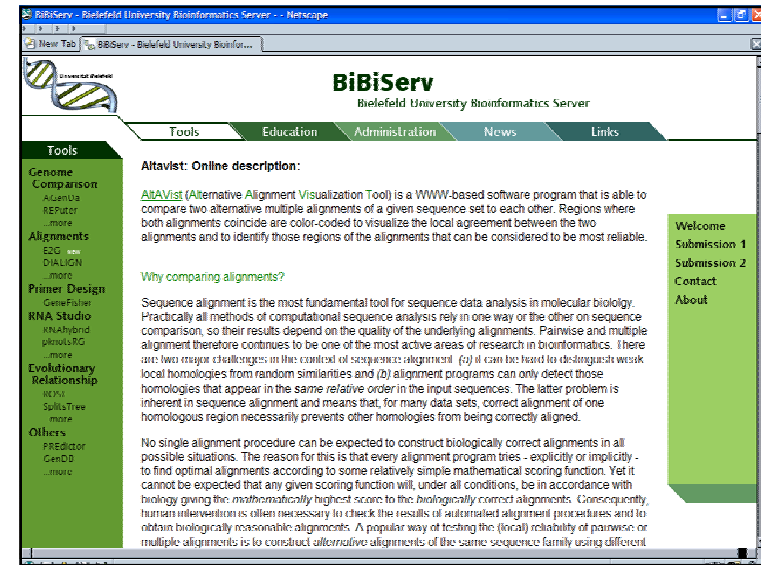
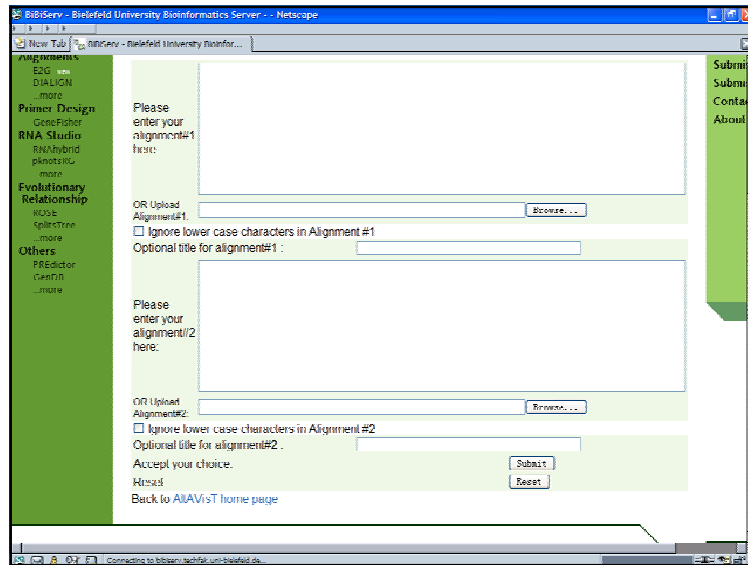
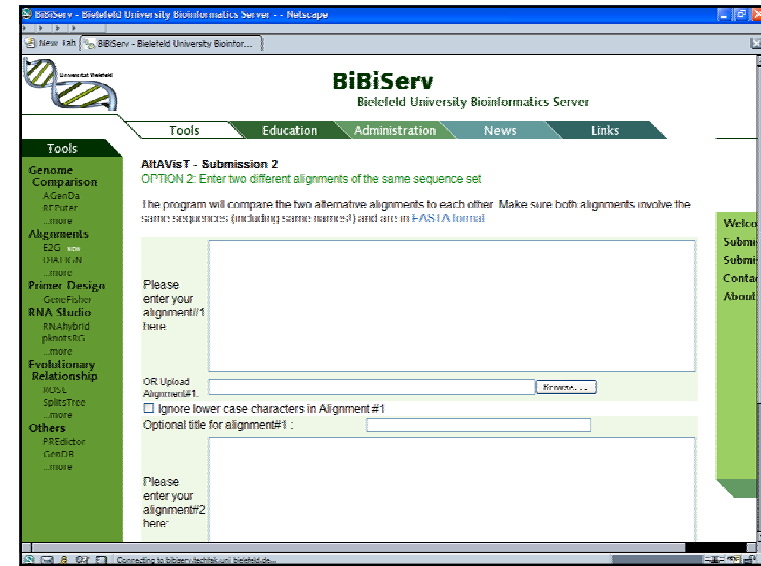
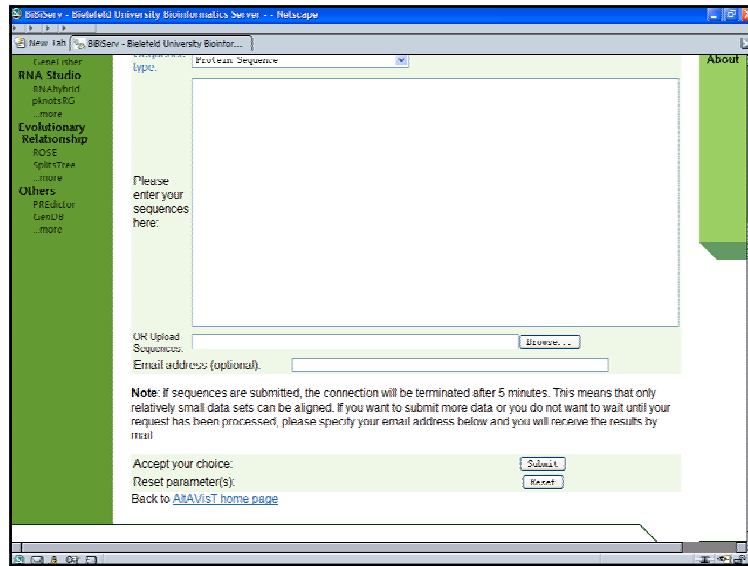
## Web server

- **Altavist**
- **CHAOS + DIALIGN**

# Altavist

- AltAVist (Alternative Alignment Visualization Tool) is a WWW-based software program that is able to compare two alternative multiple alignments of a given sequence set to each other. Regions where both alignments coincide are color-coded to visualize the local agreement between the two alignments and to identify those regions of the alignments that can be considered to be most reliable.





BBIServ - Bielefeld University Bioinformatics Server - Netscape

Alignments  
EZG  
DIALIGN  
more

Primer Design  
Gene Fisher  
RNA Studio  
RNAfold  
phits.RG  
more

Evolutionary Relationship  
ROSE  
SPITS tree  
more

Others  
PREDATOR  
(Lentini)  
more

welcome

Submission 1  
Submission 2  
Contact  
About

### Why comparing alignments?

Sequence alignment is the most fundamental tool for sequence data analysis in molecular biology. Practically all methods of computational sequence analysis rely in one way or the other on sequence comparison, so their results depend on the quality of the underlying alignments. Pairwise and multiple alignment therefore continues to be one of the most active areas of research in bioinformatics. There are two major challenges in the context of sequence alignment, (a) it can be hard to distinguish weak local homologies from random similarities and (b) alignment programs can only detect those homologies that appear in the same relative order in the input sequences. The latter problems are inherent in sequence alignment and means that, for many data sets, correct alignment of one homologous region necessarily prevents other homologies from being correctly aligned.

No single alignment procedure can be expected to construct biologically correct alignments in all possible situations. The reason for this is that every alignment program tries - explicitly or implicitly - to find optimal alignments according to some relatively simple mathematical scoring function. Yet it cannot be expected that any given scoring function will, under all conditions, be in accordance with biology giving the mathematically highest score to the biologically correct alignments. Consequently, human intervention is often necessary to check the results of automated alignment procedures and to obtain biologically reasonable alignments. A popular way of testing the (local) reliability of pairwise or multiple alignments is to construct alternative alignments of the same sequence family using different alignment methods. Notredame *et al.* (2000) used this idea systematically and proposed a software tool that integrates results from different multi alignment methods into one single output alignment.

For multiple alignment, a variety of programs are now available that rely on very different objective functions and optimization techniques. The results of these methods can therefore be quite diverse, see Notredame (2002) for an excellent review of the state-of-the-art multi-alignment algorithms and Thompson *et al.* (1999b) for a systematic evaluation of the most widely used software tools. If two alignments have been constructed by different methods, those regions where both alignments coincide are generally considered to be more reliable than regions where they disagree. However, manually comparing different multiple alignments is a tedious task.

BBIServ - Bielefeld University Bioinformatics Server - Netscape

### Input options:

AtAVis / compares two different multiple alignment of a given data set and highlights regions where both alignments coincide. Two input options are available:

- It is possible to enter a **family of sequences**. In this case, our program will run the programs DIALIGN (Morgenstern, 1998) and CLUSTAL W (Thompson *et al.*, 1994) on the input sequences and compare the resulting alignments to each other. These two programs are currently among the most popular multi-alignment methods. Since they rely on fundamentally different algorithmical approaches, those parts of the alignments where both programs agree can be considered to be reliable.
- It is possible to enter two different **pre-calculated alignments** of a sequence family that may have been produced by any method, this way the user can compare the results of arbitrary alignment methods to each other.

With either option, those residue pairs that are aligned to each other in *both* alignments are colored different colors are used to distinguish groups of residues for which the alignment coincides *within* groups but not *between* different groups. In other words, considering alignments as *consistent* *computational relationships* (introduced in Morgenstern *et al.* (1998)), residues pairs that are in the same column and have the same color belong to the set-theoretical intersection of the equivalence relations corresponding to the two alignments.

Our tool can not only be used to determine reliable regions in alignments but also to evaluate alignment programs by comparing the alignments they produce to reference alignments that are considered as a *standard of truth*. There is now a high-quality data base called BAI/BASE that has been designed as a benchmark data base for evaluation of multiple alignment methods (Thompson *et al.*, 1999a). The authors of BAI/BASE also provide software that automatically compares arbitrary alignments of their test data to the reference alignments and determines the overall degree of agreement between these two alignments. However, for the development of alignment methods, it can be interesting to know not only the overall quality of the produced alignments but to also know where exactly these alignments are in agreement with the given reference alignment and where they are not. Our method can be used for this purpose and should therefore also be useful for further development and improvement of pairwise and multiple alignment methods.

BBIServ - Bielefeld University Bioinformatics Server - Netscape

### Program output:

Below is the result of AtAVisT applied to a small test sequence set. The first alignment has been produced by DIALIGN, the second one by CLUSTAL. For each column in the first alignment, those residue pairs are colored that also appear in one column in the second alignment. Different colors are used to distinguish groups of residues where the alignment coincides *within* groups but not *between* different groups. For example, the two Ms in column 4 in the DIALIGN alignment also appear in the same column in the CLUSTAL alignment, namely in column 21; they are therefore colored the same both: blue for the two Cs in the same column of the DIALIGN alignment. These residues also appear in a common column in the CLUSTAL alignment, namely in column 4. However, the Ms and Cs belong to different columns in the CLUSTAL alignment so different colors are used. All lower case residues in the DIALIGN alignment are printed in black because they are not considered aligned by DIALIGN, regardless in which column they are.

In the second alignment, all residues have the same color as in the first alignment so the two alignments can be easily compared. This may imply, however, that residues in the second alignment appear in the same color, even though they are not aligned together in the first alignment, see for example column 21 in the second alignment.

This is what the AtAVisT output looks like:

**THE RESULT OF DIALIGN IS :**

```

dntp_mouse YQSMIS-----QTELLSQQYQLLINSQVMDCHFMGSEFVDSL
yue6_zeel1 ---MIS-----RYLNAVYDNDLQDLVNSGVDLCHALMSGQFTDFL
cbp1_zeel1 ---LNKRLKAFUNWKPFTITAVLDLNDLPLVYASQKDIICNLSGNSKIVLFL
ydy9_zeel1 ---LNDVFTGFLTFGSGKPFQYVLELNHNPFLVYASQKDIICNLSGNSKIVLFL
cbp1_ploa  YLCHNINRNLFAQDNDKPIVHIVDGLNVDLPLVYASQKDIICNLSGNSKIVLFL
cbp1_zeel1 FVSCSTSVYQMLV--DMSRNLVGIITLLEDGSLVYASQKDIICNLSGNSKIVLFL
dntp_mouse ---QMEVQR--DNLVYDQDQVAGVYEC--SHIFLITNGAC--PY-

```

BBIServ - Bielefeld University Bioinformatics Server - Netscape

This is what the AtAVisT output looks like:

**THE RESULT OF DIALIGN IS :**

```

dntp_mouse YQSMIS-----QTELLSQQYQLLINSQVMDCHFMGSEFVDSL
yue6_zeel1 ---MIS-----RYLNAVYDNDLQDLVNSGVDLCHALMSGQFTDFL
cbp1_zeel1 ---LNKRLKAFUNWKPFTITAVLDLNDLPLVYASQKDIICNLSGNSKIVLFL
ydy9_zeel1 ---LNDVFTGFLTFGSGKPFQYVLELNHNPFLVYASQKDIICNLSGNSKIVLFL
cbp1_ploa  YLCHNINRNLFAQDNDKPIVHIVDGLNVDLPLVYASQKDIICNLSGNSKIVLFL
cbp1_zeel1 FVSCSTSVYQMLV--DMSRNLVGIITLLEDGSLVYASQKDIICNLSGNSKIVLFL
dntp_mouse ---QMEVQR--DNLVYDQDQVAGVYEC--SHIFLITNGAC--PY-
yue6_zeel1 LLI-----SIFSTFTVY-----QDGGVYQKLSQVTFATVYASGNSKIVLFL
cbp1_zeel1 KKYKPKATGHWNRYASTY---DQVAVYVST---KFFVAVYVNRHIS---
ydy9_zeel1 NITFRYQKMLRPNVSKET---SELSQVWY---GFFFLIAYDGHQVY---
cbp1_ploa  VYDAGFFKAEVQVQWLV-----SFAASFFNY---SFFYLVYDAGNSKIVLFL
cbp1_zeel1 KSDKINPAAKKVPLVLD--DQKADLQKTY--KQSPKPKVYDAGNSKIVLFL

```

NOTE: lower case letters are not considered to be aligned

```

THE RESULT OF CLUSTALW IS :

prtp mouse      YQSTNSUYLGLLSSTKRYOILLKRYGVYLSAACHYHCHGWFYUSLN
vua6_cesrel     -----MTSRVLDVYDNDLFDGLYDSDVLLCHALNGRPTDFLG
cnpq_great      ---LKNVLPANWPKPPTATVTEELNQRDLPFYVYAHKQFVTEWAKRATVTFP
yby9_yeast      ---LQVDTGFLFTGDSGRFPQQTDELINMNTFVLYYAGQKQICWNGHNSHELE
cbpy_pi3pce     YESCHFEIWRNLFAGDNRGPFYRHYVSLNMRGSPVLYVAGQDFICWLNRAWTFDPE
obpx_arach      FVFSQSDVYQWHLVD  WRNRLVGLTTELDUCLSLVYAGYDULLCHNGCNSKRVNME

prtp_mouse     QDME-----YQRRPH-LVDYSEDEQVLAGFVVECDIITFLINQAS-PY---
vua6_cesrel    LILS  KSGEHP  IVK  ULHIDQVYVHK  CSQVREAVVGHCHAY
cbpy_pi3pce    NRYDEEFAEQVWVNTASITDEVAGEVKSYP-----RFTYLVFVQAGMC---
yby9_yeast     MNTSRVQSGMGRPVVYVETAPFLAQVQVYV-----PFTTITVYALNIVY---
obpy_pi3pce    NYDADCEKAEVQDM  LVMCRAGC/  oole**0000FF**YDAGMAYT
cbpy_pi3pce    NYDADFEKAEVQDM--LVMCRAGEFFQYV-----RFTYLVYDAGMAYF-
mpx_arach      WSDKYNPDAKRVYV  LVMCRAGLCTVF  QLSKRVVYDAGRVYVMS
  
```

# CHAOS + DIALIGN [job submission]

## Pair-wise and multiple alignment of genomic sequences using CHAOS and DIALIGN

By [Mike Brudno](#) and [Burkhard Morgenstern](#)

This web server calculates pair-wise and multiple alignments of genomic sequences using CHAOS and DIALIGN as described in

- M. Brudno, M. Chapman, U. Göttingen, S. Balzoglou, B. Morgenstern (2003) [Fast and sensitive multiple alignment of large genomic sequences](#) *BMC Bioinformatics* 4, 66.

If you use our software for your research, please cite this article.

CHAOS is used to rapidly identify strong sequence similarities that serve as anchor points to speed up the DIALIGN alignment procedure. More about the CHAOS/DIALIGN procedure is available [here](#).

**Program Input:**

Upload sequences in multiple **FASTA** format (as a TLF1 file, not MS-Word, not HTML)

Your email address

This web server calculates pair-wise and multiple alignments of genomic sequences using [CHAOS](#) and [DIALIGN](#) as described in

- M. Brudno, M. Chapman, U. Göttingen, S. Balzoglou, B. Morgenstern (2003) [Fast and sensitive multiple alignment of large genomic sequences](#) *BMC Bioinformatics* 4, 66.

If you use our software for your research, please cite this article.

CHAOS is used to rapidly identify strong sequence similarities that serve as anchor points to speed-up the DIALIGN alignment procedure. More about the CHAOS/DIALIGN procedure is available [here](#).

**Program Input**

Upload sequences in multiple **FASTA** format (as a TLF1 file, not MS-Word, not HTML)

Your email address

**Program Output.**

Our server creates four different output files from your input sequence set. The full alignment is returned in **DIALIGN format**. In addition, a **gap-free segment pairs** created by DIALIGN is returned, as well as **anchor points** created by CHAOS and the full alignment in **FASTA format**.

You will receive an email containing the URL of the chosen output files. These files can be accessed during the next 5 days. If small input output alignment is shown on the screen in DIALIGN format

# CHAOS + DIALIGN [manual]

## Alignment of genomic sequences using CHAOS and DIALIGN:

DIALIGN is a widely used software program for pair-wise and multiple alignment of DNA and protein sequences. It assembles global alignment from local pair-wise alignments, so called *fragments*. A number of recent studies have used DIALIGN to discover functional sites sequences, e.g.

While DIALIGN alignments are generally of high quality, the original program is too slow for large-scale sequence comparison. One way of speeding up DIALIGN without compromising on alignment quality is to use an *anchors-alignment* procedure. Here, a fast local alignment search tool identifies regions of high sequence similarity among the input sequences. These similarities are then used as *anchor points* to reduce the search space and run DIALIGN.

On our web server, we are using the program **CHAOS** to identify such anchor points. CHAOS has been developed by [Mike Brudno](#) at Stanford. In our experience, anchor points created by CHAOS speed up DIALIGN by one to two orders of magnitude without affecting the quality of



CHAOS + DIALIGN: manual - Netscape

New Tab | CHAOS + DIALIGN: manual

- H. Göttgens, J. Barton, M. Chapman, A. Sinclair, H. Knudsen, D. Grafham, J. Gilbert, J. Rogers, D. Bentley, and A. Green (2002) Transcriptional regulation of the stem cell leukemia gene (SCL) comparative analysis of five vertebrate SCL loci. *Genome Res.* 12:749-759.
- J. Filch, S. Gardner, J. Kuzmowski, S. Kurz, R. Myers, J. Liu, I. Skrzak, F. Vitale, A. Zimka, and P. McCready (2002) Rapid Development of Nucleic Acid Diagnostics. *Proceedings of the LCCC*, 90: 1700-1721.

While DIALIGN alignments are generally of high quality, the original program is too slow for large scale sequence comparison. One way of DIALIGN without compromising on alignment quality is to use an *anchored-alignment* procedure. Here, a fast local alignment search tool identifies high sequence similarity among the input sequences. These similarities are then used as *anchor points* to reduce the search space and run DIALIGN.

On our web server, we are using the program CHAOS to identify such anchor points. CHAOS has been developed by Mike Brudno at Stanford. In our experience, anchor points created by CHAOS speed-up DIALIGN by one to two orders of magnitude without affecting the quality of alignments. Details are described

In a first step, CHAOS is applied to create a list of anchor points. Then DIALIGN is run on the input sequences using some new options alignment that are described in

- H. Morgenstern, O. Hinner, S. Abdeldaim, I. Haase, K. Mayer, A. Dress, H.-W. Mewes (2002) Four Decades by Genomic Sequences Alignment. *Bioinformatics* 18, 777-787.

CHAOS + DIALIGN: manual - Netscape

New Tab | CHAOS + DIALIGN: manual

**Input sequence file:**

CHAOS/DIALIGN requires a single ASCII file containing the sequences to be aligned in FASTA format

For each sequence, the first line starts with ">" and contains the name of the sequence.

**Program Output:**

Our web server creates different output files containing

- A full alignment of the input sequences in **DIALIGN format**.
- The same alignment in **FASTA format**.
- A list of the **fragments**, i.e. the pair-wise local gap-free alignments that DIALIGN uses to create the output alignment
- A list of the anchor points created by CHAOS.

CHAOS + DIALIGN: manual - Netscape

New Tab | CHAOS + DIALIGN: manual

**Program Output:**

Our web server creates different output files containing

- A full alignment of the input sequences in **DIALIGN format**.
- The same alignment in **FASTA format**.
- A list of the **fragments**, i.e. the pair-wise local gap-free alignments that DIALIGN uses to create the output alignment.
- A list of the anchor points created by CHAOS.

**This is DIALIGN alignment format:**

```

dog_114      20866  AAGAGCGTGTGAT CTGAGAGGAAA GTTAAAGTAT AGLTGTGCTT TCTTTAGCAG
hum_114      21780  AAGAGCGTGTGAT CTGAGAGGAAA GTTAAAGTAT AGLTGTGCTT TCTTTAGCAG
mus_114      21290  ----- --SSGCGGA GCTGATGCTT ACCTGTGCTT TCTTTAGCAG
          1111111111 1111022222 2222222222 2222222222 2222222222

dog_114      20615  ATGAGATATA uGAG---TAC ACCAGTGTGAG CATGAGCTTC TCGAGCTCTA
hum_114      21780  ATGAGATATA uGAG---CAC ACCAGTGTGAG CATGAGCTTC TCGAGCTCTA
mus_114      31427  ATGAGATATA GGGGAGGAG AGAGTGTGAG CATGAGCTTC TCGAGCTCTA
          2222222000 0222000010 0034444444 4444444444 4444444444

dog_114      20862  AAGAGCGTGTGAT CTGAGAGGAAA AGTGTGTGAG CTTTGAAGTTC CAGAGCTCTG
hum_114      21927  AAGAGCGTGTGAT CTGAGAGGAAA AGTGTGTGAG CTTTGAAGTTC CAGAGCTCTG
mus_114      21477  AAGTGTGTGAT GACTTAAAGCA AATGTGTGAGC CTTTGAAGTTC CAGAGCTCTG
          4444444444 4433666666 6666666666 6666662229 9999999999
  
```

CHAOS + DIALIGN: manual - Netscape

New Tab | CHAOS + DIALIGN: manual

```

dog_114      20866  AAGAGCGTGTGAT CTGAGAGGAAA GTTAAAGTAT AGLTGTGCTT TCTTTAGCAG
hum_114      21780  AAGAGCGTGTGAT CTGAGAGGAAA GTTAAAGTAT AGLTGTGCTT TCTTTAGCAG
mus_114      21290  ----- --SSGCGGA GCTGATGCTT ACCTGTGCTT TCTTTAGCAG
          1111111111 1111022222 2222222222 2222222222 2222222222

dog_114      20615  ATGAGATATA uGAG---TAC ACCAGTGTGAG CATGAGCTTC TCGAGCTCTA
hum_114      21780  ATGAGATATA uGAG---CAC ACCAGTGTGAG CATGAGCTTC TCGAGCTCTA
mus_114      31427  ATGAGATATA GGGGAGGAG AGAGTGTGAG CATGAGCTTC TCGAGCTCTA
          2222222000 0222000010 0034444444 4444444444 4444444444

dog_114      20862  AAGAGCGTGTGAT CTGAGAGGAAA AGTGTGTGAG CTTTGAAGTTC CAGAGCTCTG
hum_114      21927  AAGAGCGTGTGAT CTGAGAGGAAA AGTGTGTGAG CTTTGAAGTTC CAGAGCTCTG
mus_114      21477  AAGTGTGTGAT GACTTAAAGCA AATGTGTGAGC CTTTGAAGTTC CAGAGCTCTG
          4444444444 4433666666 6666666666 6666662229 9999999999
  
```

- Names of the aligned sequences are shown on the left
- Numbers on the left hand side of the alignment denote the position of the first residue in a line within the respective sequence.
- Capital letters denote aligned residues; i.e. residues involved in at least one of the 'fragments' (= aligned segment points) the alignment. Lower-case letters denote residues not belonging to any of these selected 'fragments'. They are *not* considered to be DIALIGN. Thus, if a lower-case letter is standing in the same column with other letters, this is pure chance; these residues are *not* considered homologous.
- Numbers below the alignment roughly reflect the degree of local similarity among the sequences. More precisely: They represent the sum of fragments connecting residues at the respective position. The numbers are normalized such that every position gets a value between 0 and 1. Thus, these numbers reflect the *relative* degree of similarity within an alignment, since in *every* alignment, the region of maximum similarity score of 1.

CHAO5 + DIALIGN: manual - Netscape

This is FASTA alignment format:

```

>HTL2
IIdapDLFSOGS-----PQKAYVLDQDIL--QDITLFSWQDLSA
GGCELLALGGLGAAWESLNTFSDSTIIEELSDIAGLQEDDN
-----GTT-----GALAKWVLAQEFYVIAVYKCHT-----KQWYDSTY
NFVYTRGLLAp1-----
-----
>MMV
pDdntWYUCSSLLURQKQKAGAAVITKTCVIMVALLSD I SA
QAPPLTALQATMAPGK--LVVYVQVYVATATITNSVRSRLLVQF
SRIIFQDDE--IILILALFSPFSLIIFQFSLQ-----KSDLSAAS
WQWQDQARLQDITTFPDSLIL-----
-----
>FCOI
TPlLQVTAHQ-----PFWHTQWHTQWNR--KSTLALPILH----
--MELALQFALALQWQW--IISTDQ-----
-----SIVLSR-----IISTDQ-----KYISFFHLGQDNWI
LWGISFVYVLSALNADUUSSTGIGADTFLKADTIPQSTGQJYDQDQ
WQWQDQARLQDITTFPDSLIL-----
-----
>FCOI
ALkuvEITFDQSCGNPQGGVQAILRYRRE--LFTSAGVLT--TN
NRMELMALVLEALKERKEVILSTDSQYVROGIIQWIRHKKKQKRTAD
KQVYQVQVWQWQKQKQKQKQKQKQKQKQKQKQKQKQKQKQKQK
NFKQDITLILKAWPTLQdrgypp-----
-----

```

Note that, as in the DIALIGN format, **only UPPER CASE letters are considered to be aligned**, lower case letters are NOT aligned. Please take into account if you run other programs on the DIALIGN output alignment.

**Fragment file returned by DIALIGN:**

The fragments (aligned gap from segment pairs) used by DIALIGN to assemble the alignment are returned in the following format:

1) seq1	2) seq2	3) beg1	4) beg2	5) len1	6) len2	7) wgt	8) cov1	9) cov2	10) iter	11) cons			
1)	seq1	2	seq2	189955	178118	90	wgt:	42.00	cov1:	107.68	iter:	1	cons:
2)	seq1	1	seq2	201612	189948	90	wgt:	39.98	cov1:	108.66	iter:	1	cons:
3)	seq1	1	seq2	20700	22890	100	wgt:	98.95	cov1:	104.80	iter:	1	cons:
4)	seq1	1	seq2	109848	109798	50	wgt:	37.82	cov1:	104.30	iter:	1	cons:
5)	seq1	2	seq2	19403	52648	50	wgt:	39.08	cov1:	104.27	iter:	1	cons:

CHAO5 + DIALIGN: manual - Netscape

**Fragment file returned by DIALIGN:**

The fragments (aligned gap from segment pairs) used by DIALIGN to assemble the alignment are returned in the following format:

1) seq1	2) seq2	3) beg1	4) beg2	5) len1	6) len2	7) wgt	8) cov1	9) cov2	10) iter	11) cons			
1)	seq1	2	seq2	189955	178118	90	wgt:	42.00	cov1:	107.68	iter:	1	cons:
2)	seq1	1	seq2	201612	189948	90	wgt:	39.98	cov1:	108.66	iter:	1	cons:
3)	seq1	1	seq2	20700	22890	100	wgt:	98.95	cov1:	104.80	iter:	1	cons:
4)	seq1	1	seq2	109848	109798	50	wgt:	37.82	cov1:	104.30	iter:	1	cons:
5)	seq1	2	seq2	19403	52648	50	wgt:	39.08	cov1:	104.27	iter:	1	cons:

For each fragment, the file specifies the involved sequences (seq1), the starting positions of the fragment in these sequences (beg1), the fragment length (len1), the *weight score* of the fragment (wgt), its *overlap weight* (cov1), the iteration step during the program in which the fragment has been found (iter), and information about *consistency* of the fragment (cons (0 = success)).

**Anchor points produced by CHAO5:**

Anchor points produced by CHAO5 are printed in the following format:

1) seq1	2) seq2	3) beg1	4) beg2	5) len1	6) len2	7) wgt	8) cov1	9) cov2	10) iter	11) cons
1)	seq1	2	seq2	178	265	1	4165.000000			
2)	seq1	1	seq2	220	225	1	4105.000000			
3)	seq1	1	seq2	78	162	1	6297.000000			
4)	seq1	1	seq2	157	244	1	6297.000000			
5)	seq1	2	seq2	302	289	1	2619.000000			
6)	seq1	2	seq2	316	333	1	2619.000000			
7)	seq1	1	seq2	178	156	1	6305.000000			

The first two entries are the sequences involved, entries 3 and 4 are the starting points in the respective sequences. Entry 5 is the length of a segment pair, and the last entry is a quality score calculated by CHAO5. This is used to prioritize anchor points in case contradicting anchors are found. In this case, anchors with high scores are accepted first, anchors with lower scores are used only if they are consistent with the previous higher-scoring anchors.

**This is PHYLIP tree format:**

```

((HTL2:0.111024,
MMV:0.078471)
:0.082554)
:0.121218,
FCOI:0.232442)
;
```

CHAO5 + DIALIGN: manual - Netscape

For each fragment, the file specifies the involved sequences (seq1), the starting positions of the fragment in these sequences (beg1), the fragment length (len1), the *weight score* of the fragment (wgt), its *overlap weight* (cov1), the iteration step during the program in which the fragment has been found (iter), and information about *consistency* of the fragment (cons (0 = success)).

**Anchor points produced by CHAO5:**

Anchor points produced by CHAO5 are printed in the following format:

1) seq1	2) seq2	3) beg1	4) beg2	5) len1	6) len2	7) wgt	8) cov1	9) cov2	10) iter	11) cons
1)	seq1	2	seq2	178	265	1	4165.000000			
2)	seq1	1	seq2	220	225	1	4105.000000			
3)	seq1	1	seq2	78	162	1	6297.000000			
4)	seq1	1	seq2	157	244	1	6297.000000			
5)	seq1	2	seq2	302	289	1	2619.000000			
6)	seq1	2	seq2	316	333	1	2619.000000			
7)	seq1	1	seq2	178	156	1	6305.000000			

The first two entries are the sequences involved, entries 3 and 4 are the starting points in the respective sequences. Entry 5 is the length of a segment pair, and the last entry is a quality score calculated by CHAO5. This is used to prioritize anchor points in case contradicting anchors are found. In this case, anchors with high scores are accepted first, anchors with lower scores are used only if they are consistent with the previous higher-scoring anchors.

**This is PHYLIP tree format:**

```

((HTL2:0.111024,
MMV:0.078471)
:0.082554)
:0.121218,
FCOI:0.232442)
;
```

Trees can be visualized using the *d3drawtree* program contained in the [PHYML](#) software package.

Back to [CHAO5/DIALIGN submission form](#).

