

Programas de Alinhamento

Departamento de Genética FMRP- USP

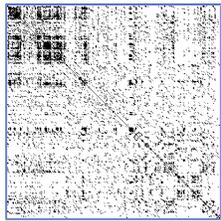
Alynne Oya Chiromatzo
alynne@lgmb.fmrp.usp.br

Sumário

- Introdução
- Programas de alinhamento para buscas em base de dados
 - Fasta
 - Blast
- Programa para alinhamento múltiplo de seqüências
 - Clustal

• Alinhamento de pares de seqüências

Programação dinâmica



Matriz de pontos

MATRIZ N-W	GAP	A	A	G
GAP	0	10	20	30
A	1	10	0	10
G	2	20	10	1
T	3	30	20	11
A	4	40	30	20
C	5	50	40	30

A = AGTAC
 B = A--AG

Alinhamento Global (Needleman & Wunsch)

	G	A	A	T	T	C	A	G	T	T	A
G	0	0	0	0	0	0	0	0	0	0	0
G	0	5	1	0	0	0	0	0	5	1	0
G	0	5	2	0	0	0	0	0	5	2	0
A	0	1	10	7	3	0	0	5	1	2	0
T	0	0	6	7	12	8	4	1	2	6	8
C	0	0	2	3	8	9	13	9	5	2	4
G	0	5	1	0	4	5	9	10	14	10	6
A	0	1	10	6	2	1	5	14	10	11	7

GAATTCAG GAATTCAG
 |||||
 GGA-TC-G GGAT-C-G

GAATTC-A GAATTC-A
 |||||
 GGA-TCGA GGAT-CGA

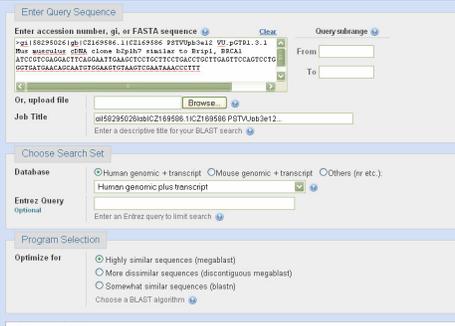
Alinhamento Local (Smith & Waterman)

• Alinhamento de pares de seqüências

Método de busca por palavras



FASTA
(Fast Alignment Search Tool All)



BLAST
(Basic Local Alignment and Search Tool)



Grupo de Bioinformática
Depto. de Genética - FMRP/USP

Introdução



Faculdade de Medicina
USP - RIBEIRÃO PRETO - SP

• Alinhamento múltiplo de seqüências

Seqüências de entrada

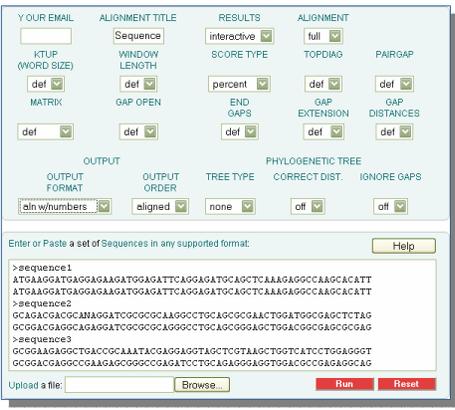
```
CGAGGCCGAAAG
GATCCTGCAGAG
CGATCGCGCCG
AGGAGAGATG
```

Combinções

```
CGAGGCCGAAAGCC
GATCCTGCAGAGGG
CGAGGCCGAAAGCC GATCCTGCAGAGGG
CGATCGCGCCGAAAG CGATCGCGCCGAAAG
CGAGGCCGAAAGCC GATCCTGCAGAGGG CGATCGCGCCGAAAG
AGGAGAGATGAGG AGGAGAGATGAGG AGGAGAGATGAGG
```

Alinhamento final

```
CGAGGCCGAAAGCC
CGATCGCGCCGAAAG
GATCCTGCAGAGGG
AGGAGAGATGAGG
```



Clustal



Grupo de Bioinformática
Depto. de Genética - FMRP/USP

Programas de alinhamento para buscas em base de dados



Faculdade de Medicina
USP - RIBEIRÃO PRETO - SP

• FASTA

- 1985 FASTP
Alinhamento de seqüências de proteínas
- 1988 FASTA
Pacote de softwares para alinhamento de seqüências de ADN e proteínas



Dr. William R. Pearson
Bioquímico
Prof. de Bioquímica e
Ciência da Computação
Universidade de Virgínia



Ms. David J. Lipman
Biólogo
Diretor do NCBI
National Center for
Biotechnology Information

- FASTA – Web

<http://www.ebi.ac.uk/fasta33/>

Fasta - Protein Similarity Search

Provides sequence similarity searching against nucleotide and protein databases using the Fasta programs. Fasta can be very specific when identifying long regions of low similarity especially for highly diverged sequences. You can also conduct sequence similarity searching against complete [proteome](#) or [genome](#) databases using the [Fasta programs](#).

[Download Software](#)

PROGRAM	DATABASES	RESULTS	SEARCH TITLE	YOUR EMAIL
fasta3	Protein UniProt Knowledgebase UniProtKB/Swiss-Prot UniProt Clusters 100%	interactive	Sequence	

MATRIX	GAP OPEN	GAP EXTEND	KTUP	EXPECTATION UPPER VALUE	EXPECTATION LOWER VALUE
BLOSUM62	-10	-2	2	10.0	default

DNA STRAND	HISTOGRAM	MOLECULE TYPE
none	no	Protein

SCORES	ALIGNMENTS	SEQUENCE RANGE	DATABASE RANGE	FILTER	STATISTICAL ESTIMATES
50	50	START-END	START-END	none	Regress

Enter or Paste a **PROTEIN** Sequence in any format: [Help](#)

Upload a file: [Browse...](#) [Run Fasta3](#) [Reset](#)

- FASTA – Web

fasta3: procura seqüências similares em uma biblioteca de proteínas ou de seqüências de DNA

fastx/y3: compara uma seqüência de DNA com um banco de proteínas nos dois sentidos

fastf3: compara uma mistura de peptídeos com um banco de dados de proteínas

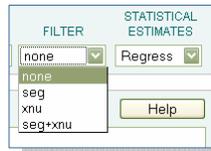
fasts3: compara peptídeos que estão ligadas com um BD de proteínas

ssearch3: procura por seqüências similares em um biblioteca de proteína ou de seqüências de DNA

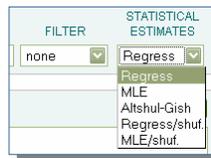
PROGRAM	DATABASES
fasta3	Protein
fasta3	UniProt Knowledgebase
fastx3	UniProtKB/Swiss-Prot
fastf3	UniProt Clusters 100%
fasts3	
ssearch3	

GAP OPEN	GAP EXTEND
-10	-2

• FASTA – Web



seg: mascara regiões de baixa complexidade
xnu: mascara regiões contendo repetições internas de curta periodicidade
seg + xnu: combinação dos anteriores



Regress: usa regressão ponderada da pontuação média contra o tamanho das seqüências da biblioteca
MLE: usa a estimativa de máxima parcimônia para λ
Regress/MLE huf: estima os parâmetros estatísticos a partir de cópias aleatórias de cada biblioteca de seqüências

• FASTA – Web

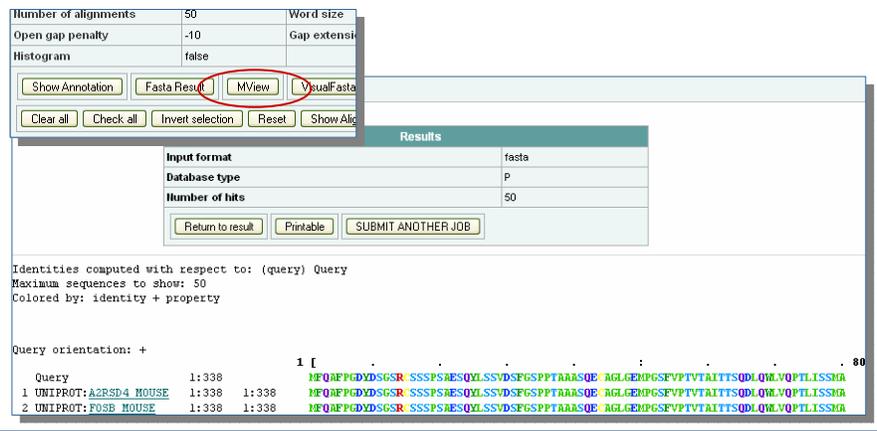
Number of alignments	50	Word size	
Open gap penalty	-10	Gap extension	
Histogram	false		

Source	Length	Identity%	Similar%	Overlap	E()
FBJ osteosarcoma oncogene	338	100.000	100.000	338	1.8e-104

```

ID A2RSD4_MOUSE Unreviewed; 338 AA.
AC A2RSD4;
DT 06-MAR-2007, integrated into UniProtKB/TrEMBL.
DT 06-MAR-2007, sequence version 1.
DT 24-JUL-2007, entry version 4.
DE FBJ osteosarcoma oncogene B.
GN Name=Fosb;
OS Mus musculus (Mouse).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Euarchontoglires; Glires; Rodentia; Sciurognathi;
OC Muridea; Muridae; Murinae; Mus.
OX NCBI_TaxID=10090;
RN [1]
RP NUCLEOTIDE SEQUENCE.
RC TISSUE=Brain;
RX MEDLINE=22388257; PubMed=12477932; DOI=10.1073/pnas.242603899;
RA Strausberg R.L., Feingold E.A., Grouse L.H., Derge J.G.,
RA Klausner R.D., Collins F.S., Wagner L., Shenmen C.M., Schuler G.D.,
    
```

• FASTA – Web



Number of alignments: 50 Word size: []
 Open gap penalty: -10 Gap extension: []
 Histogram: false

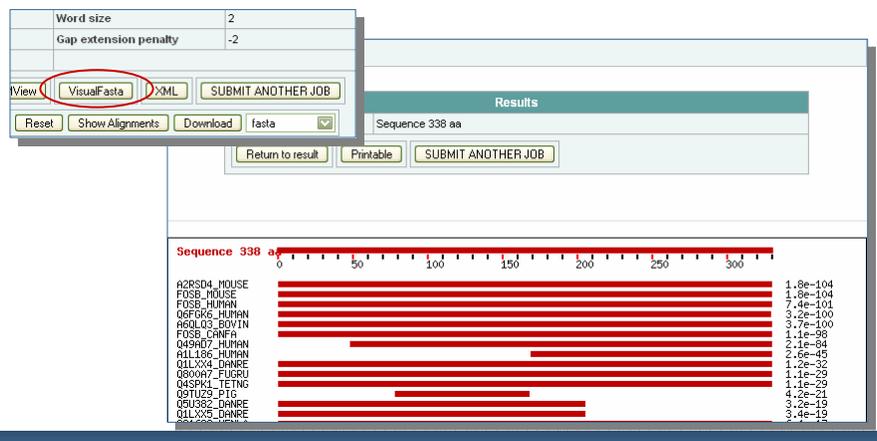
Results
 Input format: fasta
 Database type: P
 Number of hits: 50

Identities computed with respect to: (query) Query
 Maximum sequences to show: 50
 Colored by: identity + property

Query orientation: +

Query	1:338	1 [80
1 UNIPROT:A2RSD4_MOUSE	1:338	1:338	MFQAFPGDYDGSRCSSSPSAESQYLSVDSFGSPPTAASUQCAGLGEMPGSFVPTVTAITTSQDLQMLVQPTLISSKA
2 UNIPROT:FO5B_MOUSE	1:338	1:338	MFQAFPGDYDGSRCSSSPSAESQYLSVDSFGSPPTAASUQCAGLGEMPGSFVPTVTAITTSQDLQMLVQPTLISSKA

• FASTA – Web



Word size: 2
 Gap extension penalty: -2

 fasta

Results
 Sequence 338 aa

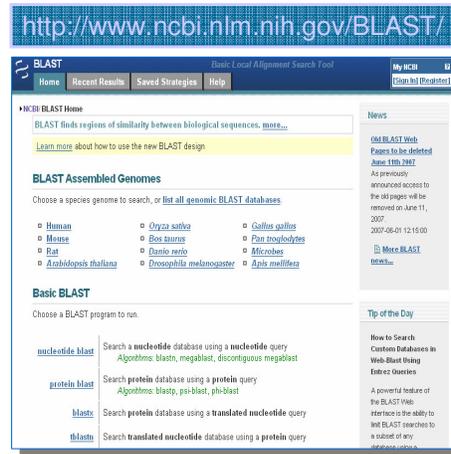
Sequence 338 aa

Sequence	Score
A2RSD4_MOUSE	1.8e-104
FO5B_MOUSE	1.8e-104
FO5B_HUMAN	7.4e-101
Q6F0K6_HUMAN	5.2e-100
AFUL03_BOVIN	3.7e-100
FO5B_CANFA	4.1e-98
Q49A07_HUMAN	2.1e-84
AL1196_HUMAN	2.6e-46
Q1LXV4_DANRE	2.6e-42
Q800A7_FUGRU	1.1e-26
Q4SRK1_TETNG	1.1e-26
Q9TUZ5_PIG	4.2e-21
Q5U352_DANRE	6.2e-19
Q1LXV4_DANRE	6.4e-19
Q1LXV4_DANRE	6.4e-19

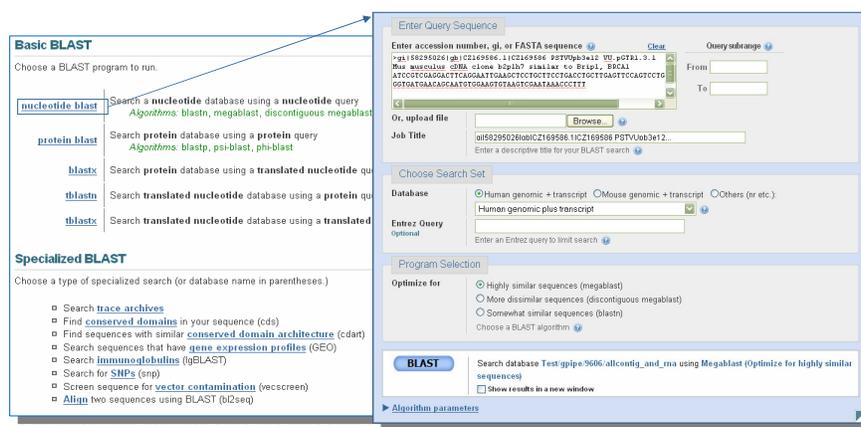
• BLAST

1990 BLAST

- Alinha seqüências de ADN ou proteínas com as seqüências da base de dados
- Mais rápido que o FASTA



• BLAST – Web



• BLAST – Web

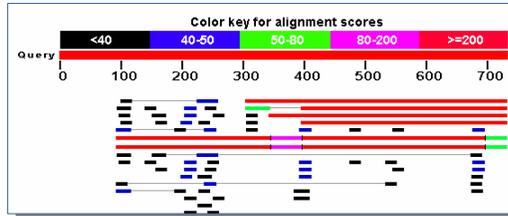
The screenshot shows the 'Choose Search Set' dropdown menu in the BLAST web interface. The menu is open, showing a list of databases. The 'Genomic plus Transcript' option is highlighted. Below it, a list of 'Other Databases' is visible, including 'Nucleotide collection (nr/nt)', 'Reference mRNA sequences (refseq_ma)', 'Reference genomic sequences (refseq_genomic)', 'Expressed sequence tags (est)', 'Non-human, non-mouse ESTs (est_others)', 'Genomic survey sequences (gss)', 'High throughput genomic sequences (HTGS)', 'Patent sequences (pat)', 'Protein Data Bank (pdb)', 'Human ALU repeat elements (alu_repeats)', 'Sequence tagged sites (dsts)', 'Whole-genome shotgun reads (wgs)', and 'Environmental samples (env_nt)'. The 'BLAST' button is visible at the bottom left of the interface.

refseq: banco de dados de seqüências não redundantes
gss: seqüências de origem genômica e não por cDNA
HTGS: *High Throughput Genomic Sequences* são dados de seqüências genômicas não terminadas

• BLAST – Web

The screenshot shows the 'Algorithm parameters' section of the BLAST web interface. The 'General Parameters' section is expanded, showing settings for 'Max target sequences' (100), 'Short queries' (Automatically adjust parameters for short input sequences), 'Expect threshold' (10), and 'Word size' (3). The 'Scoring Parameters' section is also visible, showing settings for 'Matrix' (BLOSUM62), 'Gap Costs' (Existence: 11 Extension: 1), and 'Compositional adjustments' (Composition-based statistics). The 'Filters and Masking' section is also visible, showing settings for 'Filter' (Low complexity regions) and 'Mask' (Mask for lookup table only, Mask lower case letters). The 'BLAST' button is visible at the bottom left of the interface.

• BLAST – Web



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

Sequences producing significant alignments:
 (Click headers to sort columns)

Accession	Description	Max score	Total score	Query coverage	E value	Max ident	Links
Transcripts							
NM_002287.3	Homo sapiens leukocyte-associated immunoglobulin-like receptor	803	803	58%	0.0	99%	U E G M
NM_021706.2	Homo sapiens leukocyte-associated immunoglobulin-like receptor	624	704	51%	3e-176	100%	U E G M
NM_002288.3	Homo sapiens leukocyte-associated immunoglobulin-like receptor	478	478	53%	5e-132	90%	U E G M
NM_021270.2	Homo sapiens leukocyte-associated immunoglobulin-like receptor	383	383	46%	2e-103	89%	U E G M
XM_001125918.1	PREDICTED: Homo sapiens hypothetical protein LOC727872 (LOC	40.1	40.1	2%	3.5	100%	G M

• BLAST – Local

```
$ blastall -p blastn -o hs_est -i entrada.fa -o out.entrada
```

-p Nome do programa
blastp: seq. de proteína contra BD de seq. de proteína
blastn: seq. de nucleotídeos contra banco de nucleotídeos
blastx: seq. de nucleotídeos nos 6 quadros de leitura contra BD de seq. de proteína
tblastn: seq. de proteína contra BD de nucleotídeos dinamicamente traduzidos nos 6 quadros de leitura
tblastx: seq. de nucleotídeos contra um BD de nucleotídeos traduzidos

-d Nome do banco de dados

• BLAST – Local

```

BLASTP 2.2.2 [Jan-08-2002]

Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer,
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.

Query= g1|730028|sp|P40692|MLH1_HUMAN MUTL PROTEIN HOMOLOG 1 (DNA
MISMATCH REPAIR PROTEIN MLH1)
(756 letters)

Database: sp35
69,113 sequences; 25,083,768 total letters

Searching..... done

Sequences producing significant alignments:
Score E-Value
g1|604369|sp|P40692|MLH1_HUMAN MUTL PROTEIN HOMOLOG 1 (DNA MISMA... 1293 0.0
g1|825572|sp|P38920|MLH1_YEAST MUTL PROTEIN HOMOLOG 1 (DNA MISMA... 437 e-122
g1|1573016|sp|P44494|MUTL_HAEIN DNA MISMATCH REPAIR PROTEIN MUTL... 190 4e-48
g1|304915|sp|P23367|MUTL_ECOLI DNA MISMATCH REPAIR PROTEIN MUTL... 190 5e-48

Query: 128 ASYSDGKLGKPKPKPCAGNQGQITVEDLFYNIATRRLKLNKPNSEYVSKILEVGRYSVHN 187
SY+HGK+ PKP AG GT I VEDLFNI +R +AL++ +EY KIL+VVGRY+++
Sbjct: 125 VSYAEGKMLESPPKPVAGKDGTTILVEDLFFNIPSLRLRALRSHNDEYSKILDVVGRIYIHS 184

Query: 188 AGISFSVKKQGETVADVRTLPNASTVDNIRSFIGNAVSRELIEI---GCEDKTLAFKMNG 244
I FS KK G++ + P+ + D IR++F +V+ LI ED L ++G
Sbjct: 185 KDIGFSCKKFGDSNYSLSVKPSYTVQDRIRTVFNKSVASNLITFHI SKVEDLNLE-SVDS 243
    
```

• Clustal

- 1994 Clustal
- Clustal: atribui pesos iguais a todas seqüências
- ClustalW: atribui diferentes pesos às seqüências
- ClustalX: proporciona uma interface gráfica para o ClustalW

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

YOUR EMAIL: ALIGNMENT TITLE: RESULTS: ALIGNMENT:

KTUP (WORD SIZE): WINDOW LENGTH: SCORE TYPE: TOPDIAG: PAIRGAP:

MATRIX: GAP OPEN: END GAPS: GAP EXTENSION: GAP DISTANCES:

OUTPUT FORMAT: OUTPUT ORDER: TREE TYPE: PHYLOGENETIC TREE CORRECT DIST: IGNORE GAPS:

Enter or Paste a set of Sequences in any supported format:

```

>sequence1
ATGAAAGATGAGAGAAAGATGAGATTCAAGAGATGCAAGCTCAAAGAGGCCAAGCACATT
ATCAAGGATGAGGAAAGATGAGATTCAAGAGATGCAAGCTCAAAGAGGCCAAGCACATT
>sequence2
GCAGACAGCCMNAAGATCGCCGCGCAAGCCGTCGCGCCGAACCTGGATGGCAGCTCTAG
GCGACAGGCGAGAGATCGCCGCGCAAGCCGTCGCGCCGAACCTGGATGGCAGCTCTAG
>sequence3
GCGAAGAGGCTGACCCGAAATACAGAGAGGTAGCTCTGTAAGCTGGTCTTCCTGAGGGT
GCGACAGGCGCGAAGAGCCGCGCAAGATCTGCAAGAGGAGTGGACGCCAGAGGCGAG
    
```

Upload a file:

• Clustal – Web

ClustalW Results	
Results of search	
Number of sequences	3
Alignment score	176
Sequence format	Pearson
Sequence type	aa
ClustalW version	1.83
JaView	<input type="button" value="Start JaView"/>
Output file	clustalw-20070830-04263186.output
Alignment file	clustalw-20070830-04263186.aln
Guide tree file	clustalw-20070830-04263186.dnd
Your input file	clustalw-20070830-04263186.input
<input type="button" value="SUBMIT ANOTHER JOB"/>	

Scores Table				
Sort by: <input type="button" value="Sequence Number"/> <input type="button" value="View Output File"/>				
SeqA Name	Len(aa)	SeqB Name	Len(aa)	Score
1 sequence1	120	2 sequence2	120	15
1 sequence1	120	3 sequence3	120	17
2 sequence2	120	3 sequence3	120	13

PLEASE NOTE: Some scores may be missing from the above table if the alignment

Sort by: <input type="button" value="Sequence Number"/> <input type="button" value="View Output File"/>				
---	--	--	--	--

• Clustal – Web

Alignment

CLUSTAL W (1.83) multiple sequence alignment

```

sequence1  -----NRALS AELQSVTEQLSDGGKNSAEVEKLRRLKGHENEELQIALEEAEA 49
sequence3  LENEAKASEDKAQRAHAEVARLHSELNSAQEATSTAEKSRQLVSKQVADLQSRLEDAEA 60
sequence2  -----HMELESQLESSNRVAEESQKMEKIQAIKELQSHIDDESFG 42
           : : :* . . . . : : : * : : : * * : : .

sequence1  LEQ-EECKLLKYOLEYTLRQSSDRK-----LSEKDEELEGKRMHOPQMESLQNTIDS 102
sequence3  GKGGLNQLRLEORIMELESVDVTE-----APKGAATAKARK-SEKKVELAFTIED 113
sequence2  RDD-MRDSASRSERPANDLAVQLDEARVALEQAERARKLAENKSENSDRVAELQALYNN 101
           . . . : : . : * . * . . : : . : : * : :

sequence1  ESRKAEQQKLRKKYDAD- 120
sequence3  EHKREP----- 120
sequence2  VANAKAEGDYHSLQEEIED 120
           . :
    
```

PLEASE NOTE: Showing colors on large alignments is slow.

* Idêntico
 . semi conservada
 : conservada

- Clustal – Local

```
ftp://ftp.ebi.ac.uk/pub/software/unix/clustalw/
```

```
$ wget ftp://ftp.ebi.ac.uk/pub/software/unix/clustalw/  
clustalw1.83.linux.tar.gz  
$ tar -zxvf clustalw1.83.linux.tar.gz
```

- Clustal – Local

```
clustalw -entrada.fa -align -tree
```

-infile	Arquivo de entrada
-align	Faz o alinhamento múltiplo
-bootstrap(=n)	Faz o Bootstrap em uma árvore Neighbour Join NJ (n= número de bootstraps; def. = 1000).
-convert	Converte as seqüências de entrada em um formato diferente
-help ou -check	Mostra uma visão geral sobre os parâmetros
-options	Lista os parâmetros da linha de comando
-tree	Calcula a árvore neighbour join

