**ERPIN DOCUMENTATION** Easy RNA Profile Indentification

Version 4.2.5 (C)2001, 2002, 2003, 2004, 2005

Daniel Gautheret: gautheret@esil.univ-mrs.fr

Andre Lambert: lambert@cpt.univ-mrs.fr

Please cite: Gautheret D. & Lambert A. (2001). Direct RNA definition and identification from multiple sequence alignments

using secondary structure profiles. J. Mol. Biol. 313:1003-1

**Database file** 

1. FILE FORMATS

The sequence database file is in the FASTA format. The file can hold up to 10 million sequences, each less than 300Mb long.

Training set: the .epn file

Training sets are FASTA-like files starting with a secondary structure descriptor that encodes each structure element with a

be encoded by 0.

user-defined number. Numbers that are repeated at two different places stand for helical elements, other numbers stand for single strands. Two or 3 numbering lines can be used where 10 or more numbers are required. No structure element should

2. THE ERPIN COMMAND LINE

[-long|-short|-mute]

[-logzero <logzero>] [-tablen <tablen>]

[-bgn <seqbgn>][-len <range>]

erpin trna.epn coli.fasta -4,4 -nomask

erpin trna.epn coli.fasta -4,11 -nomask

erpin trna.epn coli.fasta -2,2 -nomask

-> Search region includes the whole alignment except for strands 1 and 12

-> Search region includes helix 4 and 7, strands 5, 8 and 9, and the 5' part of stem 20

erpin trna.epn coli.fasta -2,+2 -umask 2 4 7

the training set sequences using the tstat program (see additional Tools).

[-warnings]

[-chrono]

[-pcw <pcw>] [-hpcw <hpcw>]

Regions and Masks (-mask)

[-spcw <spcw>]

-> Search region includes helix 4 and strand 5

-umask i j .. n

Score cutoff (-cutoff)

scores instead.

**Examples:** 

now).

-> Searches region -2 to +2 ignoring strand 8

-> Searches only elements 2, 4 and 7 in region -2 to +2

The default score for any region or mask is 100%.

-> Select hits with a score higher than 20

Multi-level search strategies

argument, as follows:

shown in lowercase.

frequencies to 0.25 each.

pseudocounts.

Example:

Log zero (-logzero)

should be raised to 0.2-0.5. Higher weigths are risky!

tolerance to deviations from the training set sequences.

Length of pre-aligned region (-tablen)

Reverse complement searches (-fwd, -rev)

**Output format (-long, -short, -mute)** 

-mute: only prints final number of hits.

Output format can be set as long, short or mute.

-short: for every hit, prints coordinates and final score.

Selection of database sequences (-seq1, -nseq, -bgn, -len)

- seq1 n: search begins at sequence number n in the database

- nseq n: n sequences in the database are searched

- lowest score in current search (minus epsilon) - highest score in current search (plus epsilon)

This is followed by histogram values (one integer per interval).

./parent2epn.pl <parenthesized alignment>

>secondary structure ( [ are pseudoknots )

The -chrono argument displays the elapsed CPU time after a Erpin run is completed.

- number of score intervals in histogram

- total number of solutions

**Elapsed CPU time (-chrono)** 

random sequences.

3. EXTRA TOOLS

parent2epn.pl

Syntax:

>e. coli

>t. pallidum

readerpin.pl

Syntax:

motifs).

Syntax:

3.9.4

3.9.3

3.9.2

fixed bugs

fixed bugs

E-value computing

A good general strategy could use 3 levels:

- first level to speed up search by selecting a short motif

- third level with the complete region unmasked for output.

However, you won't handle a 16S RNA with just 3 levels!

recover some "lost" solutions by lowering score cutoffs in late stages.

- second level for specificity: extend to a larger motif eliminating false positives

gapped region may cause memory/CPU overflow!

(see also sections "log-zero" and "expected frequencies")

-add i j .. n

-nomask

erpin <training-set>

Sequences should be aligned with gaps, so that all entries have the same length as the secondary structure descriptor. No gaps are tolerated within helical regions.

- Both T and U are accepted. Both uppercase and lowercase are accepted.

Training set example (tRNA):

>structure 

>DA0260 >DA0340

-GGGCGAAUAGUGUCAGC-GGG--AGCACACCAGACUUGCAAUCUGGUAG-GGAGGGUUCGAGUCCCUCUUUGUCCACCA

-GGGCUCGUAGCUCAGC--GGG--AGAGCGCCGCCUUUGCGAGGCGGAGGCCGCGGGUUCAAAUCCCGCCGAGUCCA---

>DA0380 -GGGCCCAUAGCUCAGU--GGU--AGAGUGCCUCCUUUGCAAGGAGGAUGCCCUGGGUUCGAAUCCCAGUGGGUCCA--->DA0420 -GGGCCCAUAGCUCAGU--GGU-AGAGUGCCUCCUUUGCAAGGAGGAUGCCCUGGGUUGGAAUCCCAGUGGGUCCA---

In this training set file, elements 02, 04, 07 and 20 describe the four tRNA helices, other numbers describe single strands.

A Perl script (parent2epn.pl) is provided to convert a parenthesized (Fasta-like) alignment into .epn format. This script will also check (and optionnally correct) gaps in helices and gap-only columns. Limitations: Training sets should not exceed 12,000 sequences. Helices should not exceed 64bp.

**General syntax** erpin [-h] help

<input-file> database file name (fasta) <region> region of interest -nomask|((-mask|-umask|-add) <elt1> ...) level1, [-nomask|((-mask|-umask|-add) <elt1> ...)]

training set file name

level2, default: void

default: 1, SEQ\_MAX\_LEN

pseudo-counts for helices

pseudo-counts for strands

pseudo-counts weight, default: 0.1

idem

default: -long

default: OFF

default: -Eon

default: OFF

default: -20

default: OFF

default: 1024

level.. [-cutoff <cutoff1> <cutoff2> ..] default: 100% default: -dmp [-dmp|-smp][-fwd|-rev|-fwd+rev] default: -fwd+rev

[-globstat|-locstat|-unifstat] default: -globstat [-Eon|-Eoff] [-hist] [-seq1 <seqnb1>][-nseq <nseq>] default: 1, SEQ\_MAX\_NB (all)

like a single-stranded element during search. Another compulsory argument is the Mask. When -nomask is used, the whole region is searched. Other types of masks are exposed later. Examples (using the tRNA training set above):

Compulsory argument <region> contains two comma-separated numbers <r1>,<r2> defining the boundaries of the region of the alignment that will be used for searches. These numbers refer to the structure header of the training set file. When a boundary is a helix, use "plus" or "minus" signs to specify 5' or 3' strand, respectively. The "plus" sign is optional. Output sequence alignments show the defined region only. When a region contains only one strand of a helix, this strand is treated

[-sumf <fname>] substitution matrix file name, default: none

are included or excluded. Masks do not use plus or minus signs for helices. When a helix is referred to, both strands are used. Types of masks: -mask i j .. n : elements i,j,n are excluded from the region

: only elements i,j,n are considered in the region

: all elements of the region are considered

: elements i,j,n are considered in the region, as well

as the elements already considered at the previous stage

Masks are used to restrict searches to certain elements in a region. A mask is followed by numbers indicating which elements

usages. Use multi-level searches in this case (see below). Examples using the tRNA training set above: erpin trna.epn coli.fasta -2,+2 -mask 8

CAUTION with -nomask: When a region is large or contains several gapped strands, this may result in huge memory and CPU

The cutoff argument is used to set the score below which solutions will be discarded. In the absence of mask (-nomask), the score cutoff is that of the whole "Region". When masks are used, a cutoff is provided for each mask. Erpin uses two types of scoring: absolute and relative. Absolute scores are the sum of scores obtained for each secondary structure element, based on the lod-score profiles. An integer or real number following the -cutoff tag is interpreted as an absolute score, unless followed by a "%" sign. To get help in establishing absolute score cutoffs, visualize absolute scores for

Relative scores are expressed as a percentage of training set sequences captured. Theses scores are followed with a "%" sign. A score of 100% for a given region or mask is the lowest score in the training set for this region or mask. So that when

sequences in the training set (hence higher score than 100%). If scores lower than the 100% score are required, use absolute

100% is used as a cutoff, all sequences in the training set are captured. A score of 50% is a score capturing 50% of

erpin trna.epn coli.fasta -2,+2 -nomask -cutoff 20

erpin trna.epn coli.fasta -2,+2 -nomask -cutoff 90% -> Select hits with a score higher than that of 90% of training set sequences E-value (-Eon, -Eoff)

Since Version 3.9, ERPIN computes Expect-values (E-values). For any hit of score S obtained in a given database, The Evalue is the number of hits of same or higher score that can be expected by chance in the same database. It thus provides user with a statistical significance of hits. Typically an E-value of 10e-2 or less is significant. However, many non-biological hits

E-value calculation is turned on by default. It can be turned off using the -Eoff parameter (this was useful mostly with the first implementation of E-value - ERPIN 3.9 through 4.1 - which took a toll on CPU time, but hopefully no one should need this

This is done by applying several masks consecutively. When the command line contains several masks, Erpin will conduct a

first search using the first mask, and continue with the next mask only if a solution bove cutoff has been found with the

might arise with "good" E-values, due to other factors such as the presence of a low complexity region in the motif.

previous mask. Since the search at level n is performed only around solutions found at level n-1 (within distance intervals specified in the input alignment), the search speed is increased. For instance, the following command: erpin trna.epn coli.fasta -2,+2 -umask 20 11 -nomask will run much faster than: erpin trna.epn coli.fasta -2,+2 -nomask

.. and yet produce similar result. This permits to implement multi-level search strategies, where the most significant signatures

Use "-nomask" as the last step if you want to display the complete alignment. Bear in mind however that "-nomask" on a long,

are searched first and the motif is gradually expanded thereafter, thus speeding up database searches. Any number of consecutive masks are allowed. By default, the score cutoff for each mask is 100%. This can be changed using the cutoff

erpin trna.epn coli.fasta -2,+2 -umask 20 11 -nomask -cutoff 10 20

erpin <tr-set> <databank> 1,10 -umask 1 2 -umask 3 4 -umask 5 6

erpin <tr-set> <databank> 1,10 -umask 1 2 -add 3 4 -add 5 6

-> in this example, a cutoff of 10 is used at the first level and a cutoff of 20 at the second level.

Another important thing to understand about multi-level searches: In the following command:

constructing masks incrementally, so that the whole region is matched in the end, such as: erpin <tr-set> <databank> 1,10 -umask 1 2 -umask 1 2 3 4 -umask 1 2 3 4 5 6 To simplify such a command line, one may use the -add parameter. "-add X" just updates the previous mask by unmasking element X from it. Therefore, the previous command can be written:

If the higher search level still contains masked elements (like elements 7 to 10 above), the sequences within these elements will appear UNALIGNED in the final output (left-justified in a space of same size as in training set). Unaligned sequences are

the last step of the search only seeks elements 5 and 6. Although elements 1, 2, 3 and 4 have been detected at previous steps, they are not considered anymore at this stage and will not be included in the finaloutput. Generally, users will prefer

Dynamic/Static Mask Processing (-dmp / -smp) By default, multiple-level searches imply that each level starts with the partial configurations found at the previous level. In the last example above, the search for elements 2 & 3 at step 2 will be performed for each configuration of elements 1 and 2 identified at step 1. This is called "Dynamic Mask Processing", and corresponds to the -dmp option.

Then, this part would be missed at the first level using a cutoff of X. ALWAYS use a lower cutoff for the first stage.

Beware: a 2-level search using a first mask at cutoff X and then the entire region (nomask) at cutoff X is not equivalent to a single level search for the entire region at cutoff X. Indeed, a solution of score X may have a lower score in any of its part.

With the -smp option (Static Mask Processing), each search step is performed independently. In most cases, this results in higher CPU times, since all elements in the current mask must be identified at the same time. A possible use of -smp is to

Single strand and helix profiles are computed based on observed and expected frequencies for each base and base pair. By

default, expected frequencies are those in the search database, averaged over all sequences ("-globstat" option). If the database is very heterogeneous (e.g. a mixture of sequences from different organisms), expected frequencies should be computed independently for each sequence, using the "-locstat" option. Beware however, that short sequences may have highly biased compositions, possibly resulting in spurious high-scoring solutions. The "-unifstat" (uniform) option sets A/T/G/C

simulate what could have been observed in a larger sequence alignment. Pseudocounts require some prior knowledge of "typical" mutation frequencies in RNA molecules. Let's pretend that G often mutates to A in typical RNAs. Then if a column of the initial alignment mostly contains Gs, we can expect that some As should occur too. This is the way pseudocounts work. "Typical" mutation rates for single-strands and base-pairs were evaluated from a 16S/18S ribosomal RNA alignement, and the

Users can set the level of pseudocount to be injected in ERPIN profiles, or pseudocount weight. A high pseudocount weight is necessary when training sets are really poor, but this will affect search specificity. In the extreme situation where profiles are made of 100% pseudocounts, ERPIN searches would just produce noise. Internally, pseudocount weights are comprised between 0 and 1 with a default value of 2x10e-3. At weight=0 no pseudocount is used and at weight=1, profiles are 100%

Use the -pcw parameter to set the pseudocount weight. For convenience and backward compatibility, the default value of -pcw is 0.1 (corresponding to an internal pseudocount weight of 2x10e-3). For training sets with fewer than 10 sequences, pcw

Pseudocounts (-pcw, -hpcw, -spcw, -sumf) Training sets with few sequences (i.e.20 or less) result in hollow weight matrices that strongly penalize any variation from training set sequences and thus affects the sensitivity of detection (see also "log zero"). To solve this problem, ERPIN uses pseudocounts (from version 4.2). Pseudocounts introduce some articifical base or base-pair counts in the weigth matrix, that

detailed counting procedure resembles that of Henikoff and Henikoff (CABIOS 1996, 12:135).

erpin <tr-set> <databank> 1,10 -nomask -sumf MYSUM.dat

The logzero parameter is irrelevant when pseudocounts are used (default behavior)

-long: for every hit, prints scores at each search level, coordinates and complete sequence.

Database searches can be restricted to subsets of sequence files, specified as follows:

Expected/background nucleotide frequencies (-locstat, -globstat, -unifstat)

Expert users may want to separately set pseudocount weights for helices and single strands. Parameters -hpcw (helices) and spcw (single strands) will do just this. The SUM.dat file in the ERPIN distribution contains the default substitution matrices used for pseudocount calculation. There is one 16x16 matrix for base-pairs and one 4x4 matrix for single strands. Users can provide an alternate substitution matrix file using the -sumf parameter. Base order is A,T,G,C in the 4x4 matrix and A:A,A:T,A:G,A:C,T:A,T:T,etc. in the 16x16 matrix.

The mksum program provided in the distribution constructs new substitution matrices from any alignment (see Tools section).

In the absence of pseudocounts, the lod-scores of bases or basepairs that are never observed at a given position are set at the default arbitrary value of -20. The "-logzero" parameter can change this. A higher logzero (e.g. -5) will result in a greater

Erpin precalculates alignments of single stranded elements over a sliding window of 1024 nt. This can be changed using the -

tablen argument. Raising this number could produce some CPU gains when motifs longer than 100 nt are searched.

Searches are performed on both strands by default. Use -fwd or -rev to limit search on plus or minus strand, respectively.

- bgn n : search starts at position n in each sequence - len n : n nucleotides are searched after start point **Histogram (-hist)** 

Creates histogram of scores obtained in the current search. Generates a file called epnhist.dat, the first line of which contains:

The -hist option is best used in conjunction with the epnstat program (see Extra Tools), to evaluate score distributions in

This Perl script reads a parenthesized alignment (fasta-like format) and translates it into .epn (ERPIN2) format. The program also checks that helices have no gap (otherwise change positions to single-stranded) and deletes columns with gaps only.

GGUGACGAUAGCGAGAAGGUCACACCCGUUCCCGAACACGGAAGUGAGCUUCUCAGCGCCGACGGUAGAGAGUAGGACGUUGCC

GUGGUUAA-AGAAAAGAGGAAACACCUGUUAUCGAACACAGAAGUUAGCUCUUAUUCGCUGAUGGUAGAGAGUAGG-UUAUUGC

GUUGCCAU-GGUGGAGAGGUCAUACCCGUUCCCGAACACGGAAGUCAGCUCUCCUACGCCGAUGAUAGAAAGUAGG-UAGUAGC

This Perl script reads, filters and prints out Erpin outputs. It removes all solutions with 3 undefined nucleotides or more, or with

./readerpin.pl <erpin output file> [-fasta] [-c <cutoff>] [-e <cutoff>]

-c <cutoff> : prints only sequences with score higher than <cutoff>. -e <cutoff> : prints only sequences with E-value lower than <cutoff>.

This programs evaluates the score of any region or mask in the training set sequences (.epn file).

tstat <trset> <region> [(-mask | -umask) <arg1> <arg2> ..][..]

"Mask" and "region" arguments are the same as for Erpin. For example:

compatibility. default pcw value is still 10e-1.

tstat trna.epn -7,7 -umask 7 -umask 8

: prints sequences in fasta format

>b. subtillis CGGCGGCCAUAGCGGCAGGGAAACGCCCGGAACCCGGAAGCUAGCCUGCCAGCGCCGAUGGUGGAGAGUAGGUCACCGC >h. sapiens >s. cerevisiae

a score lower than a specified threshold.

conjunction with the -cutoff argument in Erpin.

tstat trna.epn -7,7 -mask 8

Typical input (parenthesized) alignment:

4. VERSION HISTORY 4.2.5 A significance performance improvement is achieved through optimization of element score assignment. Overall program speed is doubled! A error in pseudocount calculation is corrected. Does not affect results significantly, but default pcw

3.2 Multiple alerts added for excessive CPU and memory usage 3.1.2 secured dynamic mask handler 3.1 code cleaned 3.0.3 fixed bug with very large alignments (16S) 3.0.2 fixed bug

code cleaned. fixed bug with very large alignments (16S). aug 2002 Masking must be specified: does not anymore default to -nomask Lowercase replace Asterisks when displaying masked regions. jun 2002 Fixed bug in precomputed alignment tables (affected multi-level searches) Fixed bug in boundary transmission during multi-level search Fixed bug in search boundaries

<u>1.2</u> <u>1.1</u> 1.0 JMB paper version

tview This programs pretty prints training set files (.epn), displaying all sequences and the list of structural elements. Syntax: tview <trset> [<region>] Region: optionally restricts display to a specified region. Region is expressed as in erpin (e.g. -2,2). tstat

For each region or mask in the command line, tstat prints the best/worst/mean score in the alignment, as well as the score cutoffs that should be used to retain 100%, 90%, etc. of the training set sequences. These values can then be used in

Tstat also computes the number of configurations and memory requirement for a given region (useful when dealing with large

values should now be around 2x10e-4 instead of 10e-1. A corrective factor is applied for backward

Improved E-value (convolution product), Henikoff-like pseudo-counts (-pcw parameter), discontinued 4.2.2 Markov representation of gapless single-strands 4.0 stable version w/ E-value 3.9.9.last fixed bugs 3.9.8.bis fixed bugs 3.9.5 fixed bugs

3.0.1 Revised output and parameter formats as in Erpin 2.7 3.0 Introduction of dynamic mask handling. -add parameter introduced 2.9

2.7 2.5 2.3 2.1.1

2.1 Fixed bug in ungapped single strand scores Handling of complex patterns and large alignments Changes in command line and training set format 2.0 First documented version nov 5 2001. a few bugs fixed Introduction of "Helix" element