

---

# COMPARI → ←← MOTIF

## CompariMotif: Motif-Motif comparison software

---

© Richard J. Edwards (2007)

---

<b>1 Introduction .....</b>	<b>2</b>
1.1 Version.....	2
1.2 Using this Manual.....	2
1.3 Why use CompariMotif? .....	2
1.4 Getting Help.....	2
1.4.1 Something Missing?.....	3
1.5 Citing CompariMotif.....	3
1.6 Availability and Local Installation.....	3
<b>2 Fundamentals .....</b>	<b>4</b>
2.1 Running CompariMotif.....	4
2.1.1 The Basics.....	4
2.1.2 Options.....	4
2.1.3 Running in Windows.....	4
2.2 Input.....	4
2.2.1 Motif Input Formats.....	4
2.2.2 Advanced Input Options.....	5
2.2.3 SLiM Database Files .....	5
2.2.4 Use of DNA Motifs .....	5
2.2.5 Amino acid frequencies.....	6
2.3 Output .....	7
2.3.1 Output Styles .....	7
2.3.2 Cytoscape XGMML File.....	7
2.3.3 Optional Motif Information File.....	8
2.3.4 Sortable table output (webserver only).....	8
2.4 Commandline Options .....	9
2.5 The CompariMotif Webserver.....	10
<b>3 Motif Comparisons .....</b>	<b>11</b>
3.1 How CompariMotif Works .....	12
3.1.1 Single Position Comparisons.....	12
3.1.2 Selecting Pairwise Matches.....	12
3.1.3 Defining Motif Relationships.....	13
3.2 Information Content .....	13
3.3 Score.....	15
3.4 Example Application .....	15
<b>4 Appendices.....</b>	<b>17</b>
4.1 Troubleshooting & FAQ .....	17
4.2 SLiM Definitions .....	17
4.3 References.....	18

## FIGURES

---

Figure 2.1. Partial Cytoscape visualisation of CompariMotif relationships.....	8
Figure 2.2. CompariMotif webserver. ....	10
Figure 3.1. Overview of CompariMotif .....	11
Figure 3.2. CompariMotif Match Type Examples.....	14
Figure 4.1. CompariMotif webserver output for 14-3-3 HPRD SLiMFinder analysis vs. ELM.....	15
Figure 4.2. CompariMotif XGMML output visualized with Cytoscape(Shannon <i>et al.</i> 2003) (recoloured).....	16
Figure 5.1. Anatomy of a SLiM. ....	17

## TABLES

---

Table 2.1. SLiM Databases provided for CompariMotif searches. ....	6
Table 2.2. Field headings for main CompariMotif output file.....	7
Table 2.3. CompariMotif Commandline Options. ....	9

# 1 Introduction

---

This manual gives an overview of SLiMfinder as implemented in the `comparimotif_V3.py` module. Because there are many options, this manual will probably not be fully comprehensive but aims to cover the basics and the most useful of the more advanced stuff. If anything is missing or needs clarification, please contact me. The fundamentals are covered in [Chapter 2, Fundamentals](#), including input and output details. Later sections give more details on how the methods work and statistics are generated. General details about Command-line options can be found in the [PEAT Appendices](#) document included with this download. Details of command-line options specific to Slim Pickings can be found in the distributed [readme.txt](#) and [readme.html](#) files

Like the software itself, this manual is a ‘work in progress’ to some degree. If the version you are now reading does not make sense, then it may be worth checking the website to see if a more recent version is available, as indicated by the [Version](#) section of the manual. Furthermore, many options have been added to Slim Pickings over the past few weeks and not all of them have found their way into the manual yet. Check the [readme](#) on the website for up-to-date options etc. In particular, default values for options are subject to change and should be checked in the [readme](#).

Good luck.

Rich Edwards, 2007.

## 1.1 Version

---

This manual was written to accompany [CompariMotif version 3.4](#). The manual was last edited on 17 March 2008.

## 1.2 Using this Manual

---

As much as possible, I shall try to make a clear distinction between explanatory text (this) and text to be typed at the command-prompt etc. Command prompt text will be written in Courier New to make the distinction clearer. Program options, also called ‘command-line parameters’, will be **written in bold Courier New** (and coloured red for fixed portions or dark red for user-defined portions, such as file names etc.). Command-line examples will be given in (purple) *italicised Courier New*. Optional parameters will (if I remember) be [in square brackets]. Names of files will be marked in normal text by (blue-grey) Times New Roman.

## 1.3 Why use CompariMotif?

CompariMotif takes two lists of protein motifs and compares them to each other, identifying which motifs have some degree of overlap and describing the relationships between those motifs. It can be used to compare a list of motifs with themselves, their reversed selves, or to a second list of motifs. CompariMotif outputs a table of all pairs of matching motifs, along with their degree of similarity (information content) and their relationship to each other. Details can be found in 2.3.

Short linear motifs (SLiMs) in proteins are functional microdomains of fundamental importance in many biological systems (Neduva & Russell 2005). SLiMs typically consist of a 3 to 10 amino acid stretch of the primary protein sequence, of which as few as two sites may be important for activity. SLiM can usually tolerate a number of alternative amino acids at one or more positions, making precise definitions extremely difficult. CompariMotif can therefore be extremely useful when a new SLiM has been discovered, either by high throughput SLiM discovery (Neduva *et al.* 2005; Davey *et al.* 2006; Neduva & Russell 2006; Edwards *et al.* 2007) or by low throughput experimental studies, by allowing similar motifs to be readily identified from published resources (e.g. ELM (Puntervoll *et al.* 2003) or MiniMotif (Balla *et al.* 2006)). Alternatively, a list of motifs could be compared to itself to identify recurring motifs.

## 1.4 Getting Help

---

Much of the information here is also contained in the documentation of the Python modules themselves. A full list of command-line parameters can be printed to screen using the **help** option, with short descriptions for each one:

```
python comparimotif_V3.py help
```

General details about Command-line options can be found in the [PEAT Appendices](#) document included with this download. Details of command-line options specific to Slim Pickings can be found in the distributed [readme.txt](#) and [readme.html](#) files.

If still stuck, then please e-mail me ([r.edwards@southampton.ac.uk](mailto:r.edwards@southampton.ac.uk)) whatever question you have. If it is the results of an error message, then please send me that and/or the log file (see [2.3](#)) too.

### 1.4.1 Something Missing?

As much as possible, the important parts of CompariMotif are described in detail in this manual. If something is not covered, it is generally not very important and/or still under development, and can therefore be safely ignored. If, however, curiosity gets the better of you, and/or you think that something important is missing (or badly explained), please contact me.

## 1.5 Citing CompariMotif

Until published in its own right, please cite the [SLiMDisc Webserver](#) paper ([Davey et al. 2007](#)). For analyses on the webserver using specific motif databases, please cite the listed papers for those motif databases.

## 1.6 Availability and Local Installation

CompariMotif can be run from the CompariMotif webserver, available at <http://bioware.ucd.ie/>.

CompariMotif is also distributed as a number of open source Python modules as part of the PEAT (Protein Evolution Analysis Toolkit) package. It should therefore work on any system with Python installed without any extra setup required. If you do not have Python, you can download it free from [www.python.org](http://www.python.org) at <http://www.python.org/download/>. The modules are written in Python 2.5. The Python website has good information about how to download and install Python but if you have any problems, please get in touch and I will help if I can.

All the required files should have been provided in the download zip file. Details can be found at <http://bioinformatics.ucd.ie/shields/software/peat/> and the accompanying **PEAT Appendices** document. The Python Modules are open source and may be changed if desired, although please give me credit for any useful bits you pillage. I cannot accept any responsibility if you make changes and the program stops working, however!

Note that the organisation of the modules and the complexity of some of the classes is due to the fact that most of them are designed to be used in a number of different tools. As a result, not all the options listed in the `__doc__()` (**help**) will be of relevance. If you want some help understanding the way the modules and classes are set up so you can edit them, just contact me.

## 2 Fundamentals

---

### 2.1 Running CompariMotif

#### 2.1.1 The Basics

If you have python installed on your system, you should be able to run CompariMotif directly from the command line in the form:

```
python comparimotif_V3.py motifs=FILENAME
```

For the example provided in the distribution:

```
python comparimotif_V3.py motifs=comparimotif_eg.motifs
```

If the motif file is to be compared to itself, then no other commands are needed. If a second file is to be compared, however, this should be specified using the **searchdb=FILENAME** command:

```
python comparimotif_V3.py motifs=file1 searchdb=file2
```

**IMPORTANT:** If filenames contain spaces, they should be enclosed in double quotes: **motifs="example file"**. That said, it is recommended that files do not contain spaces as function cannot be guaranteed if they do.

A CompariMotif webserver is also available at <http://bioware.ucd.ie>. See **Error! Reference source not found.** for details.

#### 2.1.2 Options

Command-line options are suggested in the following sections. General details about Command-line options can be found in the [PEAT Appendices](#) document included with this download. Details of command-line options specific to Slim Pickings can be found in the distributed [readme.txt](#) and [readme.html](#) files. These may be given after the run command, as above, or loaded from one or more \*.ini files (see [PEAT Appendices](#) for details).

#### 2.1.3 Running in Windows

If running in Windows, you can just double-click the `comparimotif_V3.py` file and use the menu prompts to navigate through the program. It is recommended to use the **win32=T** option. (Place this command in a file called `comparimotif.ini`.)

## 2.2 Input

The main input for CompariMotif is a motif file and an optional second motif file to compare these motifs to. If no second motif file is given (**searchdb=FILE**) then the first motif file (**motifs=FILE**) will be compared to itself.

#### 2.2.1 Motif Input Formats

The recommended motif input format is PRESTO format. This should have a single line per motif, with the format:

```
Name Sequence # Comments
```

Comments are optional but anything after the # will be ignored.

Alternative allowed formats include: a fasta format file with motif/peptide names and sequences in the usual fasta format; a raw list of peptides/motifs (in this case the name and sequence will be the same); SLiMDisc output; TEIRESIAS output; Slim Pickings output. Additional input formats can be added on request

In either case, the motif should be a peptide sequence using the standard single letter amino acid codes and the following regular expression rules:

- **A** = single fixed amino acid.
- **[AB]** = ambiguity, **A** or **B**. Any number of options may be given, e.g. **[ABC]** = **A** or **B** or **C**. **[^A]** = not **A**.
- **<R:m:n>** = At least **m** of a stretch of **n** residues must match **R**, where **R** is one of the above regular expression elements (single or ambiguity).
- **x** or **.** = Wildcard positions (any amino acid)
- **X{m,n}** or **.{m,n}** = At least **m** and up to **n** wildcards.
- **R{n}** = **n** repetitions of **R**, where **R** is any of the above regular expression elements.
- **(AB|CD)** = **AB** or **CD**.
- **(ABC)** = **ABC** in any order (**BAC**, **CAB** etc.).
- **^** = Beginning of sequence
- **\$** = End of sequence

*E.g.* **[IL][^P]X{3}RG** means: “leucine or isoleucine, followed by anything but proline, followed by three residues, followed by arginine followed by glycine”.

*E.g.* (2) **^<KR:3:5>P** means: “three of the first five amino acids must be lysine or arginine; the sixth amino acid must be proline”.

## 2.2.2 Advanced Input Options

If **reverse=T** is used, the first file only will be reversed before comparison. If a self-comparison is made using **reverse=T** then reversed motifs will be compared to the original “forward” motif file. Input motifs can be filtered to remove short or highly degenerate motifs (see 0 for details).

## 2.2.3 SLiM Database Files

CompariMotif provides a number of pre-defined motif datasets in the correct format for use. When using results from these datasets, please always cite the relevant paper. These datasets are given in Table 2.1. From the CompariMotif webserver, these databases are available via a drop-down list box.

If you have additional motif databases you would like to see available, please contact me. In addition, individual motifs may be submitted for inclusion in the miscellaneous literature motif file at: <http://bioware.ucd.ie/~comparimotif/MotifBrowser/index.html>.

## 2.2.4 Use of DNA Motifs

Although explicitly designed for protein motifs, there is no reason why DNA motifs cannot be compared with CompariMotif, as long as the correct regular expression notation is maintained. CompariMotif version 3.3 introduced a DNA motif option, **dna=T/F**. When **dna=T**, only the four DNA nucleotides G, A, T and C are used and information content (see 3.2) is adjusted accordingly. Uridines (U) are converted to thymidines (T). The following additional replacements (based on official IUB/IUPAC abbreviations) are made:

- N --> . (any)
- R --> G A (purine)
- Y --> T C (pyrimidine)
- K --> G T (keto)
- M --> A C (amino)
- S --> G C (strong)
- W --> A T (weak)
- B --> G T C
- D --> G A T

- H --> A C T
- V --> G C A

## 2.2.5 Amino acid frequencies

By default, all information content calculations (see 3.2) assume uniform amino acid frequencies (or nucleotide frequencies if **dna=T**.) This is because the motifs themselves are independent of amino acid bias. The chance of any given motif having a match to a database is more a reflection of the frequency of amino acids in motifs, rather than proteomes. (Indeed, rare amino acids might be more likely to occur in motifs because they are rare and therefore look unusual to discovery algorithms.) Nevertheless, it is sometimes desirable to weight results according to amino acid frequencies, and this can be done using the **aafreq=FILE** command. The **FILE** can be a FASTA sequence file from which frequencies are to be calculated, or a delimited file containing amino acid frequencies:

```
AA      FREQ
A       0.074
C       0.033
...
Y       0.033
```

The impact of weighted frequencies is discussed in 3.2.

**Table 2.1. SLiM Databases provided for CompariMotif searches.**

Database	Description	Reference	Motif File
ELM	The Eukaryotic Linear Motif database provides a number of high quality annotated SLiMs with known occurrences. <b>Note.</b> Some motifs have been split into <b>_a</b> and <b>_b</b> to be compatible with CompariMotif input formats. Such motifs are marked <b>*Modified*</b> in their descriptions.	(Puntervoll <i>et al.</i> 2003)	<a href="#">ELM.motifs</a>
MiniMotif	Another database of SLiMs from all organisms. This has less annotation than ELM but more motifs. These motifs have been reformatted to conform to standard regular expressions.	(Balla <i>et al.</i> 2006)	<a href="#">MnM.motifs</a>
Phospho-MotifFinder	Motifs from the PhosphoMotif Finder database of HPRD. Motifs are labelled KIN for Kinase / Phosphatase motifs or BIND for binding motifs. <b>_Y</b> indicates a tyrosine motif, while <b>_ST</b> indicated serine/threonine. The number part of the motif identifier is arbitrary and has no link to the website. <b>Note.</b> All these motifs are phosphorylation motifs and, as such, have a key Ser, Thr or Tyr position. These are <u>not</u> given special treatment in CompariMotif and the user should pay special attention to whether the appropriate residue is included in the match.	(Amanchy <i>et al.</i> 2007)	<a href="#">phosphomotif.motifs</a>
Misc. Literature	Miscellaneous motifs collected from the literature. (These include pubmed links to the relevant paper but we cannot guarantee the accuracy of the motifs or their descriptions.)	See file.	<a href="#">literature.motifs</a>
Combined Literature	A combined database of the above sources. The source database is indicated in the motif description: [ELM] for ELM, [MnM] for MiniMotif and [PMF] for PhosphoMotif Finder.	See above.	<a href="#">combined.motifs</a>
Neduva & Russell	Predicted interaction SLiMs from the high-throughput study of Neduva <i>et al.</i> (2005). The motif name indicates what part of the study it is from. All names begin NR, followed by a two-letter code for the species and a one-letter code denoting <b>_Domain-level datasets</b> or <b>Protein-level datasets</b> . (See paper for details.) Species codes are: Ce, <i>C. elegans</i> ; Dm, <i>D. melanogaster</i> ; Hs, <i>H. sapiens</i> ; Sc, <i>S. cerevisiae</i> .	(Neduva <i>et al.</i> 2005)	<a href="#">Ned2005_Sig.motifs</a>

## 2.3 Output

The main output for CompariMotif is delimited text file containing the following fields:

**Table 2.2. Field headings for main CompariMotif output file.**

Field	Description
File1	Name of motif file 1 ( <b>motifs=FILE</b> ). [ <b>outstyle=multi</b> only]
File2	Name of file 2 ( <b>searchdb=FILE</b> ). [ <b>outstyle=multi</b> only]
Name1	Name of motif from motif file 1.
Name2	Name of motif from motif file 2.
Motif1	Motif (pattern) from motif file 1.
Motif2	Motif (pattern) from motif file 2.
Sim1	Description of motif1's relationship to motif2.
Sim2	Description of motif2's relationship to motif1.
Match	Regular expression of match between motifs. Upper case positions indicate an exact match, while lower case positions have some degree of degeneracy difference between the two motifs. Mismatches are marked with an asterisk.
MatchPos	Number of matched positions between motif1 and motif2 ( <b>&gt;=mishare=X</b> ).
MatchIC	Information content of matched positions.
NormIC	MatchIC as a proportion of the maximum possible MatchIC. If this is 1.0 then the match is a good as could possibly be achieved.
Score	Heuristic score (MatchPos x NormIC) for ranking motif matches.
Info1	Information Content of motif1 (if <b>motific=T</b> ).
Info2	Information Content of motif2 (if <b>motific=T</b> ).
Desc1	Description of motif1 (if <b>motdesc=1</b> or <b>motdesc=3</b> ).
Desc2	Description of motif1 (if <b>motdesc=1</b> or <b>motdesc=2</b> ).

Details of relationship descriptions and information content calculations can be found in [Chapter 3](#).

Note that Info1 and Info2 are only output if the **motific=T** option is used. Desc1 and Desc2 are controlled by **motdesc=X**. This outputs motif descriptions depending on the value of X, as follows: 0 = Neither; 1 = Motif1 only; 2 = Motif2 only; 3 = both. The default is 3 (both).

### 2.3.1 Output Styles

With the exception of the file names, which are only output if **outstyle=multi**, the above is the output for the default "normal" output style. If **outstyle=single** then only statistics for motif2 (the **searchdb** motif) are output as this is designed for searches using a single motif against a motif database. If **outstyle=normalsplit** or **outstyle=multisplit** then motif1 information is grouped together, followed by motif2 information, followed by the match statistics, e.g. **[File1,] Name1, Motif1, Sim1, [Info2,] [Desc1,] [File2,] Name2, Motif2, Sim2, [Info2,] [Desc2,] Match, MatchPos, MatchIC, NormIC, Score**

### 2.3.2 Cytoscape XGMML File

CompariMotif also outputs an XGMML file (\*.compare.xgmml) that can be imported directly into Cytoscape (Shannon *et al.* 2003) (<http://www.cytoscape.org>) for visualisation of the results ([Figure 2.1](#)). This file contains all the necessary node (motif) and edge (match) data found in the results table, which can be viewed for selected nodes/edges using the Cytoscape Data Panel. The file can be uploaded into Cytoscape using the **File -> Import -> Network (Multiple File Types)** command (**CTRL+L**). When first loaded, nodes will be displayed in a simple, uninformative, grid. Use one of the Cytoscape Layouts (e.g. **Layout -> yFiles -> Organic**) to make it clearer. See Cytoscape documentation for details.







## 2.4 Commandline Options

Table 2.3 lists the commandline options for CompariMotif. Please see also the [PEAT Appendices](#) document for additional general commandline options and the [RJE\\_SEQ Manual](#) for further input data options. Beginners will probably want to leave the default settings unchanged.

**Table 2.3. CompariMotif Commandline Options.**

Option	Description	Default
<b>Basic Input/Output Options</b>		
<b>motifs=FILE</b>	Loads motifs from FILE	[None]
<b>searchdb=FILE</b>	Optional second motif file to compare.	[None]
<b>resfile=FILE</b>	Name of results file, FILE.presto.txt.	[motifsFILE-searchdbFILE.presto.txt]
<b>motinfo=FILE</b>	Filename for output of motif summary table	[None]
<b>motific=T/F</b>	Output Information Content for motifs.	[False]
<b>Motif Comparison Parameters</b>		
<b>minshare=X</b>	Min. number of non-wildcard positions for motifs to share.	[2]
<b>normcut=X</b>	The minimum normalise IC allowed for a match.	[0.5]
<b>matchfix=T/F</b>	If >0 must exactly match *all* fixed positions in the motifs from: - 1: input (motifs=FILE) motifs - 2: searchdb motifs - 3: *both* input and searchdb motifs	[0]
<b>Advanced Input Parameters</b>		
<b>minic=X</b>	Min information content for motif	[2]
<b>minfix=X</b>	Min number of fixed positions for a motif to contain	[0]
<b>minpep=X</b>	Min no. of defined positions	[2]
<b>trimx=T/F</b>	Whether to trim leading and trailing wildcards	[False]
<b>nrmotif=T/F</b>	Remove redundancy in input motifs (partial/full identities)	[False]
<b>reverse=T/F</b>	Reverse the first set of motifs.	[False]
<b>mismatches=X</b>	<= X mismatches of positions can be tolerated.	[0]
<b>dna=T/F</b>	Whether motifs should be considered as DNA motifs	[False]
<b>aafreq=FILE</b>	Use FILE to replace uniform AAFreqs (FILE can be sequences or aafreq)	[None]



## 3 Motif Comparisons

This chapter gives more details on the inner workings of CompariMotif.

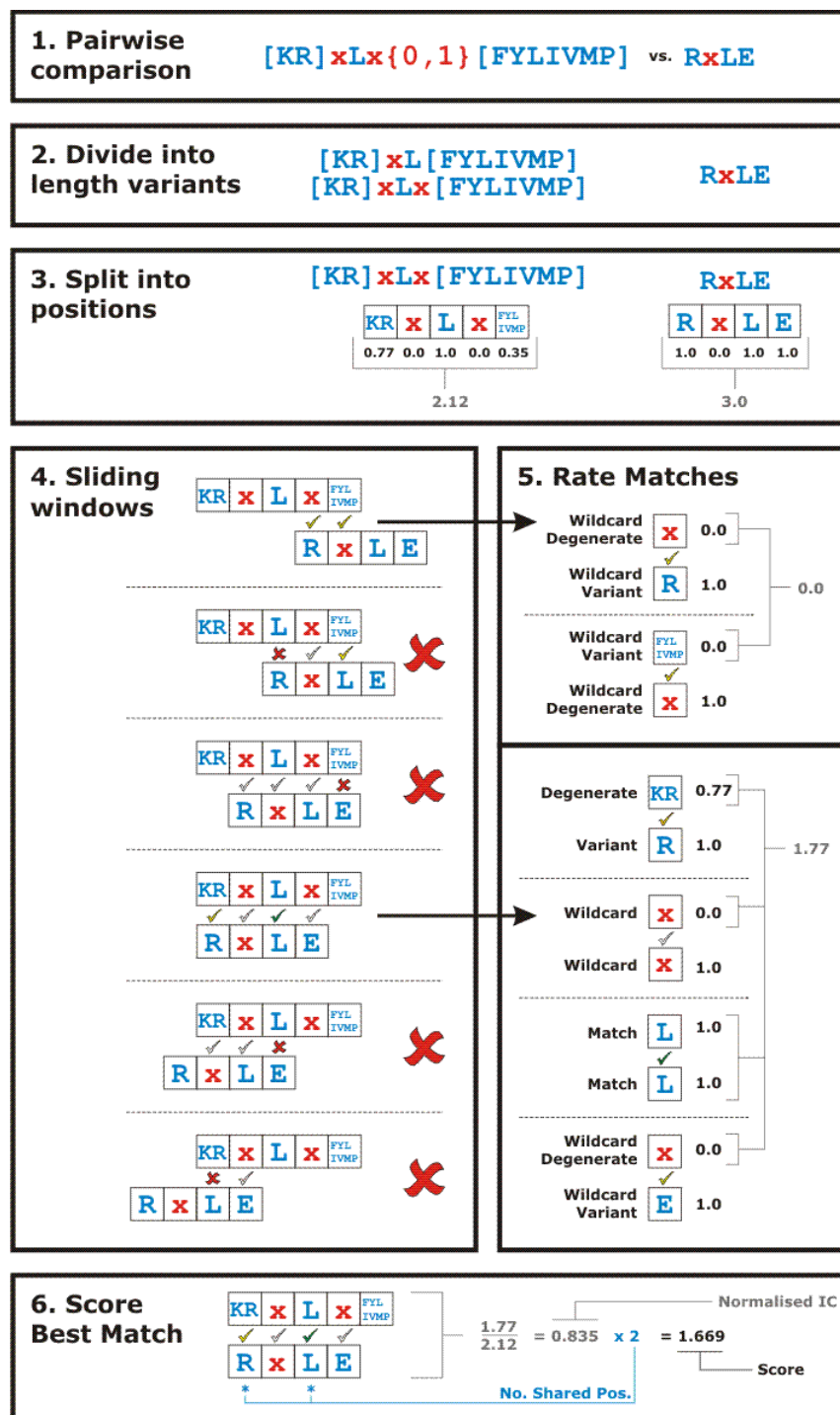


Figure 3.1. Overview of CompariMotif.  
Details can be found in the text.

### 3.1 How CompariMotif Works

An overview of CompariMotif is given in [Figure 3.1](#). Motifs are first compared for precise matches. If these are not found, then CompariMotif adopts a sliding window comparison in which every possible overlap between (variants of) the two motifs are compared against each other. Matches must meet a minimum match requirement determined by the **minshare=X**, **normcut=X** and **matchfix=X** options (see [3.1.2](#)). Fixed positions in motifs are often more important than ambiguous ones, especially when the motif has been experimentally determined. For this reason, it is also possible to stipulate that all fixed positions in one or other motif (or both) match exactly to fixed positions in the compared motifs. This is controlled using the **matchfix=X** option.

#### 3.1.1 Single Position Comparisons

For every comparison, each position in each motif is then rated according to its relationship with the compared position in the other motif:

- Match = perfect fixed position match
- Wildcard = wildcard in both motifs
- Wildcard variant = wildcard in compared motif but not in focal motif
- Wildcard degenerate = wildcard in focal motif but not in compared motif
- Ambiguous Match = ambiguities in both motifs comparing the same amino acids
- Degenerate Ambiguity = ambiguity in both motifs but the compared site is a subset of the ambiguity in the focal site
- Variant Ambiguity = ambiguity in both motifs but the focal site is a subset of the ambiguity in the compared site
- Degenerate = ambiguity in focal motif but fixed position in compared motif
- Variant = a fixed variant in the focal motif of a degenerate position in the compared motif
- Overlapping ambiguity = ambiguity in both motifs where 1+ amino acids overlap but each ambiguity also contains amino acids not in the other
- Bad ambiguity = ambiguity in focal motifs sharing no amino acids in common with compared motif
- Ambiguity mismatch = fixed position in focal motif that does not match ambiguity in compared motif
- Mismatch = different fixed positions in each motifs

Each positional comparison is then given an information content (IC) rating, if it is a “good” match. This is the lower IC out of the two positions compared. *E.g.* a fixed variant matching an ambiguity will take the IC of the ambiguity.

#### 3.1.2 Selecting Pairwise Matches

The entire pairwise comparison is then rated for:

- Number of matching positions, allowing for degeneracy
- Number of exactly matching fixed positions
- Match Information content (IC), which is the sum of IC over all matched positions
- Number of incompatible positions (*e.g.* Bad ambiguities and mismatches)

The comparison is then rejected as a potential match if one of the following conditions is met:

- There are any incompatible positions. (If the **mismatches=X** option is used, this is relaxed.)
- The number of matched positions is less than that stipulated by **minshare=X**.

- The **matchfix=X** option is used and the relevant motif(s) in the comparison does not have exact matches at all its fixed positions.
- The normalised IC is below that set by **normcut=X**.

When a motif has multiple length variants and/or “NofM” elements, each possible variant is compared.

Multiple variants and/or sliding windows can produce multiple matches that meet the acceptance criteria. In this case, the match with the best information content is used. In the case of tied information content, matches are assessed by the number of matching positions and then the number of exactly matching fixed positions. The earlier comparison made is considered “best” if all these stats tie.

### 3.1.3 Defining Motif Relationships

The best match is then considered to define the relationship between the two motifs. These relationships are comprised of the following keywords:

- Match type keywords identify the type of relationship seen:
  - **Exact** = all the matches in the two motifs are precise
  - **Variant** = the focal motif contains only exact matches and subvariants of degenerate positions compared to the other motif
  - **Degenerate** = the focal motif contains only exact matches and degenerate versions of positions in the other motif
  - **Complex** = some positions in the focal motif are degenerate versions of positions in the compared motif, while others are subvariants of degenerate positions
- Match length keywords identify the length relationships of the two motifs:
  - **Match** = both motifs are the same length and match across their entire length
  - **Parent** = the focal motif is longer and entirely contains the compared motif
  - **Subsequence** = the focal motif is shorter and entirely contained within the compared motif
  - **Overlap** = neither motif is entirely contained within the other

This gives sixteen possible classifications for each motif’s relationship to the compared motif (Figure 3.2).

## 3.2 Information Content

Information content (IC) is calculated for each position based on a modification of Shannon’s classical Information Content (Shannon 1997):

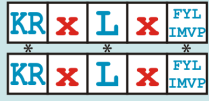
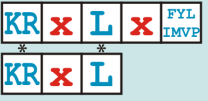
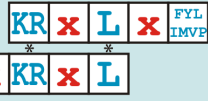
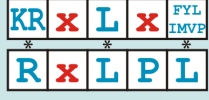


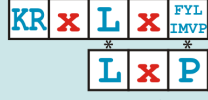



$$IC_i = -\log_N(f_a)$$

where  $IC_i$  is the information content for position  $i$ ,  $f_a$  is the summed frequency for the amino acids (or nucleotides) at position  $i$  and  $N$  is number of amino acids (or nucleotides) in the alphabet, *i.e.*  $N=4$  for DNA and  $N=20$  for proteins. The IC for the motif is simply this score summed over all positions.

This score is essentially a rescaling of Shannon’s self-information such that a wildcard receives 0.0 and a fixed position scores 1.0 when a uniform frequency distribution is used. Ambiguous positions are given a value between 0.0 and 1.0. When non-uniform frequencies are used (see 2.2.5), fixed rare amino acids ( $f_a < 1/N$ ) will score above 1.0, while fixed common amino acids ( $f_a > 1/N$ ) will score less than 1.0.

For motif matches, the *lowest* information content for any compared pair of elements is used as the information content for that part of the match. As a result, “matches” at wildcard position contribute nothing to the IC for the match, while for fixed sites matching ambiguous positions, the IC for the degenerate (ambiguous) position is used. The match IC,  $IC_m$ , is the sum of the lower  $IC_i$  values for each position. The “best” IC for any given pair of motifs is therefore equal to the lower IC of the two motifs. For the “Normalised IC” output, the IC of the match is divided by this number. This normalises

the  $IC_m$  against the fact that longer and less degenerate motifs will tend to produce higher IC matches. Matches with a normalised IC below that set by **normcut=X** will not be returned.

	Match	Parent/ Subsequence	Overlap
Exact	<b>Exact Match</b>  Exact Match	<b>Exact Parent</b>  Exact Subsequence	<b>Exact Overlap</b>  Exact Overlap
Variant/ Degenerate	<b>Degenerate Match</b>  Variant Match	<b>Degenerate Parent</b>  Variant Subsequence <b>Variant Parent</b>  Degenerate Subsequence	<b>Degenerate Overlap</b>  Variant Overlap
Complex	<b>Complex Match</b>  Complex Match	<b>Complex Parent</b>  Complex Subsequence	<b>Complex Overlap</b>  Complex Overlap

**Figure 3.2. CompariMotif Match Type Examples.**

Examples for each of the sixteen match types. In each case, the “query” motif [KR]xL[FYLIMVP] is compared to an invented motif for illustration. Because of the natural relationship between parent/subsequence and variant/degenerate matches, these have been grouped in the figure. Matched positions that contribute towards the number of matched positions (i.e. those not involving a wildcard position) are marked with an asterisk. **Exact Match**: All positions are identical and the match spans the full length of both motifs; **Variant/Degenerate Match**: The match spans the full length of both motifs. All of the positions of the query are either the same as the match (X v. X and L v. L) or more degenerate ([KR] v. R, X v. P & [FYLIMVP] v. L) and so it is classed degenerate. Likewise, all positions in the other motif, RxLPL, are either identical to the query or variants of the query positions, so it is classed as variant; **Complex Match**: Again, the match spans the full length of both motifs. This time, each motif has some positions that are more degenerate than in the other motif. *i.e.* The query is a variant for the L v. X position but more degenerate at all other positions; **Exact Parent/Subsequence**: The [KR]xL motif is entirely and exactly contained within the query; **Degenerate Parent/Variant Subsequence**: The RxLE motif is entirely contained within the query. At two positions, however, ([KR] v. R & X v. E) the query is more degenerate and is hence a “degenerate parent”, while RxLE is a “variant subsequence”; **Variant Parent/Degenerate Subsequence**: This time it is the query that is the variant in one position (L v. IL) and so the variant/degenerate labels are swapped; **Complex parent/subsequence**: the L[IL]xL motif is less degenerate at two positions (X v. L & [FYLIMVP] v. L) but more degenerate at one (L v. [IL]) and so the relationship is “complex”; **Exact overlap**: Neither motif is entirely contained within the other but the positions overlapping match exactly; **Degenerate/variant overlap**: Neither motif is entirely contained within the other. The first P of LxPP is a variant of the [FYLIMVP] in the query, while the other two matches (an L and an X) are exact, therefore the query is “degenerate” and LxPP is “variant”; **Complex overlap**: Neither motif is entirely contained within the other and both contain positions that are degenerate when compared to the matching position in the other motif ([KR] v. R & X v. S are degenerate in the query, L v. [ILMV] is degenerate in RxRS[ILMV]).

### 3.3 Score

The Score assigned to a match is a simple heuristic of the match IC multiplied by the normalised IC. This is a useful metric for ranking matches, as the best matches tend to get the best scores.

### 3.4 Example Application

A typical application for CompariMotif is given in the SLiMFinder paper ((Edwards *et al.* 2007), Example 1), in which HPRD interaction datasets for 14-3-3 proteins (Mishra *et al.* 2006) were analysed using SLiMFinder, returning several significant motifs ( $p < 0.05$ , see Table 2 in (Edwards *et al.* 2007)). These motifs in the recommended CompariMotif input format would be as follows:

```
YWHAE_1 R..S.P..L # Sig. SLiMFinder motif for HPRD 14-3-3 Epsilon interactors
YWHAH_1 GR.[ST]..P # Sig. SLiMFinder motif for HPRD 14-3-3 Eta interactors
YHWAG_1 ^.[AS].[AGS] # Sig. SLiMFinder motif for HPRD 14-3-3 Gamma interactors
YHWAG_2 KE..K # Sig. SLiMFinder motif for HPRD 14-3-3 Gamma interactors
YWHAQ_1 P..P..P # Sig. SLiMFinder motif for HPRD 14-3-3 Theta interactors
YWHAZ_1 [AGS]..P..P..P # Sig. SLiMFinder motif for HPRD 14-3-3 Zeta interactors
YWHAZ_2 ^.[AGS].[GS] # Sig. SLiMFinder motif for HPRD 14-3-3 Zeta interactors
YWHAZ_3 FR..[ST].S # Sig. SLiMFinder motif for HPRD 14-3-3 Zeta interactors
YWHAZ_4 [ST]P.[ST]P # Sig. SLiMFinder motif for HPRD 14-3-3 Zeta interactors
YWHAZ_5 Y.C.PG.L # Sig. SLiMFinder motif for HPRD 14-3-3 Zeta interactors
```

Raw SLiMFinder and SLiMDisc delimited text results can be loaded directly into CompariMotif for analysis without any reformatting.

These motifs were compared to the ELM database (Puntervoll *et al.* 2003) using CompariMotif with a “normalized IC” cut-off of 0.4. (Figure 3.3, Figure 3.4). Results were constrained such that fixed positions in an ELM must match a fixed position in the SLiMFinder motif. In total, eight out of ten SLiMFinder motifs had matches with seventeen ELMs. Matches fell into three main clusters: (1) Three motifs matching known 14-3-3 motifs, (2) Three motifs matching SH3 binding motifs, and (3) Two motifs matching the highly degenerate LIG\_PCNA\_1 motif (Figure 3.4). In addition to the 14-3-3 and SH3 ELMs, matches to five phosphorylation ELMs were also identified; phosphorylation of the 14-3-3 motif is important for ligand recognition. These comparisons took less than two seconds to run on an Intel(R) Xeon(TM) dual 3.20GHz processor with 3Gb RAM.

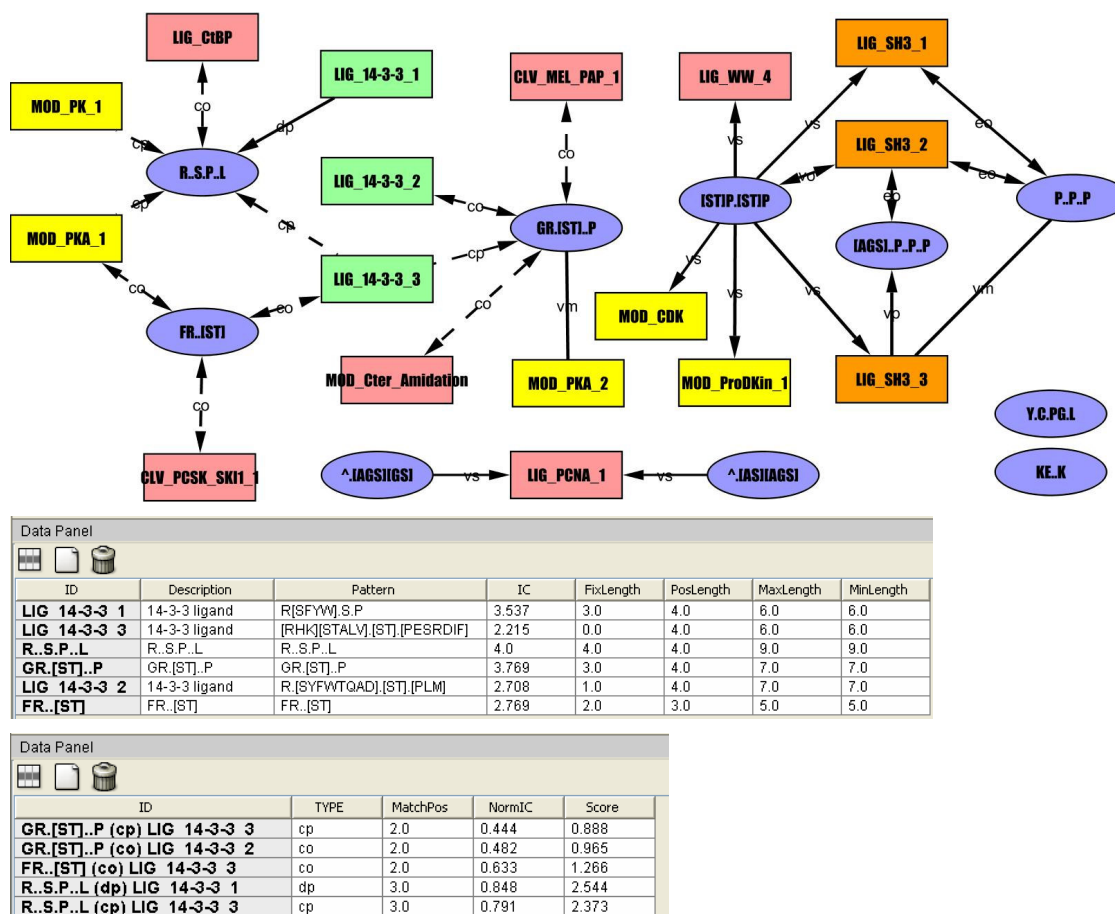
It is beyond the scope of this manual to discuss these results in detail. They do, however, highlight the ease with which CompariMotif can help to make sense of motif discovery results. As a simple, quick and high-throughput tool, CompariMotif can be an invaluable initial step in making sense of such data. Because of this, CompariMotif is now directly linked to both SLiMDisc and SLiMFinder web implementations (Davey *et al.* 2007).

#### Results

Name1	Name2	Motif1	Motif2	Sim1	Sim2	Match	Pos	MatchIC	NormIC	Score	Info1	Info2
YWHAE_1	LIG_14-3-3_1	R.S.P..L	R[SFYW].S.P	DP	VS	R[fswy].S.P	3	3.000	0.848	2.544	4.00	3.54
YWHAE_1	LIG_14-3-3_3	R.S.P..L	[RHK][STALV].[ST].[PESRDIF]	CP	CS	r[alstv].s.p	3	1.752	0.791	2.373	4.00	2.22
YWHAZ_3	LIG_14-3-3_3	FR.[ST].S	[RHK][STALV].[ST].[PESRDIF]	CP	CS	r[alstv].[ST].s	3	1.752	0.791	2.373	3.77	2.22
YWHAQ_1	LIG_SH3_1	P..P..P	[RKY].P..P	EO	EO	P..P	2	2.000	0.760	1.519	3.00	2.63
YWHAZ_4	LIG_SH3_1	[STP].[STP]	[RKY].P..P	VS	DP	[st]P.[st]P	2	2.000	0.760	1.519	3.54	2.63
YWHAQ_1	LIG_SH3_2	P..P..P	P..P.[KR]	EO	EO	P..P	2	2.000	0.722	1.445	3.00	2.77
YWHAZ_1	LIG_SH3_2	[AGS].P..P	P..P.[KR]	EO	EO	P..P	2	2.000	0.722	1.445	3.63	2.77
YWHAZ_4	LIG_SH3_2	[STP].[STP]	P..P.[KR]	VO	DO	P.[st]P	2	2.000	0.722	1.445	3.54	2.77
YWHAE_1	MOD_FK_1	R.S.P..L	[RK].[S][M]..	CP	CS	r..S[iv]p.	2	1.769	0.697	1.394	4.00	2.54
YWHAH_1	MOD_Cter_Amidation	GR.[ST].P	([G][RK][RK]	CO	CO	Gr[kr]	2	1.769	0.697	1.394	3.77	2.54
YWHAZ_4	MOD_CDK	[STP].[STP]	...([ST]P).[KR]	VS	DP	[st]p.[st]P	2	1.769	0.697	1.394	3.54	2.54
YWHAE_1	MOD_PKA_1	R.S.P..L	[RK][RK].[ST]..	CP	CS	r[kr].s.p.	2	1.537	0.667	1.333	4.00	2.31
YWHAZ_3	MOD_PKA_1	FR.[ST].S	[RK][RK].[ST]..	CO	CO	r[kr].[ST].s	2	1.537	0.667	1.333	3.77	2.31
YWHAE_1	LIG_CIBP	R.S.P..L	[PG][LVPM][DENS][LVASTRGE]	CO	CO	p[elimpv][dens]L	2	1.769	0.578	1.157	4.00	3.06
YWHAH_1	CLV_MEL_PAP_1	GR.[ST].P	[ILV].[R][VF][GS]	CO	CO	gR[v]s.	2	1.769	0.558	1.116	3.77	3.17
YWHAH_1	LIG_14-3-3_2	GR.[ST].P	R.[SYFWTGAD].[ST].[PLM]	CO	CO	R.[st].[st]p	2	1.306	0.482	0.965	3.77	2.71
YHWAG_1	LIG_PCNA_a	^[AS][AGS]	^[0,3].[FHWY][ILM][*P][*FHILWY][DHFM][FMY]..	VS	DP	^[as][ags]	2	1.074	0.447	0.895	2.40	3.07
YWHAZ_2	LIG_PCNA_a	^[AGS][GS]	^[0,3].[FHWY][ILM][*P][*FHILWY][DHFM][FMY]..	VS	DP	^[ags][gs]	2	1.074	0.447	0.895	2.40	3.07
YWHAH_1	LIG_14-3-3_3	GR.[ST].P	[RHK][STALV].[ST].[PESRDIF]	CP	CS	r[alstv][st][st]p	2	0.984	0.444	0.888	3.77	2.22

Figure 3.3. CompariMotif webserver output for 14-3-3 HPRD SLiMFinder analysis vs. ELM.





**Figure 3.4. CompariMotif XGML output visualized with Cytoscape(Shannon *et al.* 2003) (recoloured).**

Motifs returned by SLiMFinder analysis of 14-3-3 interaction datasets are shown as blue ellipses. ELMs with CompariMotif matches are shown as rectangles. These are pale red by default but the following groups have been manually recoloured: 14-3-3 ligands, green; SH3 ligands, orange; phosphorylation motifs, yellow. Arrows proceed from parent to subsequence motifs (bidirectional where equal); Data Panel details for 14-3-3 ELMs and connected nodes and edges are shown.

## 4 Appendices

### 4.1 Troubleshooting & FAQ

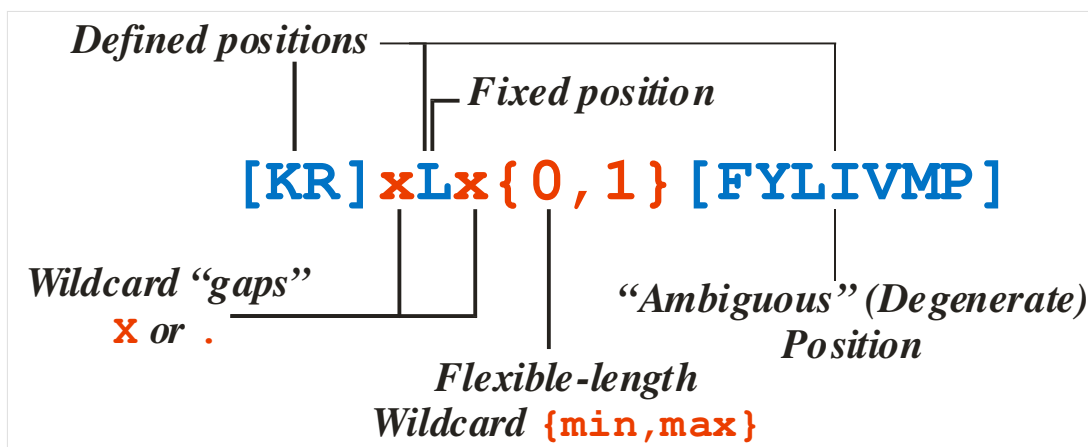
There are currently no specific Troubleshooting issues arising with SLiMfinder. Please see general items in the [PEAT Appendices](#) document and contact me if you experience any problems not covered.

### 4.2 SLiM Definitions

This covers the basic definitions needed to understand this manual. The term “motif” can be used in a number of different contexts with different meanings. In this manual, I use motif to mean a short, linear motif (SLiM) in a protein. In biology, SLiMs are functional microdomains with three main properties:

- *Short* – generally less than 10aa long with five or less defined residues.
- *Linear* – comprised of adjacent amino acids in a protein’s primary sequence. While three-dimensional conformation may be important for function, it is not necessary for definition.
- *Motifs* – there are some defined sequence patterns that are necessary for function and will therefore recur in the relevant proteins, allowing identification.

The basic anatomy of a SLiM is shown in [Figure 4.1](#).



**Figure 4.1. Anatomy of a SLiM.**

Definitions of different properties of SLiM have been marked on the example ELM, LIG\_CYCLIN\_1 (Puntervoll *et al.* 2003). This motif has three defined positions (one fixed and two degenerate) and two wildcard spacers (one fixed, one flexible-length) for a total length of 4-5aa.

## 4.3 References

- Amanchy R, Periaswamy B, Mathivanan S, Reddy R, Tattikota SG & Pandey A (2007). "A curated compendium of phosphorylation motifs." *Nat Biotechnol.* **25**(3): 285-6.
- Balla S, Thapar V, Verma S, Luong T, Faghri T, Huang CH, Rajasekaran S, del Campo JJ, Shinn JH, Mohler WA, Maciejewski MW, Gryk MR, Piccirillo B, Schiller SR & Schiller MR (2006). "Minimotif miner: A tool for investigating protein function." *Nat Methods.* **3**(3): 175-7.
- Davey NE, Shields DC & Edwards RJ (2006). "Slimdisc: Short, linear motif discovery, correcting for common evolutionary descent." *Nucleic Acids Res.* **34**(12): 3546-54.
- Davey NE, Edwards RJ & Shields DC (2007). "The slimdisc server: Short, linear motif discovery in proteins." *Nucleic Acids Res* **35**(Web Server issue): W455-9.
- Edwards RJ, Davey NE & Shields DC (2007). "Slimfinder: A probabilistic method for identifying over-represented, convergently evolved, short linear motifs in proteins." *PLoS ONE* **2**(10): e967.
- Mishra GR, Suresh M, Kumaran K, Kannabiran N, Suresh S, Bala P, Shivakumar K, Anuradha N, Reddy R, Raghavan TM, Menon S, Hanumanthu G, Gupta M, Upendran S, Gupta S, Mahesh M, Jacob B, Mathew P, Chatterjee P, Arun KS, Sharma S, Chandrika KN, Deshpande N, Palvankar K, Raghavnath R, Krishnakanth R, Karathia H, Rekha B, Nayak R, Vishnupriya G, Kumar HG, Nagini M, Kumar GS, Jose R, Deepthi P, Mohan SS, Gandhi TK, Harsha HC, Deshpande KS, Sarker M, Prasad TS & Pandey A (2006). "Human protein reference database-2006 update." *Nucleic Acids Res.* **34**(Database issue): D411-4.
- Neduva V, Linding R, Su-Angrand I, Stark A, de Masi F, Gibson TJ, Lewis J, Serrano L & Russell RB (2005). "Systematic discovery of new recognition peptides mediating protein interaction networks." *PLoS Biol.* **3**(12): e405.
- Neduva V & Russell RB (2005). "Linear motifs: Evolutionary interaction switches." *FEBS Lett.* **579**(15): 3342-5 Epub 2005 Apr 18.
- Neduva V & Russell RB (2006). "Dilimot: Discovery of linear motifs in proteins." *Nucleic Acids Res.* **34**(Web Server issue): W350-5.
- Punternvoll P, Linding R, Gemund C, Chabanis-Davidson S, Mattingsdal M, Cameron S, Martin DM, Ausiello G, Brannetti B, Costantini A, Ferre F, Maselli V, Via A, Cesareni G, Diella F, Superti-Furga G, Wyrwicz L, Ramu C, McGuigan C, Gudavalli R, Letunic I, Bork P, Rychlewski L, Kuster B, Helmer-Citterich M, Hunter WN, Aasland R & Gibson TJ (2003). "Elm server: A new resource for investigating short functional sites in modular eukaryotic proteins." *Nucleic Acids Res* **31**(13): 3625-30.
- Shannon CE (1997). "The mathematical theory of communication. 1963." *MD. Comput.* **14**: 306-317.
- Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B & Ideker T (2003). "Cytoscape: A software environment for integrated models of biomolecular interaction networks." *Genome Res* **13**(11): 2498-504.