# RepeatMatcher

http://repeatmasker.org/

INSTITUTE FOR
Systems
Biology

User Manual
version 1.0

Juan Caballero
jcaballero@systemsbiology.org
2012

# What is RepeatMatcher?

RepeatMatcher is a simple analysis tool to help in the manual annotation of RepeatModeler sequences.

# Motivation

After running RepeatModeler you probably will have a well identified sequences and a bunch of sequences to be manually annotated. RepeatMatcher is designed to help in this hard and repetitive task. We provide a simple pipeline to generate all the required files and a graphical user interface to manipulate and annotate the sequences. The GUI is intended to show all the information to make the decisions quick and easy for each sequence.

# The pipeline

## *Installation*

The pipeline is a Perl script (RepatMatcher.pl) that control all the processes, to be able to run the pipeline, first you need to resolve all the dependencies:
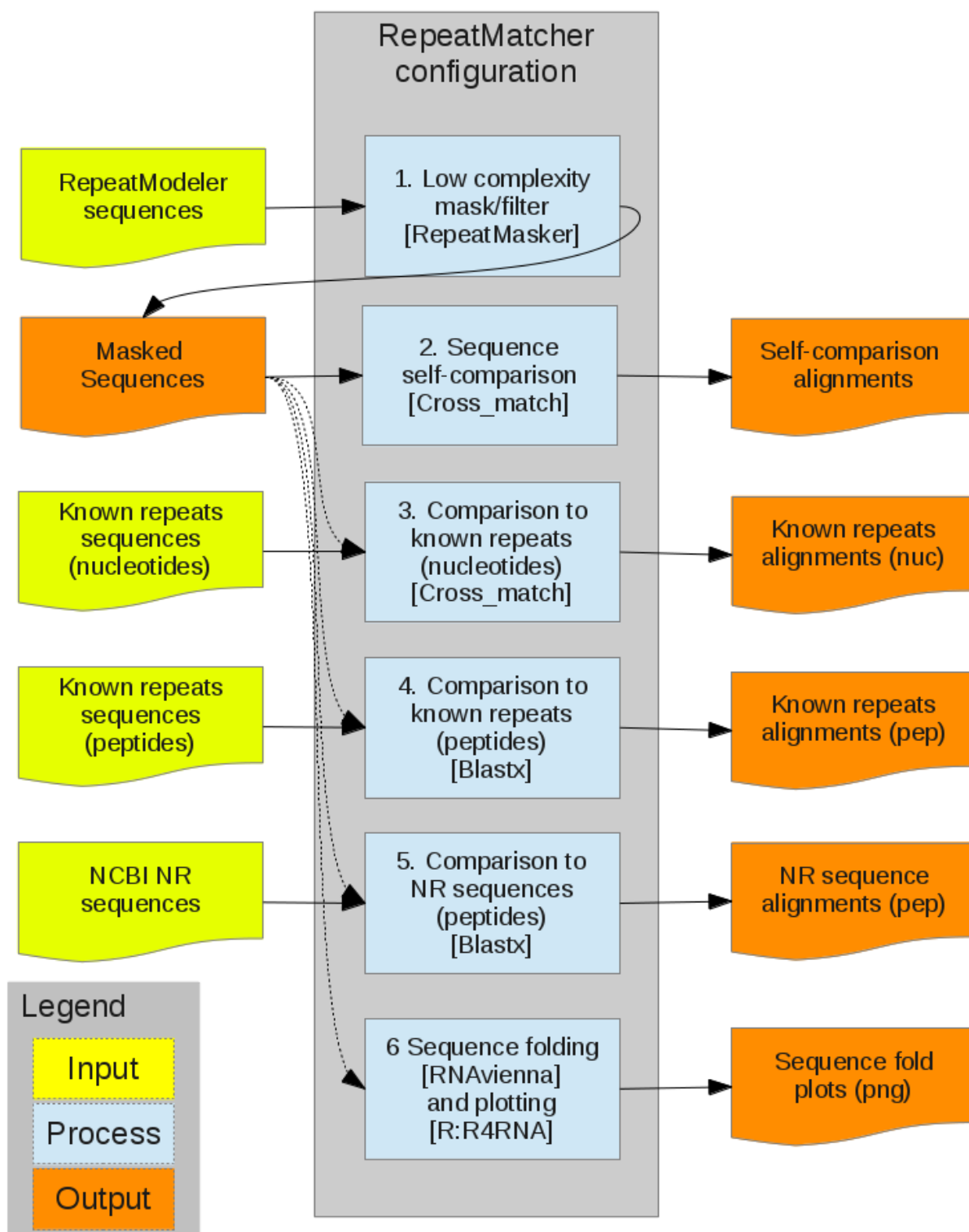- RepeatMasker
- Blast
- Cross_match
- RNAVienna
- R
- R package R4RNA
- Databases: NCBI NR, known repeats (in nucleotides and proteins).

## *Process description*

We started with the RepeatModeler consensi output file (fasta), the steps are:
1. Masking all the low complexity regions with RepeatMasker and discard sequences with filters for a minimal composition.
2. All sequences will be compared with itself using Cross_match.
3. All sequences are compared in nucleotide space with known repeats consensi with Cross_match.
4. All sequences are compared in protein space with known repeats consensi with Blastx.
5. All sequences are compared in protein space with the sequences in NCBI-NR data base with Blastx.
6. All sequences are folded with RNA-Vienna and the principal structure is plotted with R:R4RNA.

# RepeatMatcher Pipeline

**RepeatMatcher configuration**

RepeatModeler sequences → 1. Low complexity mask/filter [RepeatMasker]

Masked Sequences

2. Sequence self-comparison [Cross_match] → Self-comparison alignments

Known repeats sequences (nucleotides) → 3. Comparison to known repeats (nucleotides) [Cross_match] → Known repeats alignments (nuc)

Known repeats sequences (peptides) → 4. Comparison to known repeats (peptides) [Blastx] → Known repeats alignments (pep)

NCBI NR sequences → 5. Comparison to NR sequences (peptides) [Blastx] → NR sequence alignments (pep)

6 Sequence folding [RNAvienna] and plotting [R:R4RNA] → Sequence fold plots (png)

**Legend**

Input

Process

Output

## *Configuration file*

RepeatMatcher uses a plain text file (`RepeatMatcher.conf`) to define some parameters, but if it's required, you can adjust the parameters before running the pipeline editing this file or activating the `-e` option.

The follow is an example of this file:

```
# RepeatMatcher configuration file
# Juan Caballero, Institute for Systems Biology @ 2012
# This is a simple configuration file, syntax is:
#          VARIABLE: VALUE
# One variable per line, blank lines are omitted.

# STEP 1. Masking low complexity sequences
# min_mask => minimal % of masked sequence
# min_size => minimal size after masking
min_mask: 90
min_size: 30

# STEP 2. Self-comparison
# crossmatch_self => cross_match parameters
crossmatch_self: -M data/nt_sub.matrix -gap_init -25 -gap_ext -5 -minscore 200 -
minmatch 9 -minscore 200

# STEP 3. Known annotation comparison
# crossmatch_comp => cross_match parameters
crossmatch_comp: -M data/nt_sub.matrix -gap_init -25 -gap_ext -5 -minscore 200 -
minmatch 9 -minscore 200

# STEP 4. Blastx to known repeat proteins
# blastx_rep => blast parameters
blastx_rep: -d data/repeats -W 2 -v 5 -b 5 -F F -e 0.01

# STEP 5. Blastx to NR database
# blastx_nr => blast parameters
blastx_nr: -d data/nr -W 2 -v 5 -b 5 -F F -e 1e-6

# STEP 6. Fold DNA
# fold => folding (RNAVienna) parameters
fold: --noconv --noGU --ImFeelingLucky
plot_w = 600
plot_h = 300
```

## *Launching the pipeline*

The basic operation is:
```
perl RepeatMatcher.pl  -s sequences.fa -k knownrepeats.fa -o NewAnnotation
```

for a complete list of options:
```
perl RepeatMatcher.pl --help
```

# THE GUI

## *Installation*

The program is a Perl:Tk development, so you need to install it from CPAN.

After you ran the pipeline, you have all the files required for `RepeatMatcherGUI.pl`. The GUI uses a text file to keep all the parameters and annotations made (the LOG file). The first lines of that file record the files used, the after that, the rest of lines records all changes for each sequence.

## *Launching the GUI*

```
perl RepeatMatcherGUI.pl -o OUT -i SEQS -s SELF -a ALIGN -b BLAST -n BLAST \
     -f FOLD -l LOG
```

Launching the GUI with a previous generated LOG:
```
perl RepeatMatcherGUI.pl -r LOG
```

for a complete list of options:
```
perl RepeatMatcherGUI.pl --help
```

## *Description of the GUI*

The program will call a series of windows, this are:
1. Annotation window
2. Sequence window
3. Self-alignments window
4. Alignments to known repeats (nucleotides) window
5. Alignments to known repeats (protein) window
6. Alignments to NR proteins
7. Sequence folding window

The annotation window is used to check each repeat and corroborate or edit the annotation, the sequence also can be marked as "Excluded", or can be flagged to be reverse complement with the "Reverse" mark. Each time the "Update" button is pushed, the LOG file is updated too, to close the session just close all windows.

The Annotation window also shows the status of the sequence using a color code in the labels:
- black – sequence is not saved or modified yet
- blue – sequence is saved without modifications
- red – sequence is saved with modifications

When the annotation is done, you can write the final sequences in a new fasta file pushing the "Export fasta" button.