

# psRobot manual

## DESCRIPTION

psRobot is a smRNA analysis software package, which so far contains four functions: (1) psRobot\_map is designed to find all the perfect matching locations of short sequences (less than 40bp) in longer reference sequences; (2) psRobot\_mir is designed to find potential miRNAs or small hairpin RNAs (shRNAs) from high throughput sequencing data. (3) psRobot\_tar is designed to find potential small RNA targets on a large scale. (4) psRobot\_deg is designed to identify which small RNA targets are supported by user specified degradome data.

## REQUIREMENTS

**perl 5** or greater

**gcc 4.0** or greater

psRobot\_mir requires “nafold” command in **Mfold-3.5** software to get the potential folding structure of precursor RNAs.

## Installation

The simplest way to install this package is (require root permission):

```
./configure
make
make install (log as root)
```

Install psRobot in an alternative path with no root permission:

```
./configure -p <your full path> -l <your full path>
make
make install
```

Configure options:

```
-h          display this help and exit
-p PREFIX  install architecture-independent files in PREFIX
           [complete path needed, default /usr/bin]
-l LIBS    install perl libs in LIBS
           [complete path needed, default standard
           perl lib, the path next to last in your @INC]
```

## COMMANDS AND OPTIONS

**psRobot\_map**      **psRobot\_map** <short\_sequences> <reference\_sequences> <output>

<short\_sequences> is the first input file which contains the short sequences separated by <Tab>:

```
<ID>        <SEQ> (Don't include this row!)
SrID001     TGACAGAAGAGAGTGAGCAC
SrID002     TCGCTTGGTGCAGATCGGGAC
SrID003     TTGACAGAAGAGAGTGAGCAC
SrID004     TGAAGCTGCCAGCATGATCTGA
SrID005     GGCGGATGTAGCCAAGTGGA
```

(Don't use “blank” in your short sequence IDs!)

<reference\_sequences> is the second input file which contains the reference sequences in fasta format.

(Don't use blank in your reference sequence IDs!)

<output> is the output file containing the mapping results:

ID	ref	strand	start	end	seq
SrID003	ref01	+	28	48	TTGACAGAAGAGAGTGAGCAC
SrID001	ref01	+	29	48	TGACAGAAGAGAGTGAGCAC
SrID001	ref01	+	60	79	TGACAGAAGAGAGTGAGCAC

Field	Description
ID	Short sequence ID in <short_sequences>
ref	Reference sequence ID in <reference_sequences>
strand	Mapping strand
start	Start position
end	End position
seq	Short sequence

## psRobot\_mir

**psRobot\_mir -s <smRNA> -g <genome>**

**Note:** Different psRobot\_mir run must be in **different folders**.

### Options:

**-s** input file name which contains smRNA deep sequencing data in tab-delimited (tsv) format. [default: smRNA]

Suggested format:

<smRNA sequence><Tab><reads1><Tab><reads2><Tab><reads3>...

The first column must be smRNA sequences. The following columns are the smRNA clone counts in various samples/conditions. Columns must be separated by <Tab>.

**-g** input file name which contains reference genome/contig sequences in FASTA format. [default: genome]

**-k** input file name which contains known miRNA GFF3 file corresponding to the reference genome ([GFF3 format](#)). This file is needed only if one wants to exclude known miRNAs from prediction results. [default: kmiRNA]

**-r** clone counts selection: minimal smRNA clone counts. SmRNAs with less than this clone counts will be excluded from the following analysis. [default: 2]

**-l** loci selection: the maximal number of smRNA mapping locations in reference genome. SmRNAs with more than this number of mapping locations will be excluded from the following analysis. [default: 20]

**-cg** minimal number of nucleotides between adjacent smRNA clusters. [default: 200]

**-cl** maximal length limitation of smRNA clusters selected to predict stem-loop precursors. [default: 300]

**-cr** all the clone counts of smRNAs in one cluster will be summed up to represent the smRNA cluster's clone counts. SmRNA clusters with more than this clone counts will be selected to predict stem-loop precursors. [default: 10]

**-mr** the highest expressed smRNA in one cluster higher than this clone counts will be selected to predict stem-loop precursors. [default: 10]

**-mml** minimal number of mismatches in supposed miRNA mature region. [default: 1]

**-mmh** maximal number of mismatches in supposed miRNA mature region. [default: 5]

**-ll** retain large loop miRNAs (T/F). [default: F]

**Output:**

<\*.StarInfo> contains the major information of predicted miRNAs:

```

Col_1      Col_2      Col_3      Col_4      Col_5      Col_6      Col_7
Sr7_1_1_21_188_60  CTGAAAGTITGGGGGACTC  2      Sr4_1_1_21_30_3  gttccotttaacgcttcattg  ref01:20:132:+  attgacttcaaaaatatgagttccotttaacgcttcattggtgaactca
Sr9_1_1_21_2000_10  AGAAACTTCTGAGACCAA  2      Sr5_1_1_21_30_1  tggctctcagaagttctcttg  ref02:708:910:-  gaggagcaagccttgatgctgcagaagtgagggtttgggttcagaaag
Sr3_1_1_21_200_100  ACAATGATCTGCAICTTTCATT  2      Sr1_1_1_7_3      tgcasaagatgcagatcatatgccc  ref03:6953:7055:+  attctctgaggcttctgatacaccggACAATGATCTGCAICTTTCATTtctc
Sr12_1_1_21_9_33  AAKAGTGCAGTCTATATGTC  2      Sr19_1_1_21_2_7  acaatgactctgactctctgccc  ref04:4972:7034:-  atggaactatgactctgactctctgactctgaaaggaAAKAGTGCAGTCTAT
Sr2_1_1_21_3_14  AACTAATTTTATTGGACGTTTA  2      -              -              -              ref04:11:173:-  atatacgggaaaaactcgaaaaaacctcAACTAATTTTATTGGACGTTT
  
```

Col	Description
1	Information of mature miRNA. Format: [miRNA ID] _ [No. of miRNA mapping loci] _ [order of miRNA mapping location] _ [miRNA length] _ [Reads1] _ [Reads2]...
2	miRNA sequence
3	No. of smRNA clusters in miRNA precursor. The canonical miRNAs always have two clusters on their precursors, one is at miRNA mature region and the other is at miRNA* region.
4	Information of miRNA*. "." represents no miRNA* sequence can be detected. Format: [miRNA* ID] _ [No. of miRNA* mapping loci] _ [order of miRNA* mapping location] _ [miRNA* length] _ [Reads1] _ [Reads2]...
5	miRNA* sequence. "-" represents no miRNA* sequence can be detected.
6	Location of miRNA precursor in reference genome. Format: [reference genome ID] _ [start] _ [end] _ [strand]
7	miRNA precursor sequence

<\*.Struc> contains the stem-loop folding structures of predicted miRNA precursors.

```

Sequence      1 Structure      1
Folding bases 1 to 113 of Sr7_1_1_21_188|177|115 (a)
dG = -44.30 (b)
          10      20      30      40      50      (c)
at----| t   aaaaat      tt      t   a   tca   ac
      tgac tcaa      gagttccct aacgcttca tgttg atac aagcc \
      actg ggtt      CTCAGGGGG ITGTGAAGT acaat tatg tttgg a
gttcct^ t   -----  GI      C   a   ---  tt
110      100      90      80      70
  
```

Field	Description
(a)	Information of mature miRNA. Format: [miRNA ID] _ [No. of miRNA mapping loci] _ [order of miRNA mapping location] _ [miRNA length] _ [Max Reads in all samples]    [No. of upstream nucleotides to consist precursor]    [No. of downstream nucleotides to consist precursor]
(b)	Free energy of the precursor folding structure
(c)	Precursor folding structure

<\*.Reads> contains the small RNA reads distribution on predicted miRNA precursors.

```

(a) >Sr7_1_1_21_188_60 ref01|20|132|+
(b) attgacttcaaaaatatgagttccotttaacgcttcattggtgaactcaaaagccacattggtttgtatataacaCTGAAAGTGTTTGGGGGACTCtttggtcactccttg 212 69
(c) *****ctgaagtgtttgggggactc***** 1 21 188 60
*****ctgaagtgtttgggggactc***** 1 21 188 60
(d) -----atgagttccotttaacgcttc----- 1 21 10 10
-----gttccotttaacgcttc----- 1 19 1 3
-----gttccotttaacgcttc----- 1 20 5 8
-----gttccotttaacgcttc----- 1 21 3 3
-----ctgaagtgtttgggggac----- 1 18 10 7
-----ctgaagtgtttgggggac----- 1 19 9 1
-----ctgaagtgtttgggggactc----- 1 21 188 60
-----tgaagtgtttgggggactc----- 1 20 5 1
  
```

Field	Description
-------	-------------



<*degradome\_data*> is the first input file which contains the degradome sequences data in tab-delimited (tsv) format:

```

<Sequence>                                <Counts> (Don't include this row!)
TGACAGAAGAGAGTGAGCAC                      108
TCGCTTGGTGCAGATCGGGAC                    373
TTGACAGAAGAGAGTGAGCAC                      10

```

<*target\_sequences*> is the second input file which contains the target sequences in FASTA format as prediction library and **must be identical to** <*target*> file in psRobot\_tar.

<*smRNA-target.gTP*> is the third input file containing the target prediction results, which could be the direct output of psRobot\_tar.

<*output*> is the output file containing the degradome data supported smRNA target prediction results:

```

>ath_miR156a   Score: 1.0   (a) Deg: 285:2378:285:1635 AT1G27370.1
Query:         1 TGACAGAAGAGAGTGAGCAC 20
                |||*|||
Sbjct:        2387 ACTGTCTTCTCTCTCGTG 2368
(b)           |-----2
                |-----1
                |-----1
                |-----3
                |-----285
                |-----152
                |-----8
                |-----6
                |-----1
                |-----2
                |-----3
                |-----1

```

Field	Description
(a)	Degradome support. Format: Deg: [reads of predicted cleavage site] : [position of predicted cleavage site on target sequence] : [maximal reads of all cleavage sites on target sequence] : [total reads of all cleavage sites on target sequence]
(b)	Degradome sequence distribution on smRNA target sites. Number on each row represents the reads of degradome sequence on each cleavage site ( <b>cleaved on the right-hand of the pointing nucleotide</b> ).