

# ReRCoP

Recombination Removal for Core-genome Phylogeny

## Prerequisites

1. Linux environment
2. Python 2.7
3. Python package: 'scipy'
4. R (if removal plot is needed)
5. BLAST package (makeblastdb, blastn)

## Usage

```
1. python ReRCoP.py [options] Genomes.fasta
2.
3. Options:
4.     --version          show program version number and exit
5.     -h, --help         show this help message and exit
6.
7. Input Options:
8.     -a, --aligned       Set this if genome sequences in the input file are
9.                         already aligned.
10.                        # For aligned genomes
11.
12.     --gbk=GBK           Input GenBank file of the reference genome.
13.                        # For aligned genomes using core genome approach
14.
15.     -w, --window        Set this if sliding windows instead of genes are to be
16.                         considered.
17.                        # For aligned genomes using complete genome approach
18.
19.     --fSize=FSIZE       Fragment size if using sliding window. [Default: 1000]
20.                        # For aligned genomes using complete genome approach
21.
22.     --sSize=SSIZE       Step size if using sliding window. [Default: 500]
23.                        # For aligned genomes using complete genome approach
24.
25.     --cds=CDS           Input coding sequences in fasta format. Used to determine
26.                         the core genome when input genomes are not aligned.
27.                        # For unaligned genomes using core genome approach
28.
29. Core Gene Identification Options:
30.     --cov=COV           Minimum sequence coverage to regard genes as present.
31.                        [Default: 0.7]
32.                        # For aligned or unaligned genomes using core genome
33.                        # approach
34.
35.     --sim=SIM           Minimum sequence similarity to regard genes as present.
36.                        [Default: 70]
```

```

37.                                     # For unaligned genomes using core genome approach
38.
39.   Outlier Removal Options:
40.     -m METHOD, --method=METHOD
41.                                     Outlier removal method. Can be 'Grubbs', 'kNN', or
42.                                     'DBSCAN', or can be multiple methods separated by ','.
43.
44.     --alpha=ALPHA                   For 'Grubbs' method: Significance level in Grubbs test.
45.                                     [Default: 0.05]
46.
47.     --radius=RADIUS                 For 'kNN' method: Maximum number of differences for
48.                                     a point to be considered as a neighbor (in the unit of
49.                                     standard deviation of all pair-wise number of differences
50.                                     ). [Default: 1.5]]
51.
52.     --k=K                           For 'kNN' method: Minimum number of neighbors for a
53.                                     non-outlier point (in the unit of total number of points
54.                                     ). [Default: 0.2]
55.
56.     --eps=EPS                       For 'DBSCAN' method: Maximum number of differences
57.                                     between two points for them to be considered as in the
58.                                     same neighborhood (in the unit of standard deviation of
59.                                     all pair-wise number of differences). [Default: 1]
60.
61.     --minP=MINP                     For 'DBSCAN' method: Minimum number of
62.                                     points required to form a dense region (in the unit of
63.                                     total number of points). [Default: 0.2]
64.
65.   Output Options:
66.     -o OUTDIR, --outdir=OUTDIR      Output directory. [Default: running directory]
67.     -p PREFIX, --prefix=PREFIX      Output prefix. [Default: ReRCoP]

```

## Input files

Input files for aligned genomes are straightforward and thus not stated again here. This part will be focused on unaligned genomes, where core genomes are to be identified and extracted by ReRCoP.

The following input files are required:

1. A file in multiple nucleotide fasta format with gene coding sequences (All the coding sequences from any one of the samples will do).
  - If one of the input is complete genome with annotation from the public database:
    - 1> Download coding sequence in multiple nucleotide fasta format.
    - 2> Remove duplicated genes.
    - 3> Remove phage genes.
    - 4> Remove genes with CRISPR sequence.
  - Else if one of the input is complete genome without annotations from the public database:
    - 1> Predict coding sequences with software like [prodigal](#).
    - 2> Remove duplicated genes.

- 3> Remove phage genes.
- 4> Remove genes with CRISPR sequence.
- Else if none of the input files are complete genomes but are raw sequencing reads, do the following:
  - 1> Use assembly tools for *de novo* assembly to get files of contigs for each sample.
  - 2> Use one of the samples with good assembly quality, predict coding sequences with software like [prodigal](#).
  - 3> Remove duplicated genes.
  - 4> Remove phage genes.
  - 5> Remove genes with CRISPR sequence.
2. A file in multiple nucleotide fasta format with genome sequences.
  - Complete sequences: Should be in fasta format.
  - Assembled contigs: Concatenate the contigs to form one fake genome sequence, which can be done with the script `FormatContig.pl` in `./scripts`. Then concatenate all genome sequences or fake genome sequences to form a multiple-fasta file.

```
1.          perl FormatContig.pl <contig fasta file> <header of the output fasta file>
          > <output fasta file>
```

## Cautions

Take note about giving different header names for the fasta file. If the header file contains blanks, the first column should all be different.

## Output files

- **.core.fasta** The concatenated core genomes before recombination removal.
- **.concatenation.log** The concatenation log of the core.fasta file that is composed of each gene sequence name and the respective start and end position in the concatenated core genome.
- **.snpmat** A matrix of scaled number of SNPs in each gene in each genomic sequence.
- **.DBSCAN.outliermat** A matrix of recombinant genes identified by DBSCAN with '1' denoting recombinant while '0' denoting non-recombinant.
- **.DBSCAN.removal.fasta** The concatenated core genomes after DBSCAN recombination removal.
- **.Grubbs.outliermat** A matrix of recombinant genes identified by Grubbs with '1' denoting recombinant while '0' denoting non-recombinant.
- **.Grubbs.removal.fasta** The concatenated core genomes after Grubbs recombination removal.
- **.kNN.outliermat** A matrix of recombinant genes identified by kNN with '1' denoting

recombinant while '0' denoting non-recombinant.

- **.kNN.removal.fasta** The concatenated core genomes after kNN recombination removal.