

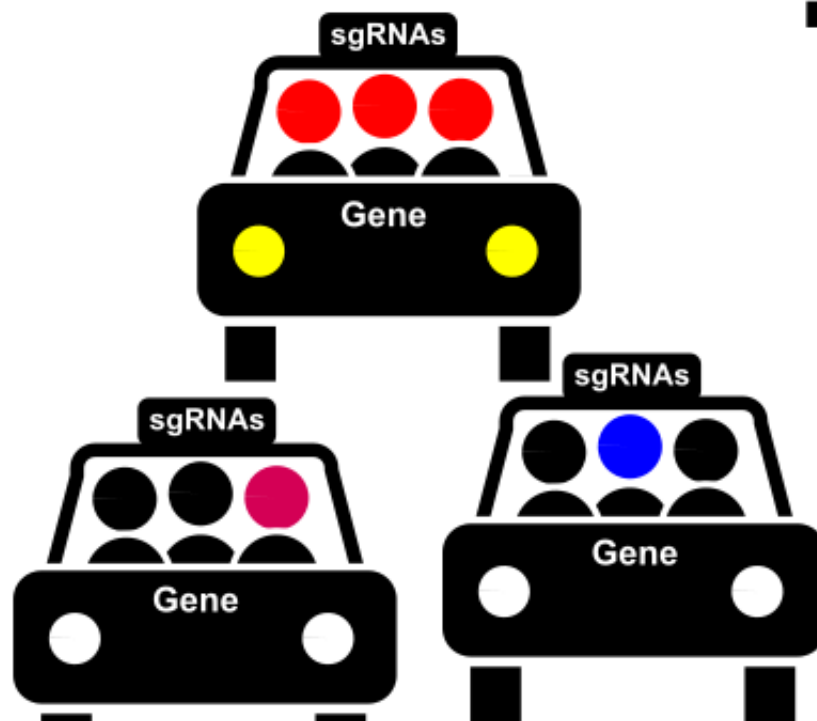
caRpools Shortcut User Guide

Jan Winter

Contents

| | | |
|----------|---|----------|
| 1 | Files and Folder Structure to use CaRpools | 3 |
| 1.1 | Setup Files and R-Studio | 6 |
| 1.2 | Check Setup | 6 |
| 1.3 | Example of a MIACCS File entry for FASTQ files | 6 |
| 1.4 | Example of a MIACCS File entry for Read-Count files | 7 |
| 2 | Start CaRpools Report Generation | 7 |
| 2.1 | Start CaRpools using R-Studio | 7 |
| 2.2 | Start CaRpools using R console | 8 |

CRISPR-AnalyzeR for Pooled Screens



Transparent. Reproducible.

caRpools...

Exploratory data analysis of CRISPR/CAS screens

1 Files and Folder Structure to use CaR pools

Please note: the MAIN FOLDER must be the R working directory!

Data and Script paths can be adjusted in the MIACCS file.

The following files are necessary to use CaR pools for report generation:

MIACCS.xls

Minimum Information About CRISPR/Cas Screens. This file needs to be filled out to provide all necessary informations about the screen.

R Markdown Template files

Either CaR pools-extended-PDF.rmd, CaR pools-PDF.rmd or CaR pools-extended-HTML.rmd or CaR pools-HTML.rmd. Is the template for report generation.

Data Files

Two replicates per Control and Treated. Can be FASTQ files OR already mapped, not normalized read count files.

CRISPR-mapping.pl

PERL script to map your extracted FASTQ files, if desired (as indicated in the MIACCS.xls)

CRISPR-extract.pl

PERL script to extract 20 nt target sequence from FAST files, if desired (as indicated in the MIACCS.xls)

CaR pools.png

The logo file

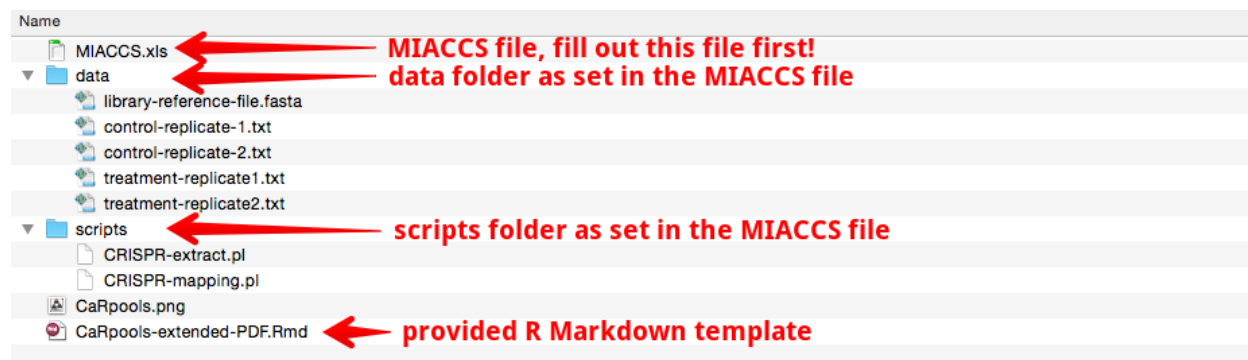
The following files are necessary to use *single* CaR pools functions:

Data Files

Either raw read count files or FASTQ files (that need to be extracted and mapped using CaR pools)

Please note that CaR pools always starts with loading data files. For raw-readcount files, use `load.file`. For FASTQ files, please see the sections below.

CaR pools folder structure for Report Generation using raw Read Count files:



CaR pools folder structure for Report Generation using raw Read Count files AFTER REPORT GENERATION:

| Name | |
|--|--|
| MIACCS.xls | |
| CaRpoools-TRAIL | ← Folder with all created plots |
| data | |
| library-reference-file.fasta | |
| CaRpoools-TRAIL_ANNOTATION.xls | |
| CaRpoools-TRAIL_COMPARE-HITS.xls | |
| CaRpoools-TRAIL_DROPOUT.xls | ← Output Tables in DATAPATH |
| CaRpoools-TRAIL_FINAL.xls | |
| CaRpoools-TRAIL_HIT-CALLING.xls | |
| CaRpoools-TRAIL_HITS-sgRNA-depleted.xls | |
| CaRpoools-TRAIL_HITS-sgRNA-enriched.xls | |
| CaRpoools-TRAIL_STATS.xls | |
| control-replicate-1.txt | |
| control-replicate-2.txt | |
| treatment-replicate1.txt | |
| treatment-replicate2.txt | |
| scripts | |
| CaRpoools-extended-PDF.pdf | ← The PDF Report |
| CaRpoools.png | |
| CaRpoools-TRAIL-ANALYSIS-MAGeCK-RAW.log | |
| CaRpoools-TRAIL-ANALYSIS-MAGeCK-RAW.R | ← RAW MAGeCK analysis files |
| CaRpoools-TRAIL-ANALYSIS-MAGeCK-RAW_summary.Rnw | |
| CaRpoools-extended-PDF.Rmd | |
| CaRpoools-extended-PDF.tex | ← TEX file used for PDF generation |
| CaRpoools-TRAIL-ANALYSIS-MAGeCK-RAW.gene_summary.txt | |
| CaRpoools-TRAIL-ANALYSIS-MAGeCK-RAW.sgrna_summary.txt | |
| CaRpoools-TRAIL-ANALYSIS-DESeq2-sgRNA.tab_DESeq2_sgRNA.tab | ← RAW MAGeCK and DESeq2 analysis files |
| CaRpoools-TRAIL-ANALYSIS-MAGeCK-RAW_MAGeCK_sgRNA.tab | |

CaRpoools folder structure for Report Generation using FASTQ files:

| Name | |
|------------------------------|---|
| MIACCS.xls | |
| data | |
| library-reference-file.fasta | |
| control-replicate-1.fastq | ← FASTQ files instead of read count files |
| control-replicate-2.fastq | |
| treatment-replicate1.fastq | |
| treatment-replicate2.fastq | |
| scripts | |
| CRISPR-extract.pl | |
| CRISPR-mapping.pl | |
| CaRpoools.png | |
| CaRpoools-extended-PDF.Rmd | |

CaRpoools folder structure for Report Generation using FASTQ files AFTER REPORT GENERATION:

| Name | |
|---|--|
| MIACCS.xls | |
| CaRpoools-TRAIL | |
| data | |
| library-reference-file.1.bt2 | |
| library-reference-file.2.bt2 | |
| library-reference-file.fasta | Bowtie2 Index files |
| library-reference-file.rev.1.bt2 | |
| library-reference-file.rev.2.bt2 | |
| CaRpoools-TRAIL_ANNOTATION.xls | |
| CaRpoools-TRAIL_COMPARE-HITS.xls | |
| CaRpoools-TRAIL_DROPOUT.xls | |
| CaRpoools-TRAIL_FINAL.xls | Output Tables |
| CaRpoools-TRAIL_HIT-CALLING.xls | |
| CaRpoools-TRAIL_HITS-sgRNA-depleted.xls | |
| CaRpoools-TRAIL_HITS-sgRNA-enriched.xls | |
| CaRpoools-TRAIL_STATS.xls | |
| control-replicate-1_extracted.fastq | extracted 20 nt target sequences |
| control-replicate-1_extracted.sam | Bowtie2 alignment file |
| control-replicate-1-designs.txt | |
| control-replicate-1-genes.txt | |
| control-replicate-1.fastq | |
| control-replicate-2_extracted.fastq | |
| control-replicate-2_extracted.sam | |
| control-replicate-2-designs.txt | Raw read count files for every sgRNA |
| control-replicate-2-genes.txt | Raw read count files summed up for genes |
| control-replicate-2.fastq | |
| treatment-replicate1_extracted.fastq | |
| treatment-replicate1_extracted.sam | |
| treatment-replicate1-designs.txt | |
| treatment-replicate1-genes.txt | |
| treatment-replicate1.fastq | |
| treatment-replicate2_extracted.fastq | |
| treatment-replicate2_extracted.sam | |
| treatment-replicate2-designs.txt | |
| treatment-replicate2-genes.txt | |
| treatment-replicate2.fastq | |
| scripts | |
| CaRpoools-extended-PDF.pdf | |
| CaRpoools.png | |
| CaRpoools-TRAIL-ANALYSIS-MAGeCK-RAW.log | |
| CaRpoools-TRAIL-ANALYSIS-MAGeCK-RAW.R | |
| CaRpoools-TRAIL-ANALYSIS-MAGeCK-RAW_summary.Rnw | |
| CaRpoools-extended-PDF.Rmd | |

1.1 Setup Files and R-Studio

All packages and software tools need be installed correctly as shown before.

1. Copy all files in the designated folders as shown above.

- **Please note: the MAIN FOLDER must be R working directory!**
 - The MIACCS.xls as well the R markdown template and CaRpools.png must be in the same folder as the R working dir.
2. Adjust the path to the data and scripts folder if necessary in the MIACCS.xls . Use the absolute path. If the folder structure is as shown above, you do not need to make any adjustments.
 3. Adjust and fill out the **MIACCS.xls** file.
 4. You can use `CarPools(type="check")` to check for the correct folder structure and data file presence as it is indicated in the MIACCS.xls file.
 5. You can check for your R working directory by `getwd()` and set it to any directory you want by `setwd("/PATH")`.

1.2 Check Setup

You can verify that the MIACCS.xls file as well as the used template file and all necessary scripts are found by calling `check.caRpools()`.

See below for more information about the arguments.

By default, it requires a correct MIACCS file + the script files + all packages installed + MAGeCK + Bowtie2 + Pandoc.

1.3 Example of a MIACCS File entry for FASTQ files

| | A | B | C | D | E | F |
|----|---|------------------------|--|---|---|---|
| 1 | Files And Storage | | | | | |
| 2 | Absolute Path to CRISPR-extract.pl and CRISPR-mapping.pl | /PATH/TO/scripts | ← absolute path | | | |
| 3 | Absolute Path to Data files | /PATH/TO/data | | | | |
| 4 | Filename Untreated Replicate 1 | untreated_NGS_file1 | ← provide FASTQ filename WITHOUT ending! | | | |
| 5 | Name of Untreated Replicate 1 | Untreated #1 | | | | |
| 6 | Filename Untreated Replicate 2 | untreated_NGS_file2 | | | | |
| 7 | Name of Untreated Replicate 2 | Untreated #2 | ← Name of Sample | | | |
| 8 | Filename Treated Replicate 1 | treated_NGS_file1 | | | | |
| 9 | Name of Treated Replicate 1 | Treated #1 | | | | |
| 10 | Filename Treated Replicate 2 | treated_NGS_file2 | | | | |
| 11 | Name of Treated Replicate 2 | Treated #2 | | | | |
| 12 | Name of library reference file (without.fasta extension) | library-reference | ← library-reference fasta file, WITHOUT .fasta ending | | | |
| 13 | In which column is the gene identifier? | 1 | | | | |
| 14 | In which columns is the read count? | 2 | | | | |
| 15 | Gene identifier of positive controls (comma separated) | positive1,positive | ← Positive controls (gene identifiers) | | | |
| 16 | Gene Identifier of non-targeting control | random | | | | |
| 17 | Gene Identifier Extraction | | | | | |
| 18 | Regular Expression to extract Gene from sgRNA identifier | ^(.+?)(_+) | ← PERL RegEx pattern to extract gene from sgRNA identifier | | | |
| 19 | Data Extraction and Bowtie2 Mapping/Alignment | | | | | |
| 20 | Do you want to extract the sgRNA target sequence from FASTQ? | TRUE | ← If you provide Fastq files, this must be TRUE | | | |
| 21 | Regular expression to extract target sequence from FASTQ file | ACC(.{20})GT{2,4}AGAGC | ← PERL RegEx to extract 20 nt target sequence from reads | | | |
| 22 | Regular Expression to extract machine ID from Reads | @(M01100.+) | ← Machine ID from your sequencer | | | |
| 23 | Is the data within the FASTQ file in Reverse Complement? | FALSE | | | | |
| 24 | Do you want to map the reads to the reference file? | TRUE | ← FASTQ data needs to be mapped to reference | | | |
| 25 | Do you want to create the Bowtie2 index files? | TRUE | ← For this we need a bowtie2 index from your library | | | |
| 26 | How many threads shall bowtie2 use? | 4 | | | | |
| 27 | Bowtie2 Sensitivity? | very-sensitive-local | | | | |
| 28 | Additional Bowtie2 parameters? | | | | | |
| 29 | Alignment Quality? | perfect | | | | |

1.4 Example of a MIACCS File entry for Read-Count files

| A | B | C | D | E | F | G | H |
|---|-------------------------------|---|---|---|---|---|---|
| Files And Storage | | | | | | | |
| Absolute Path to CRISPR-extract.pl and CRISPR-mapping.pl | /PATH/TO/scripts | absolute path to script files and data files | | | | | |
| Absolute Path to Data files | /PATH/TO/data | | | | | | |
| Filename Untreated Replicate 1 | untreated_readcount_file1.txt | Read-Count file, WITH extension | | | | | |
| Name of Untreated Replicate 1 | Untreated #1 | | | | | | |
| Filename Untreated Replicate 2 | untreated_readcount_file2.txt | | | | | | |
| Name of Untreated Replicate 2 | Untreated #2 | | | | | | |
| Filename Treated Replicate 1 | treated_readcount_file1.txt | | | | | | |
| Name of Treated Replicate 1 | Treated #1 | | | | | | |
| Filename Treated Replicate 2 | treated_readcount_file2.txt | | | | | | |
| Name of Treated Replicate 2 | Treated #2 | library reference FASTA file, WITHOUT extension | | | | | |
| Name of library reference file (without.fasta extension) | library-reference | | | | | | |
| In which column is the gene identifier? | 1 | | | | | | |
| In which columns is the read count? | 2 | | | | | | |
| Gene Identifier of positive controls (comma separated) | positive1,positive | positive control gene identifier | | | | | |
| Gene Identifier of non-targeting control | random | | | | | | |
| Gene Identifier Extraction | | | | | | | |
| Regular Expression to extract Gene from sgRNA identifier | ^(.+?)(_.+) | | | | | | |
| Data Extraction and Bowtie2 Mapping/Alignment | | | | | | | |
| Do you want to extract the sgRNA target sequence from FASTQ? | FALSE | since already mapped read-count files are provided, set this to FALSE | | | | | |
| Regular expression to extract target sequence from FASTQ file | ACC(.{20})GT{2,4}AGAGC | | | | | | |
| Regular Expression to extract machine ID from Reads | M01100 | | | | | | |
| Is the data within the FASTQ file in Reverse Complement? | FALSE | | | | | | |
| Do you want to map the reads to the reference file? | FALSE | | | | | | |
| Do you want to create the Bowtie2 index files? | FALSE | | | | | | |
| How many threads shall bowtie2 use? | 4 | | | | | | |
| Bowtie2 Sensitivity? | very-sensitive-local | | | | | | |
| Additional Bowtie2 parameters? | | | | | | | |
| Alignment Quality? | perfect | | | | | | |

2 Start CaRpoools Report Generation

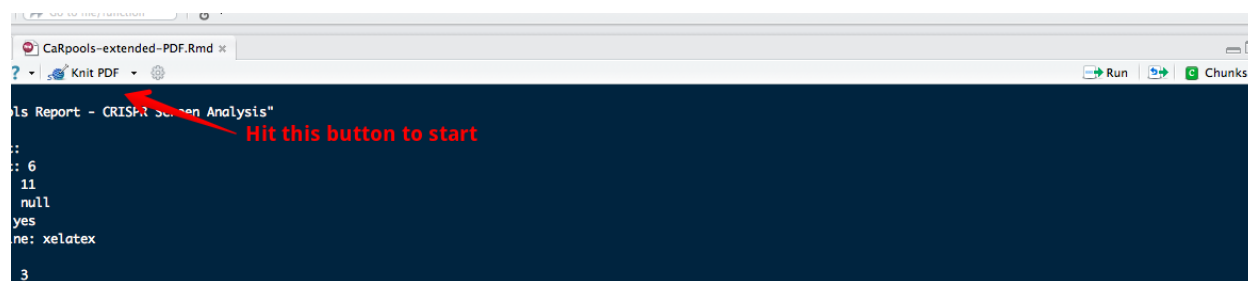
You can start caRpoools Report Generation after you did the following steps:

- Installed all required software and R packages (use `check.caRpoools(files=FALSE)` to verify)
- Put every file in the correct folder (MIACCS, data files, script files, Rmd templates)
- Put everything in the R working directory or set the working directory to the folder of your files
- Filled out the MIACCS file with all information, e.g. correct filenames, reference, data analysis options

You can check for all requirements by calling `check.caRpoools`.

2.1 Start CaRpoools using R-Studio

In the case you use R-Studio, you can start caRpoools by just opening the corresponding Rmd template file. At the top, you will find the **Knit PDF** or **Knit HTML** button, so you just need to press that and caRpoools will generate the report.



As an alternative, you can start caRpoools via `use.caRpoools` and provide additional parameters (see below).

2.2 Start CaRpools using R console

Moreover, caRpools report generation can also be initiated without R-studio installation, so that this can be done via R command line even on remote computers.

In this case, caRpools report generation can be started via `use.caRpools` with additional parameters, which are described below.

2.2.0.1 `use.caRpools()`

Usage:

use.caRpools(*type*=NULL, *file*="CaRpools-extended-PDF.Rmd", *miaccs*="MIACCS.xls", *check*=TRUE, *work.dir*=NULL)

type

Description If you provide a custom Rmd template that can generate both, PDF and HTML reports you can indicate which version you want to generate.

Default NULL

Values "PDF", "HTML"

file

Description The file name of your custom Rmd template file (with extension).

Default "CaRpools-extended-PDF.Rmd"

Values filename as character

miaccs

Description The filename of your MIACCS file.

Default "MIACCS.xls"

Values filename as character

check

Description Indicates whether caRpools will check for correct installation and file access.

Default TRUE

Values TRUE or FALSE (boolean)

work.dir

Description You can provide the absolute path to the working directory in which all files are placed (e.g. the MIACCS.xls and Rmd template).

Default NULL *Values* absolute path (character) or NULL if standard R working directory is used