

## Introduction to LTRpred

2020-10-27

Source: &lt;github&gt;:Introduction\_LTRpred [https://github.com/hajkd/LTRpred/blob/master/Introduction.md]

## Table of Contents

- 1. Introduction
- 2. Getting Started
  - 2.1 Installation
    - 2.1.1 LTRpred Docker Container
    - 2.1.2 Install tool dependencies on Linux
  - 2.2 Quick Start
  - 2.3 LTRpred output
  - 2.4 Import LTRpred output
  - 2.5 Output file format of LTRpred()
- 3. Detailed LTRpred run
  - 3.1 HMM Models
  - 3.2 Detailed description of adjustable LTRpred parameters
- 4. Metagenome scale annotations

## Introduction

The LTRpred package implements a software pipeline and provides an integrated workflow to screen for intact and potentially active LTR retrotransposons in any genomic sequence of interest. For this purpose, this package provides a rich set of analytics tools to allow researchers to quickly start annotating and explore their own genomes.

The *de novo* prediction of LTR transposons in LTRpred is based on the command line tools LTRharvest [http://www.zhu.uni-hamburg.de/~tdm/LTRharvest] and LTRdigest [http://www.zhu.uni-hamburg.de/infomung/infomung/sequence-genome/formats/kobayashi/digest.html] and extends these search strategies by additional analytics modules to filter the search space of putative LTR transposons for biologically meaningful candidates.

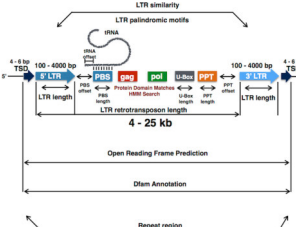
Please make sure that all command line tools are installed properly before running any LTRpred function.

## Getting Started

The rationale for implementing LTRpred was to implement an R based pipeline combining the most sensitive, accurate, and conservative state-of-the-art LTR retrotransposon detection tools and extend their inference by additional analyses and quality filtering steps to screen for functional and structurally intact elements. Hence, LTRpred aims to provide a high-level *de novo* prediction infrastructure to generate high quality annotations of intact LTR retrotransposons. All LTRpred functions are generic and their parameters can be changed to detect any form of LTR retrotransposon in any genome.

Internally, LTRpred is based on the *de novo* annotation tools LTRharvest and LTRdigest which use prior knowledge about DNA sequence features (also referred to as structure-based methods) such as the homology of Long Terminal Repeats (LTRs), Primer Binding Sites (PBS), gag and pol protein domains, and Target Site Duplications (TSDs) that are known to enable the process of transposition (Bergman and Quisenberry, 2007 [http://lib.elsevier.com/locate/S0966360307000000]) to infer LTR retrotransposons in any genome.

Hence, these *de novo* annotation tools are designed to screen the genome systematically and efficiently for these structural DNA features. Figure 1 shows the structural features of LTR retrotransposons that are used for predicting putative LTR transposons *de novo* in any genome of interest.



Sequence features of LTR retrotransposons. The structural element (LTR length, LTR similarity, TSD, max. size of full LTR retrotransposon, PBS binding, RNA binding, protein domain search, etc.) can be modified and controlled separately by corresponding arguments implemented in the LTRpred() function. In addition to controlling the structural features of candidate LTR retrotransposons, an open reading frame prediction, ORF annotation, copy number clustering, and LTR copy number

Based on the optimized output of these tools, the LTRpred package aims to provide researchers with maximum flexibility of adjustable parameters to detect any type of functional LTR retrotransposon. LTRpred package allows users to modify a vast range of parameters to screen for potential LTR transposons, so having this template shown in Figure 1 in mind will help researchers to modify structural parameters in a biologically meaningful manner.

## Installation

The fastest way to run LTRpred is to download the LTRpred Docker container which includes all pre-installed tool dependencies. For users who cannot use a Docker environment the individual installation instructions for each dependency tool are listed below (only for Linux).

## LTRpred Docker Container

Please be aware that the drosstlab/ltrpred container (current version) is suitable for Linux, macOS, and Windows users. Whereas the rstudio/ltrpred\_rstudio is not suitable for Windows users, because a port bridge cannot be established.

Please make sure to create a Docker [https://www.docker.com/get-started] account and to install Docker [https://docs.docker.com/install/] on your system.

## Download drosstlab/ltrpred container for use with R command line

```
# retrieve docker image from dockerhub
docker pull drosstlab/ltrpred
# run LTRpred container
docker run --rm -ti drosstlab/ltrpred
# start R prompt within LTRpred container
--/app# R
```

Users who wish to run the LTRpred Docker container in a console [https://docs.docker.com/get-started] environment can use the following approach based on UDocker [https://github.com/hajkd/LTRpred/issues/18] (Many thanks to Ilya Buznikov).

Within the LTRpred container R prompt run the LTRpred example:

```
LTRpred::LTRpred (./reference/LTRpred.html){genome.file = system.file (https://rdrr.io/r/base/system.file.html){"haploms_chr1.fa"}, package = "LTRpred"}
```

To exit R in the container run:

```
q (https://rdrr.io/r/base/qut.html){}
```

And to exit the LTRpred container run:

```
--/app# exit
```

Now, users can add their own genome data as well as the ORF database for further annotation to the LTRpred container by following these steps (in a different Terminal window):

```
# go to the folder path to which you want to
# store all genome and ORF data you want to
# mount to the LTRpred container and then run:
# create a new folder which will store
# all files that will be required in the
# LTRpred container
mkdir LTRpred_data
# create a ORF database folder
mkdir ORF
cd ORF
```

Now users can download and format the ORF database as follows (within the ORF folder created above). Unfortunately, the ORF database size is too large to make it part of the drosstlab/ltrpred container. In addition, the database is frequently updated and modified. Thus, it is recommended that users download and format the ORF database to their local hard drive and mount it to the running drosstlab/ltrpred container. To format the ORF database locally, users need to install HAMMER [http://hammer.cs.cmu.edu/download.html] on their local machine (see user manuals). However, within the drosstlab/ltrpred and drosstlab/ltrpred\_rstudio containers HAMMER is already pre-installed and does not need to be installed by the user. An example installation of HAMMER for Linux machines is listed below:

For each, users please install wget on your local machine using Homebrew [http://brew.sh/]:

```
wget https://raw.githubusercontent.com/rdrr/r/master/inst/doc/usingR.R
gunzip ORF.ham.gz
# format database by running hammer
hammer ORF.ham
cd ..
```

Next, make sure to also store the genome assembly file (in fasta format) you want to *de novo* annotate with LTRpred in the LTRpred\_data folder you just created.

A possible way to retrieve such a genome (in fasta format) using biomaRt [https://github.com/jharrison/biomaRt] if you use bioconductor please make sure to install all bioconductor package dependencies [https://github.com/jharrison/biomaRt#installation] before running the following code:

```
# install packages ("biomaRt")
biomaRt::getGenome (https://drosstlab.github.io/biomaRt/reference/getGenome.html){db = "ensembl",
  organism = "Saccharomyces cerevisiae",
  path = "yeast_genome",
  gunzip = TRUE}
```

The respective genome assembly file is now stored at yeast\_genome/Saccharomyces\_cerevisiae\_R64-1-1.dna.toplevel.fa and needs to be copied into the LTRpred\_data folder you just created.

Now users can mount the LTRpred\_data folder to the LTRpred Docker container (using the -v option). This -v mounting option is also available for the rstudio container version and can also be run within the RStudio Terminal:

```
docker run --rm -p 8787:8787 -v /path/to/your/path/to/LTRpred_data:/app/LTRpred_data -ti drosstlab/ltrpred
```

Within the LTRpred Docker container the LTRpred\_data folder is now stored in the working directory /app.

When running the LTRpred Docker container you should be able to see the mounted LTRpred\_data folder as following:

```
--/app# ls
```

```
bioRxivExtra_8-4-28.tar.gz
```

```
bioRxiv_data
```

```
bioRxiv_downloads
```

Next, users can again run the R prompt within the LTRpred Docker container to run LTRpred with the local data that was mounted:

```
--/app# R
```

```
# run LTRpred on the yeast genome that was mounted
# to the LTRpred container from your local "LTRpred_data" folder
LTRpred (./reference/LTRpred.html){genome.file = "LTRpred_data/yeast_genome/Saccharomyces_cerevisiae_R64-1-1.dna.toplevel.fa",
  genome = "ORF",
  ORF_db = "LTRpred_data/ORF",
  cores = 2}
```

As you can see, within the LTRpred container R prompt the current working directory is /app.

To also include the ORF database for further annotation users can specify the path to the ORF database:

```
# run LTRpred on the yeast genome that was mounted
# to the LTRpred container from your local "LTRpred_data" folder
LTRpred (./reference/LTRpred.html){genome.file = "LTRpred_data/yeast_genome/Saccharomyces_cerevisiae_R64-1-1.dna.toplevel.fa",
  genome = "ORF",
  ORF_db = "LTRpred_data/ORF",
  cores = 2}
```

Please be aware that using the ORF database for further annotation significantly increases the computation time of the LTRpred pipeline.

## Retrieve LTRpred output files from Docker container

Next, users can retrieve the LTRpred generated results from the docker container by opening a second terminal window (while the drosstlab/ltrpred container remains running) and perform the following steps:

1. Close R in the running docker container using ctrl + c.
2. This should bring you back to the docker bash prompt: root@drosstlab/ltrpred:~#.
3. Copy the docker ID, in this case ac389c7a88f (this docker ID is just an example, please use the docker ID shown on your system).
4. Type in the newly opened Terminal window:

```
# copy haploms_chr1.LTRpred output from docker container to hard drive
docker cp ac389c7a88f:/app/haploms_chr1.LTRpred path/to/your/host/hard/drive/haploms_chr1.LTRpred
```

This example assumes that you ran the example LTRpred::LTRpred (genome.file = system.file("haploms\_chr1.fa", package = "LTRpred")) which created the output folder haploms\_chr1.LTRpred in the docker folder /app/.

Please note, that if you specify different file paths when creating files within the docker container, these must be adjusted when running:

```
# copy files from docker container to hard drive
docker cp ac389c7a88f:/app/your/inside/docker/path/here/haploms_chr1.LTRpred path/to/your/host/hard/drive/haploms_chr1.LTRpred
```

## Download drosstlab/ltrpred\_rstudio container for use with RStudio Server

```
# retrieve docker image from dockerhub
docker pull drosstlab/ltrpred_rstudio
# run LTRpred container
docker run -e RStudioID=LTRpred --rm -p 8787:8787 -ti drosstlab/ltrpred_rstudio
```

To open RStudio and interact with the container go to your standard web browser and type in the following url:

```
http://localhost:8787
```

Username: rstudio Password: LTRpred

Within RStudio you can now run the example:

```
LTRpred::LTRpred (./reference/LTRpred.html){genome.file = system.file (https://rdrr.io/r/base/system.file.html){"haploms_chr1.fa"}, package = "LTRpred"}
```

Users can exit the container by pressing ctrl + c multiple times.

Next, users can mount their LTRpred\_data folder to the RStudio server run the same way they mounted folders in the command line container version (using -v). This folder mounting can also be run within the RStudio Terminal of the drosstlab/ltrpred\_rstudio container:

```
# retrieve docker image from dockerhub
docker pull drosstlab/ltrpred_rstudio
# run LTRpred container
docker run -e RStudioID=LTRpred --rm -p 8787:8787 -v /path/to/your/path/to/LTRpred_data/home/rstudio/LTRpred_data -ti drosstlab/ltrpred_rstudio
```

Now go to your standard web browser and type in the following url:

```
http://localhost:8787
```

Username: rstudio Password: LTRpred

In RStudio type:

Users can *exit* the container by pressing `Ctrl + c` multiple times.

1. Close it in the running docker container using `q()`.
2. This should bring you back to the docker bash: `root@ec388ef5a880:/app#`.
3. Copy the docker ID, in this case `a388ef5a880` (this docker ID is just an example, please use the docker ID shown on your system).
4. Type in the newly opened Terminal window:

```
# copy hisapiens_chrV_itrprcd output from docker container to hard drive
docker cp ac1809ffa0080:/app/hisapiens_chrV_itrprcd path/to/your/host/hard/disk/hisapiens_chrV_itrprcd
```

This example assumes that you ran the example `LTpred_Rsh.LTpred(genome.file = system.file("hsapiens_chrV.fa", package = "LTpred"))` which created the output folder `hsapiens_chrV_ltpred` in the docker folder `/app`. Please note that if you specify different file paths when creating files within the docker container these must be adjusted when running

```
# copy files from docker container to hard drive
docker cp ac309gf1a0089:/app/your/inside/docker/path/here/Hisapiens_ChRV_itrprd path/to/your/host/hard/disk/Hisapiens_ChRV_itrprd
```

Please read more details about how to transfer genome files and the Dfam database in the previous section Download Itrprd container for use with R command line

## Programming Languages

Please make sure that the following programming languages are installed on your system

- Perl (<https://www.perl.org/get.html>)
- Ruby (<https://www.ruby-lang.org/en/documentation/installation/>)
- Python (<https://www.python.org/downloads/>)
- C/C++ (<https://www.oplusplus.com/>)
- R (<http://lib.stat.cmu.edu/R/CRAN/>)

## Install Programming languages and Linux Tools

```
apt-get update \
$S apt-get -y install apt-utils \
$S apt-get -y install gcc \
$S apt-get -y install python\ \
$S apt-get -y install perl \
$S apt-get -y install make \
$S apt-get -y install vim \
$S apt-get -y install wget \
$S apt-get -y install git \
$S apt-get -y install git-svn \
$S apt-get -y install gitk \
$S apt-get -y install libssl-dev \
$S apt-get -y install build-essential \
$S apt-get -y install libreadline-dev \
$S apt-get -y install libncurses-dev \
$S apt-get -y install libpq-dev \
$S apt-get -y install software-properties-common
```

Install HMMER (<http://hmmer.org/download.html>)

A detailed description of how to install `HEERS` for several operating systems can be found here (<http://hmmmer.crcidownload.html>)

```
mkdir software_downloads \
  cd software_downloads \
  wget http://fedy1ab.org/software/hmmer/hmmer-3.2.tar.gz \
  tar xf hmmer-3.2.tar.gz \
  cd hmmer-3.2 \
  ./configure \
  make \
  make check \
  sudo make install \
  cd ..
```

Install USEARCH (<http://drive5.com/usearch/download.html>)

A detailed description of how to install `ustarack` for several operating systems can be found here (<http://drive5.com/search/manual/install.html>).

First, users will need to register and download USEARCH for their operating system from <http://drive5.com/usearch/download.html> (<http://drive5.com/usearch/download.html>)

After downloading USEARCH you will need to install it as a command line tool in your `/usr/local/` directory and you should be able to execute the following command in your Terminal

```
cd software_downloads \
$ wget https://www.drive5.com/downloads/usagearch-v1.0.667.linux32.gz
$ gunzip usagearch1.0.667.linux32.gz \
$ chmod +x usagearch1.0.667.linux32 \
$ sudo mv usagearch1.0.667.linux32 usagearch \
$ sudo cp usagearch /usr/local/bin/usagearch \
$ cd ..
usagearch -version

usagearch v1.0.667.linux32
```

Install VSEARCH (<https://github.com/lorognes/vsearch>)

A detailed description of how to install VSEARCH for several operating systems can be found here (<https://github.com/torognesi/vsearch>).

Please install git (<https://git-scm.com/book/en/v2/Getting-Started-Installing-Git>) before running the following commands:

```
cd software_downloads \
## wget https://github.com/torognex/vsearch/archive/v2.14.2.tar.gz
## tar xzf v2.14.2.tar.gz \
## cd vsearch-2.14.2 \
## sudo ./autogen.sh \
## sudo ./configure \
## sudo make \
## sudo make install \
## cd
```

Install dfamscan.pl ([http://www.dfam.org/web\\_download/Current\\_Release/dfamscan.pl](http://www.dfam.org/web_download/Current_Release/dfamscan.pl))

dfamscan.pl ([https://www.dfam.org/releases/Dfam\\_3.1/infrastructure/dfamscan.pl.gz](https://www.dfam.org/releases/Dfam_3.1/infrastructure/dfamscan.pl.gz)) needs to be unzipped and stored at `/usr/local/bin/dfamscan.pl` and executable via `perl /usr/local/bin/dfamscan.pl -help`. This is important to be able to run the hmmer search against the Dfam database

```
cd software_downloads \
## wget https://www.dfam.org/releases/Dfam_3.2/infrastructure/dfamscan.pl.gz \
## gunzip dfamscan.pl.gz \
## sudo cp dfamscan.pl /usr/local/bin/dfamscan.pl \
## -d
```

```
wget https://www.dfan.org/releases/Dfan_3.1/families/Dfan.hmm.gz
gunzip Dfan.hmm.gz
# run hmppress
hmppress Dfan.hmm
```

The path to the folder where the formatted `orac` database can be found can then be passed as argument `Orac.db` to the `LTPred::LTPred()` function

## Install R packages

Please make sure that Bioconductor (<https://www.bioconductor.org/install/>) and all package dependencies are installed on the system on which you would like to run `LTfpred`.

### Install prerequisite CRAN and Bioconductor packages

```
install.packages("devtools")
install.packages("tidyverse")
install.packages("BioParas")
library(devtools)
library(BioParas)
install("https://raw.githubusercontent.com/rtackai/GenomicFeatures/0.99.0/BioParas/inst/extdata/trackviz")
install.packages(c("tidyverse", "data.table", "magrittr", "biomart", "r", "png", "dplyr", "devtools"))
devtools::install("https://raw.githubusercontent.com/rtackai/GenomicFeatures/0.99.0/BioParas/inst/extdata/trackviz")
install.packages(c("RGG", "ggrepel", "ggfortify"))
https://cran.r-project.org/web/packages/trackviz/trackviz_0.6-20.tar.gz
install.packages("trackviz_0.6-20.tar.gz", type = "source")
install.packages("survival")
```

Now users may install LTRpred as follows:

```
# install.packages("devtools")
devtools::install_github("https://devtools.r-lib.org/reference/remote-reexports.html")("HajkQ/L78pred")
```

## Quick Start

The fastest way to generate a LTR retrotransposon prediction for a genome of interest (after installing <https://github.com/LTRpred/articles/Introduction.html#installation>) all prerequisite command line tools) is to use the `LTRpred()` function and relying on the default parameters. In the following example, a LTR transposon prediction is performed for parts of the Human Y chromosome

```
# load LTPred package
library (https://rdrr.io/r/base/library.html){LTPred (https://github.com/WajKD/LTPred)}
# de novo LTR transposon prediction for the Human V chromosome
LTPred (.../reference/LTPred.html){
  genome.file = system.file (https://rdrr.io/r/base/system.file.html){("HsapL1n_Chrv.fa", package = "LTPred")},
  cores = 4
}
```

```
Running LTRpred on genome ~/Library/Frameworks/Python.frameworks/Versions/3.6/Resources/Library/LTRpred/Haploids_Chry.fa' with 4 core(s) and searching for retrotransposons using the overlaps option (overlaps = "no") ...

No bow files were specified, thus the internal BWA library will be used! See ~/Library/Frameworks/Python.frameworks/Versions/3.6/Resources/Library/LTRpred/BWA/bwa.exe" for details.
No BWA files were specified, thus the internal BWA library will be used! See ~/Library/Frameworks/Python.frameworks/Versions/3.6/Resources/Library/LTRpred/BWA/BWA_library.fa' for details.
The output folder 'Haploids_Chry_LTRpred' seems to exist already and will be used to store LTRpred results ...

LTRpred - Step 1:
Run LTRharvest...
LTRharvest: Generating index file Haploids_Chry_LTRharvest/Haploids_Chry_index.fa with gt suffixeraster...
Running LTRharvest and writing results to Haploids_Chry_LTRharvest...
LTRharvest analysis finished!

LTRpred - Step 2:
Generating index file Haploids_Chry_LTRidgest/Haploids_Chry_index_LTRidgest.fa with suffixeraster...
LTRidgest: Sort index file...
Running LTRidgest and writing results to Haploids_Chry_LTRidgest...
LTRidgest analysis finished!

LTRpred - Step 3:
Import LTRidgest Predictions...

Input: Haploids_Chry_LTRidgest/Haploids_Chry_LTRidgestPrediction.gff -> Row Number: 170
Remove "NA" -> Row Row Number: 170
(LTR) Filtering For repeat regions has been finished.
(LTR) Filtering For LTR retrotransposons has been finished.
(LTR) Filtering For inverted repeats has been finished.
(LTR) Filtering For LTRs has been finished.
(LTR) Filtering for target site duplication has been finished.
(LTR) Filtering for primer binding site has been finished.
(LTR) Filtering for protein match has been finished.
(LTR) Filtering for 4R tract has been finished.

LTRpred - Step 4:
Perform ODF Prediction using "search -fasta_Finder" ...
search -v 1.1861_1861212_4.400 NA6 (17-206 total), 8 cores
(C) Copyright 2013-15 Robert C. Edgar, all rights reserved.
http://bratn.com/search

00:00 1.00% 100.0% working
Data ODF Prediction table: row(df) = 24 candidates.
unique(ID) = 24 candidates.
unique(orf.id) = 24 candidates.

LTRpred - Step 5:
Perform methylation context quantification...
Join methylation context CG, CHG, CHH, CHG) count table: row(df) = 24 candidates.
unique(ID) = 24 candidates.
unique(orf.id) = 24 candidates.
Copy files to result folder 'Haploids_Chry_LTRpred'.

LTRpred - Step 6:
Starting retrotransposon evolutionary age estimation by comparing the 3' and 5' LTRs using the molecular evolution model "EB" and the mutation rate "1.3e-07" (please make sure the mutation rate can be assumed for your species of interest!) for 24 predicted elements ...

Please be aware that evolutionary age estimation based on 3' and 5' LTR comparisons are only very rough time estimates and don't take reverse transcription mediated retrotransposon recombination between family members of retrotransposons into account! Please consult Sanchez et al., 2017 Nature Communications and Drost & Sencher, 2019 Genome Biology and Evolution for more details on retrotransposon recombination.

LTRpred - Step 7:
The LTRpred prediction table has been filtered (default) to remove potential false positive. Predicted LTRs must have an PDS or Protein Domain and must fulfill thresholds: min = 70%, max = 70%. Furthermore, LTRs having more than 100 of N's in their sequence have also been removed.
Input #file: 24
Output #file: 21

LTRpred finished all analyses successfully. All output files were stored at 'Haploids_Chry_LTRpred'.
[1] "Successful job 1."
```

LTRpred output

The LTRpred() function internally generates a folder named "LTRpred" which stores all output annotation and sequence files. In detail, the following files and folders are generated by the LTRpred() function:

- **Folder "LTRpred"**
  - **\*\_odf\_prediction\_\*.fa** : Stores the predicted open reading frames within the predicted LTR transposons as DNA sequence.
  - **\*\_odf\_prediction\_aa\_\*.fa** : Stores the predicted open reading frames within the predicted LTR transposons as protein sequence.
  - **\*\_LTRpred.gff** : Stores the LTRpred predicted LTR transposons in GFF format.
  - **\*\_LTRpred.bed** : Stores the LTRpred predicted LTR transposons in BED format.
  - **\*\_LTRpred\_DataSheet.csv** : Stores the output table as data sheet.
  - **Folder "LTRharvest"**
    - **"LTRharvest"/BetweenTSSeqs\_\*.fa** : DNA sequences of the region between the LTRs in fasta format.
    - **"LTRharvest"/Details.txt** : A spread sheet containing detailed information about the predicted LTRs.
    - **"LTRharvest"/FullLTRRetrotransposonSeqs\_\*.fa** : DNA sequences of the entire predicted LTR retrotransposons.
    - **"LTRharvest"/Index\_\*.fa** : The sufficiency index file used to predict putative LTR retrotransposons.
    - **"LTRharvest"/Prediction.gff** : A spread sheet containing detailed additional information about the predicted LTRs (partially redundant with the "Details.txt" file).
    - **"LTRharvest"/Index\_LTRidgest\_\*.fa** : The sufficiency index file used to predict putative LTR retrotransposons with LTRidgest.
  - **Folder "LTRidgest"**
    - **"LTRidgest"/LTRidgestPrediction.gff** : A spread sheet containing detailed information about the predicted LTRs.
    - **"LTRidgest"/LTRidgest\_about.csv** : A spread sheet containing additional detailed information about the predicted LTRs.
    - **"LTRidgest"/LTRidgest\_complete\_\*.fa** : The full length DNA sequences of all predicted LTR transposons.
    - **"LTRidgest"/LTRidgest\_conditions.csv** : Contains information about the parameters used for a given LTRidgest run.
    - **"LTRidgest"/LTRidgest\_pds\_\*.fa** : Stores the predicted PDS sequences for the putative LTR retrotransposons.
    - **"LTRidgest"/LTRidgest\_ppt\_\*.fa** : Stores the predicted PPT sequences for the putative LTR retrotransposons.
    - **"LTRidgest"/LTRidgest\_ltr\_\*.fa** and **"LTRidgest\_ltr\_\*.fa** : Stores the predicted 5' and 3' LTR sequences. Note: If the direction of the putative retrotransposon could be predicted, these files will contain the corresponding 3' and 5' LTR sequences. If no direction could be predicted, forward direction with regard to the original sequence will be assumed by LTRidgest, i.e. the 'left' LTR will be considered the 5' LTR.
    - **"LTRidgest"/LTRidgest\_pdm\_domains\_\*.fa** : Stores the DNA sequences of the HMM matches to the LTR retrotransposon candidates.
    - **"LTRidgest"/LTRidgest\_pdm\_domains\_aa\_\*.fa** : Stores the concatenated protein sequences of the HMM matches to the LTR retrotransposon candidates.
    - **"LTRidgest"/LTRidgest\_pdm\_domains\_aa\_all\_\*.fa** : Stores the alignment information for all matches of the given protein domain model to the translations of all candidates.

Import LTRpred output

```
The LTRpred() output table "LTRpred_DataSheet.csv" is in Sdy (https://hajkd.github.io/ncbi/sdy-data.html) format and can then be imported using read.LTRpred(). The sdy output format is designed to work seamlessly with the tidyverse (https://www.tidyverse.org/) and R data science (http://h4ds.had.co.nz/) framework.

# Import LTRpred prediction output
Haploids_Chry <- read.LTRpred("~/reference/read.LTRpred.html")\Haploids_Chry_LTRpred/Haploids_Chry_LTRpred_DataSheet.csv")
# Look at some results
dplyr::select (https://dplyr.tidyverse.org/reference/select.html)\Haploids_Chry", ltr_similarity:and, vRNA_posit, Clust_cn)

# A tibble: 24 x 9
  ltr_similarity similarity protein_domain orfs chromosome
  <dbl> <dbl> <chr> <chr> <chr>
1 89.73 (89,82) RVT_3 1 NC008024.3|Hemera
2 89.85 (89,96) RVT_3 1 NC008024.3|Hemera
3 79.71 (79,86) cdo 8 NC008024.3|Hemera
4 86.63 (84,84) RVT_3 8 NC008024.3|Hemera
5 75.52 cdo RVT_3 8 NC008024.3|Hemera
6 76.47 (76,76) RVT_3 1 NC008024.3|Hemera
7 88.28 (88,82) cdo 8 NC008024.3|Hemera
8 76.47 (76,76) RVT_3 8 NC008024.3|Hemera
9 89.53 (89,96) RVT_3 8 NC008024.3|Hemera
10 82.55 (82,84) RVT_3 1 NC008024.3|Hemera
11 82.35 (82,84) RVT_3 2 NC008024.3|Hemera
12 79.51 (79,86) RVT_3 8 NC008024.3|Hemera
13 78.42 (78,86) RVT_3 1 NC008024.3|Hemera
14 82.71 (82,82) RVT_3 1 NC008024.3|Hemera
# ... with 4 more variables: start_cdist, and_cdist, vRNA_posit<chr>,
# Clust_cn<chr>

Looking at all columns:
dplyr::glimpse (https://tidyverse.org/reference/glimpse.html)\Haploids_Chry")
```

```
library (https://rdrr.io/r/src/library.html)[?lqpred (https://github.com/ma3d/lqpred)]
# de novo LTV transcription prediction of 'G. simulans'
LQPred ("./reference/lqpred.html",
  genome.file = "Sereveling_genome",
  tssas = paste0 ("https://rdrr.io/r/src/library.html")[?HMMs"], package = "LQPred", "sacCer-HMMs.R",
  hms = paste0 ("https://rdrr.io/r/src/library.html")[?HMMs"], package = "LQPred", "hms.R",
  cluster = TRUE,
  client.size = 0.5,
  copy.number.est = TRUE,
  cores = 4
)
```

# Introduction to LTRpred

```
Running LTRpred on genome "_nckl_download/genomes/Saccharomyces_cerevisiae_genomic_refseq.fasta.gz" with 4 cores(x) and searching for retrotransposons using the overlaps option (overlaps = "nr") ...

The output folder "Saccharomyces_cerevisiae_genomic_refseq_ltrpred" does not seem to exist yet and will be created ...

LTRpred - Step 1:
Run LTRharvest...
LTRharvest: Generating index file Saccharomyces_cerevisiae_genomic_refseq_ltrharvest/Saccharomyces_cerevisiae_genomic_refseq_index.fasta with gt suffixseparator...
Running LTRharvest and writing results to Saccharomyces_cerevisiae_genomic_refseq_ltrharvest...
LTRharvest analysis finished!

LTRpred - Step 2:
Generating index file Saccharomyces_cerevisiae_genomic_refseq_ltrdigest/Saccharomyces_cerevisiae_genomic_refseq_index_ltrdigest.fasta with suffixseparator...
LTRdigest: Sort index file...
Running LTRdigest and writing results to Saccharomyces_cerevisiae_genomic_refseq_ltrdigest...
LTRdigest analysis finished!

LTRpred - Step 3:
Import LTRdigest Predictions...

Input: Saccharomyces_cerevisiae_genomic_refseq_ltrdigest/Saccharomyces_cerevisiae_genomic_refseq_ltrdigestPrediction.gtf -> Row Number: 283
Remove NA's -> Row Number: 283
(L1R) Filtering for repeat regions has been finished.
(L1R) Filtering for LTR retrotransposons has been finished.
(L1R) Filtering for Inverted repeats has been finished.
(L1R) Filtering for LTRs has been finished.
(L1R) Filtering for target site duplication has been finished.
(L1R) Filtering for primer binding site has been finished.
(L1R) Filtering for protein match has been finished.
(L1R) Filtering for OR tract has been finished.

LTRpred - Step 4:
Perform OR Prediction using "blastx -dust,fasta" ...
blastx v2.10.0+100m211: 4.400 MB00 MB00 (37.20k total), 8 cores
(C) Copyright 2003-15 Robert C. Edgar, all rights reserved.
http://blast.ncbi.nlm.nih.gov/Blast.cgi

00:00 2.20k 100.0% Working
Join OR Prediction table: row(df) = 36 candidates.
unique(df) = 36 candidates.
unique(df$id) = 36 candidates.
Perform clustering of similar LTR transposons using "search --cluster_Fast" ...
Running CLUSTALW with 36 as sequence similarity threshold using 4 cores ...
Reading file (search/clustalw/cluster/Fastq/seqs/uniqueLTRpred/Saccharomyces_cerevisiae_genomic_refseq_ltrdigest/Saccharomyces_cerevisiae_genomic_refseq_ltrdigest_complete.fasta 100k
201812 at 10 36 seqs, min 100, max 20180, avg 4004
Sorting by length 1000
Counting unique k-mers 1000
Clustering 1000
Sorting clusters 1000
Writing clusters 1000
Clusters: 10 size min 1, max 10, avg 3.6
Singletons: 7, 10.4% of seqs, 70.4% of clusters
Sorting clusters by abundance 1000
blastx v2.10.0+100m211: 4.400 MB00 MB00, 8 cores
https://github.com/timgates/search

CLUSTALW output has been stored in: Saccharomyces_cerevisiae_genomic_refseq_ltrpred
Join Cluster table: row(df) = 36 candidates.
unique(df) = 36 candidates.
unique(df$id) = 36 candidates.
Join Cluster Copy Number table: row(df) = 36 candidates.
unique(df) = 36 candidates.
unique(df$id) = 36 candidates.

LTRpred - Step 5:
Perform methylation context quantification...
Join methylation context (CG, CHG, CHH, CH2) count table: row(df) = 36 candidates.
unique(df) = 36 candidates.
unique(df$id) = 36 candidates.
Copy files to result folder "Saccharomyces_cerevisiae_genomic_refseq_ltrpred".

LTRpred - Step 6:
Starting retrotransposon evolutionary age estimation by comparing the 3' and 5' LTRs using the molecular evolution model "R8B" and the mutation rate "1.3e-07" (please make sure the mutation rate can be assumed for your species of interest!) for 36 predicted elements ...

Please be aware that evolutionary age estimation based on 3' and 5' LTR comparisons are only very rough time estimates and don't take reverse-transcription mediated retrotransposon recombination between family members of retroelements into account! Please consult Sanchez et al., 2007 Nature Communications and Drost & Sanchez, 2009 Genome Biology and Evolution for more details on retrotransposon recombination.

LTRpred - Step 7:
The LTRpred prediction table has been filtered (default) to remove potential false positives. Predicted LTRs must have an PDS or Protein Domain and must fulfill thresholds: min = 70%; kmerfs = 0. Furthermore, LTRs having more than 10% of N's in their sequence have also been removed.

Input files: 10
Output files: 11
Perform solo LTR Copy Number Estimation...
Run subelement of the genome assembly...

Building a new OR, current time: 02/24/2020 16:58:53
New OR name: "_nckl_download/genomes/Saccharomyces_cerevisiae_genomic_refseq.fasta
New OR title: "_nckl_download/genomes/Saccharomyces_cerevisiae_genomic_refseq.fasta
Sequence type: Not Initiated
Keep isoforms: 1
Maximum file size: 1000000000
Adding sequences from FASTA added 37 sequences in 0.120386 seconds.
Perform BLAST searches of 3' prime LTRs against genome assembly...
Perform BLAST searches of 5' prime LTRs against genome assembly...
Import BLAST results...
Filter hits results...
Estimate OR for each LTR sequence...
Finished LTR OR estimation...

LTRpred finished all analyses successfully. All output files were stored at "Saccharomyces_cerevisiae_genomic_refseq_ltrpred".
[1] "Successful job 1."

Warning message:
The LTR copy number estimation returned an empty file. This suggests that there were no solo LTRs found in the input genome sequence.

The output can then be imported using:

# Import LTRpred output for S. cerevisiae
file <- file.path (https://ndr.io/hajk/file-path.html)"Saccharomyces_cerevisiae_genomic_refseq_ltrpred",
               "Saccharomyces_cerevisiae_genomic_refseq_ltrpred.table.tsv")
Sacreviasia_LTRpred <- readLTRpred (.../reference/readLTRpred.html)(file)
# Look at output
dplyr::glimpse (https://tbls.io/tdgpcwera.org/reference/glimpse.html)(Sacreviasia_LTRpred)

Observations: 31
Variables: 82
# chr1
# chr2
# chr3 chr4 chr5 chr6 chr7 chr8 chr9 chr10 chr11 chr12 chr13 chr14 chr15 chr16 chr17 chr18 chr19 chr20 chr21 chr22 chr23 chr24 chr25 chr26 chr27 chr28 chr29 chr30 chr31 chr32 chr33 chr34 chr35 chr36 chr37 chr38 chr39 chr40 chr41 chr42 chr43 chr44 chr45 chr46 chr47 chr48 chr49 chr50 chr51 chr52 chr53 chr54 chr55 chr56 chr57 chr58 chr59 chr60 chr61 chr62 chr63 chr64 chr65 chr66 chr67 chr68 chr69 chr70 chr71 chr72 chr73 chr74 chr75 chr76 chr77 chr78 chr79 chr80 chr81 chr82 chr83 chr84 chr85 chr86 chr87 chr88 chr89 chr90 chr91 chr92 chr93 chr94 chr95 chr96 chr97 chr98 chr99 chr100 chr101 chr102 chr103 chr104 chr105 chr106 chr107 chr108 chr109 chr110 chr111 chr112 chr113 chr114 chr115 chr116 chr117 chr118 chr119 chr120 chr121 chr122 chr123 chr124 chr125 chr126 chr127 chr128 chr129 chr130 chr131 chr132 chr133 chr134 chr135 chr136 chr137 chr138 chr139 chr140 chr141 chr142 chr143 chr144 chr145 chr146 chr147 chr148 chr149 chr150 chr151 chr152 chr153 chr154 chr155 chr156 chr157 chr158 chr159 chr160 chr161 chr162 chr163 chr164 chr165 chr166 chr167 chr168 chr169 chr170 chr171 chr172 chr173 chr174 chr175 chr176 chr177 chr178 chr179 chr180 chr181 chr182 chr183 chr184 chr185 chr186 chr187 chr188 chr189 chr190 chr191 chr192 chr193 chr194 chr195 chr196 chr197 chr198 chr199 chr200 chr201 chr202 chr203 chr204 chr205 chr206 chr207 chr208 chr209 chr210 chr211 chr212 chr213 chr214 chr215 chr216 chr217 chr218 chr219 chr220 chr221 chr222 chr223 chr224 chr225 chr226 chr227 chr228 chr229 chr230 chr231 chr232 chr233 chr234 chr235 chr236 chr237 chr238 chr239 chr240 chr241 chr242 chr243 chr244 chr245 chr246 chr247 chr248 chr249 chr250 chr251 chr252 chr253 chr254 chr255 chr256 chr257 chr258 chr259 chr260 chr261 chr262 chr263 chr264 chr265 chr266 chr267 chr268 chr269 chr270 chr271 chr272 chr273 chr274 chr275 chr276 chr277 chr278 chr279 chr280 chr281 chr282 chr283 chr284 chr285 chr286 chr287 chr288 chr289 chr290 chr291 chr292 chr293 chr294 chr295 chr296 chr297 chr298 chr299 chr300 chr301 chr302 chr303 chr304 chr305 chr306 chr307 chr308 chr309 chr310 chr311 chr312 chr313 chr314 chr315 chr316 chr317 chr318 chr319 chr320 chr321 chr322 chr323 chr324 chr325 chr326 chr327 chr328 chr329 chr330 chr331 chr332 chr333 chr334 chr335 chr336 chr337 chr338 chr339 chr340 chr341 chr342 chr343 chr344 chr345 chr346 chr347 chr348 chr349 chr350 chr351 chr352 chr353 chr354 chr355 chr356 chr357 chr358 chr359 chr360 chr361 chr362 chr363 chr364 chr365 chr366 chr367 chr368 chr369 chr370 chr371 chr372 chr373 chr374 chr375 chr376 chr377 chr378 chr379 chr380 chr381 chr382 chr383 chr384 chr385 chr386 chr387 chr388 chr389 chr390 chr391 chr392 chr393 chr394 chr395 chr396 chr397 chr398 chr399 chr400 chr401 chr402 chr403 chr404 chr405 chr406 chr407 chr408 chr409 chr410 chr411 chr412 chr413 chr414 chr415 chr416 chr417 chr418 chr419 chr420 chr421 chr422 chr423 chr424 chr425 chr426 chr427 chr428 chr429 chr430 chr431 chr432 chr433 chr434 chr435 chr436 chr437 chr438 chr439 chr440 chr441 chr442 chr443 chr444 chr445 chr446 chr447 chr448 chr449 chr450 chr451 chr452 chr453 chr454 chr455 chr456 chr457 chr458 chr459 chr460 chr461 chr462 chr463 chr464 chr465 chr466 chr467 chr468 chr469 chr470 chr471 chr472 chr473 chr474 chr475 chr476 chr477 chr478 chr479 chr480 chr481 chr482 chr483 chr484 chr485 chr486 chr487 chr488 chr489 chr490 chr491 chr492 chr493 chr494 chr495 chr496 chr497 chr498 chr499 chr500 chr501 chr502 chr503 chr504 chr505 chr506 chr507 chr508 chr509 chr510 chr511 chr512 chr513 chr514 chr515 chr516 chr517 chr518 chr519 chr520 chr521 chr522 chr523 chr524 chr525 chr526 chr527 chr528 chr529 chr530 chr531 chr532 chr533 chr534 chr535 chr536 chr537 chr538 chr539 chr540 chr541 chr542 chr543 chr544 chr545 chr546 chr547 chr548 chr549 chr550 chr551 chr552 chr553 chr554 chr555 chr556 chr557 chr558 chr559 chr560 chr561 chr562 chr563 chr564 chr565 chr566 chr567 chr568 chr569 chr570 chr571 chr572 chr573 chr574 chr575 chr576 chr577 chr578 chr579 chr580 chr581 chr582 chr583 chr584 chr585 chr586 chr587 chr588 chr589 chr590 chr591 chr592 chr593 chr594 chr595 chr596 chr597 chr598 chr599 chr600 chr601 chr602 chr603 chr604 chr605 chr606 chr607 chr608 chr609 chr610 chr611 chr612 chr613 chr614 chr615 chr616 chr617 chr618 chr619 chr620 chr621 chr622 chr623 chr624 chr625 chr626 chr627 chr628 chr629 chr630 chr631 chr632 chr633 chr634 chr635 chr636 chr637 chr638 chr639 chr640 chr641 chr642 chr643 chr644 chr645 chr646 chr647 chr648 chr649 chr650 chr651 chr652 chr653 chr654 chr655 chr656 chr657 chr658 chr659 chr660 chr661 chr662 chr663 chr664 chr665 chr666 chr667 chr668 chr669 chr670 chr671 chr672 chr673 chr674 chr675 chr676 chr677 chr678 chr679 chr680 chr681 chr682 chr683 chr684 chr685 chr686 chr687 chr688 chr689 chr690 chr691 chr692 chr693 chr694 chr695 chr696 chr697 chr698 chr699 chr700 chr701 chr702 chr703 chr704 chr705 chr706 chr707 chr708 chr709 chr710 chr711 chr712 chr713 chr714 chr715 chr716 chr717 chr718 chr719 chr720 chr721 chr722 chr723 chr724 chr725 chr726 chr727 chr728 chr729 chr730 chr731 chr732 chr733 chr734 chr735 chr736 chr737 chr738 chr739 chr740 chr741 chr742 chr743 chr744 chr745 chr746 chr747 chr748 chr749 chr750 chr751 chr752 chr753 chr754 chr755 chr756 chr757 chr758 chr759 chr760 chr761 chr762 chr763 chr764 chr765 chr766 chr767 chr768 chr769 chr770 chr771 chr772 chr773 chr774 chr775 chr776 chr777 chr778 chr779 chr780 chr781 chr782 chr783 chr784 chr785 chr786 chr787 chr788 chr789 chr790 chr791 chr792 chr793 chr794 chr795 chr796 chr797 chr798 chr799 chr800 chr801 chr802 chr803 chr804 chr805 chr806 chr807 chr808 chr809 chr810 chr811 chr812 chr813 chr814 chr815 chr816 chr817 chr818 chr819 chr820 chr821 chr822 chr823 chr824 chr825 chr826 chr827 chr828 chr829 chr830 chr831 chr832 chr833 chr834 chr835 chr836 chr837 chr838 chr839 chr840 chr841 chr842 chr843 chr844 chr845 chr846 chr847 chr848 chr849 chr850 chr851 chr852 chr853 chr854 chr855 chr856 chr857 chr858 chr859 chr860 chr861 chr862 chr863 chr864 chr865 chr866 chr867 chr868 chr869 chr870 chr871 chr872 chr873 chr874 chr875 chr876 chr877 chr878 chr879 chr880 chr881 chr882 chr883 chr884 chr885 chr886 chr887 chr888 chr889 chr890 chr891 chr892 chr893 chr894 chr895 chr896 chr897 chr898 chr899 chr900 chr901 chr902 chr903 chr904 chr905 chr906 chr907 chr908 chr909 chr910 chr911 chr912 chr913 chr914 chr915 chr916 chr917 chr918 chr919 chr920 chr921 chr922 chr923 chr924 chr925 chr926 chr927 chr928 chr929 chr930 chr931 chr932 chr933 chr934 chr935 chr936 chr937 chr938 chr939 chr940 chr941 chr942 chr943 chr944 chr945 chr946 chr947 chr948 chr949 chr950 chr951 chr952 chr953 chr954 chr955 chr956 chr957 chr958 chr959 chr960 chr961 chr962 chr963 chr964 chr965 chr966 chr967 chr968 chr969 chr970 chr971 chr972 chr973 chr974 chr975 chr976 chr977 chr978 chr979 chr980 chr981 chr982 chr983 chr984 chr985 chr986 chr987 chr988 chr989 chr990 chr991 chr992 chr993 chr994 chr995 chr996 chr997 chr998 chr999 chr1000 chr1001 chr1002 chr1003 chr1004 chr1005 chr1006 chr1007 chr1008 chr1009 chr1010 chr1011 chr1012 chr1013 chr1014 chr1015 chr1016 chr1017 chr1018 chr1019 chr1020 chr1021 chr1022 chr1023 chr1024 chr1025 chr1026 chr1027 chr1028 chr1029 chr1030 chr1031 chr1032 chr1033 chr1034 chr1035 chr1036 chr1037 chr1038 chr1039 chr1040 chr1041 chr1042 chr1043 chr1044 chr1045 chr1046 chr1047 chr1048 chr1049 chr1050 chr1051 chr1052 chr1053 chr1054 chr1055 chr1056 chr1057 chr1058 chr1059 chr1060 chr1061 chr1062 chr1063 chr1064 chr1065 chr1066 chr1067 chr1068 chr1069 chr1070 chr1071 chr1072 chr1073 chr1074 chr1075 chr1076 chr1077 chr1078 chr1079 chr1080 chr1081 chr1082 chr1083 chr1084 chr1085 chr1086 chr1087 chr1088 chr1089 chr1090 chr1091 chr1092 chr1093 chr1094 chr1095 chr1096 chr1097 chr1098 chr1099 chr1100 chr1101 chr1102 chr1103 chr1104 chr1105 chr1106 chr1107 chr1108 chr1109 chr1110 chr1111 chr1112 chr1113 chr1114 chr1115 chr1116 chr1117 chr1118 chr1119 chr1120 chr1121 chr1122 chr1123 chr1124 chr1125 chr1126 chr1127 chr1128 chr1129 chr1130 chr1131 chr1132 chr1133 chr1134 chr1135 chr1136 chr1137 chr1138 chr1139 chr1140 chr1141 chr1142 chr1143 chr1144 chr1145 chr1146 chr1147 chr1148 chr1149 chr1150 chr1151 chr1152 chr1153 chr1154 chr1155 chr1156 chr1157 chr1158 chr1159 chr1160 chr1161 chr1162 chr1163 chr1164 chr1165 chr1166 chr1167 chr1168 chr1169 chr1170 chr1171 chr1172 chr1173 chr1174 chr1175 chr1176 chr1177 chr1178 chr1179 chr1180 chr1181 chr1182 chr1183 chr1184 chr1185 chr1186 chr1187 chr1188 chr1189 chr1190 chr1191 chr1192 chr1193 chr1194 chr1195 chr1196 chr1197 chr1198 chr1199 chr1200 chr1201 chr1202 chr1203 chr1204 chr1205 chr1206 chr1207 chr1208 chr1209 chr1210 chr1211 chr1212 chr1213 chr1214 chr1215 chr1216 chr1217 chr1218 chr1219 chr1220 chr1221 chr1222 chr1223 chr1224 chr1225 chr1226 chr1227 chr1228 chr1229 chr1230 chr1231 chr1232 chr1233 chr1234 chr1235 chr1236 chr1237 chr1238 chr1239 chr1240 chr1241 chr1242 chr1243 chr1244 chr1245 chr1246 chr1247 chr1248 chr1249 chr1250 chr1251 chr1252 chr1253 chr1254 chr1255 chr1256 chr1257 chr1258 chr1259 chr1260 chr1261 chr1262 chr1263 chr1264 chr1265 chr1266 chr1267 chr1268 chr1269 chr1270 chr1271 chr1272 chr1273 chr1274 chr1275 chr1276 chr1277 chr1278 chr1279 chr1280 chr1281 chr1282 chr1283 chr1284 chr1285 chr1286 chr1287 chr1288 chr1289 chr1290 chr1291 chr1292 chr1293 chr1294 chr1295 chr1296 chr1297 chr1298 chr1299 chr1300 chr1301 chr1302 chr1303 chr1304 chr1305 chr1306 chr1307 chr1308 chr1309 chr1310 chr1311 chr1312 chr1313 chr1314 chr1315 chr1316 chr1317 chr1318 chr1319 chr1320 chr1321 chr1322 chr1323 chr1324 chr1325 chr1326 chr1327 chr1328 chr1329 chr1330 chr1331 chr1332 chr1333 chr1334 chr1335 chr1336 chr1337 chr1338 chr1339 chr1340 chr1341 chr1342 chr1343 chr1344 chr1345 chr1346 chr1347 chr1348 chr1349 chr1350 chr1351 chr1352 chr1353 chr1354 chr1355 chr1356 chr1357 chr1358 chr1359 chr1360 chr1361 chr1362 chr1363 chr1364 chr1365 chr1366 chr1367 chr1368 chr1369 chr1370 chr1371 chr1372 chr1373 chr1374 chr1375 chr1376 chr1377 chr1378 chr1379 chr1380 chr1381 chr1382 chr1383 chr1384 chr1385 chr1386 chr1387 chr1388 chr1389 chr1390 chr1391 chr1392 chr1393 chr1394 chr1395 chr1396 chr1397 chr1398 chr1399 chr1400 chr1401 chr1402 chr1403 chr1404 chr1405 chr1406 chr1407 chr1408 chr1409 chr1410 chr1411 chr1412 chr1413 chr1414 chr1415 chr1416 chr1417 chr1418 chr1419 chr1420 chr1421 chr1422 chr1423 chr1424 chr1425 chr1426 chr1427 chr1428 chr1429 chr1430 chr1431 chr1432 chr1433 chr1434 chr1435 chr1436 chr1437 chr1438 chr1439 chr1440 chr1441 chr1442 chr1443 chr1444 chr1445 chr1446 chr1447 chr1448 chr1449 chr1450 chr1451 chr1452 chr1453 chr1454 chr1455 chr1456 chr1457 chr1458 chr1459 chr1460 chr1461 chr1462 chr1463 chr1464 chr1465 chr1466 chr1467 chr1468 chr1469 chr1470 chr1471 chr1472 chr1473 chr1474 chr1475 chr1476 chr1477 chr1478 chr1479 chr1480 chr1481 chr1482 chr1483 chr1484 chr1485 chr1486 chr1487 chr1488 chr1489 chr1490 chr1491 chr1492 chr1493 chr1494 chr1495 chr1496 chr1497 chr1498 chr1499 chr1500 chr1501 chr1502 chr1503 chr1504 chr1505 chr1506 chr1507 chr1508 chr1509 chr1510 chr1511 chr1512 chr1513 chr1514 chr1515 chr1516 chr1517 chr1518 chr1519 chr1520 chr1521 chr1522 chr1523 chr1524 chr1525 chr1526 chr1527 chr1528 chr1529 chr1530 chr1531 chr1532 chr1533 chr1534 chr1535 chr1536 chr1537 chr1538 chr1539 chr1540 chr1541 chr1542 chr1543 chr1544 chr1545 chr1546 chr1547 chr1548 chr1549 chr1550 chr1551 chr1552 chr1553 chr1554 chr1555 chr1556 chr1557 chr1558 chr1559 chr1560 chr1561 chr1562 chr1563 chr1564 chr1565 chr1566 chr1567 chr1568 chr1569 chr1570 chr1571 chr1572 chr1573 chr1574 chr1575 chr1576 chr1577 chr1578 chr1579 chr1580 chr1581 chr1582 chr1583 chr1584 chr1585 chr1586 chr1587 chr1588 chr1589 chr1590 chr1591 chr1592 chr1593 chr1594 chr1595 chr1596 chr1597 chr1598 chr1599 chr1600 chr1601 chr1602 chr1603 chr1604 chr1605 chr1606 chr1607 chr1608 chr1609 chr1610 chr1611 chr1612 chr1613 chr1614 chr1615 chr1616 chr1617 chr1618 chr1619 chr1620 chr1621 chr1622 chr1623 chr1624 chr1625 chr1626 chr1627 chr1628 chr1629 chr1630 chr1631 chr1632 chr1633 chr1634 chr1635 chr1636 chr1637 chr1638 chr1639 chr1640 chr1641 chr1642 chr1643 chr1644 chr1645 chr1646 chr1647 chr1648 chr1649 chr1650 chr1651 chr1652 chr1653 chr1654 chr1655 chr1656 chr1657 chr1658 chr1659 chr1660 chr1661 chr1662 chr1663 chr1664 chr1665 chr1666 chr1667 chr1668 chr1669 chr1670 chr1671 chr1672 chr1673 chr1674 chr1675 chr1676 chr1677 chr1678 chr1679 chr1680 chr1681 chr1682 chr1683 chr1684 chr1685 chr1686 chr1687 chr1688 chr1689 chr1690 chr1691 chr1692 chr1693 chr1694 chr1695 chr1696 chr1697 chr1698 chr1699 chr1700 chr1701 chr1702 chr1703 chr1704 chr1705 chr1706 chr1707 chr1708 chr1709 chr1710 chr1711 chr1712 chr1713 chr1714 chr1715 chr1716 chr1717 chr1718 chr1719 chr1720 chr1721 chr1722 chr1723 chr1724 chr1725 chr1726 chr1727 chr1728 chr1729 chr1730 chr1731 chr1732 chr1733 chr1734 chr1735 chr1736 chr1737 chr1738 chr1739 chr1740 chr1741 chr1742 chr1743 chr1744 chr1745 chr1746 chr1747 chr1748 chr1749 chr1750 chr1751 chr1752 chr1753 chr1754 chr1755 chr1756 chr1757 chr1758 chr1759 chr1760 chr1761 chr1762 chr1763 chr1764 chr1765 chr1766 chr1767 chr1768 chr1769 chr1770 chr1771 chr1772 chr1773 chr1774 chr1775 chr1776 chr1777 chr1778 chr1779 chr1780 chr1781 chr1782 chr1783 chr1784 chr1785 chr1786 chr1787 chr1788 chr1789 chr1790 chr1791 chr1792 chr1793 chr1794 chr1795 chr1796 chr1797 chr1798 chr1799 chr1800 chr1801 chr1802 chr1803 chr1804 chr1805 chr1806 chr1807 chr1808 chr1809 chr1810 chr1811 chr1812 chr1813 chr1814 chr1815 chr1816 chr1817 chr1818 chr1819 chr1820 chr1821 chr1822 chr1823 chr1824 chr1825 chr1826 chr1827 chr1828 chr1829 chr1830 chr1831 chr1832 chr1833 chr1834 chr1835 chr1836 chr1837 chr1838 chr1839 chr1840 chr1841 chr1842 chr1843 chr1844 chr1845 chr1846 chr1847 chr1848 chr1849 chr1850 chr1851 chr1852 chr1853 chr1854 chr1855 chr1856 chr1857 chr1858 chr1859 chr1860 chr1861 chr1862 chr1863 chr1864 chr1865 chr1866 chr1867 chr1868 chr1869 chr1870 chr1871 chr1872 chr1873 chr1874 chr1875 chr1876 chr1877 chr1878 chr1879 chr1880 chr1881 chr1882 chr1883 chr1884 chr1885 chr1886 chr1887 chr1888 chr1889 chr1890 chr1891 chr1892 chr1893 chr1894 chr1895 chr1896 chr1897 chr1898 chr1899 chr1900 chr1901 chr1902 chr1903 chr1904 chr1905 chr1906 chr1907 chr1908 chr1909 chr1910 chr1911 chr1912 chr1913 chr1914 chr1915 chr1916 chr1917 chr1918 chr1919 chr1920 chr1921 chr1922 chr1923 chr1924 chr1925 chr1926 chr1927 chr1928 chr1929 chr1930 chr1931 chr1932 chr1933 chr1934 chr1935 chr1936 chr1937 chr1938 chr1939 chr1940 chr1941 chr1942 chr1943 chr1944 chr1945 chr1946 chr1947 chr1948 chr1949 chr1950 chr1951 chr1952 chr1953 chr1954 chr1955 chr1956 chr1957 chr1958 chr1959 chr1960 chr1961 chr1962 chr1963 chr1964 chr1965 chr1966 chr1967 chr1968 chr1969 chr1970 chr1971 chr1972 chr1973 chr1974 chr1975 chr1976 chr1977 chr1978 chr1979 chr1980 chr1981 chr1982 chr1983 chr1984 chr1985 chr1986 chr1987 chr1988 chr1989 chr1990 chr1991 chr1992 chr1993 chr1994 chr1995 chr1996 chr1997 chr1998 chr1999 chr2000 chr2001 chr2002 chr2003 chr2004 chr2005 chr2006 chr2007 chr2008 chr2009 chr2010 chr2011 chr2012 chr2013 chr2014 chr2015 chr2016 chr2017 chr2018 chr2019 chr2020 chr2021 chr2022 chr2023 chr2024 chr2025 chr2026 chr2027 chr2028 chr2029 chr2030 chr2031 chr2032 chr2033 chr2034 chr2035 chr2036 chr2037 chr2038 chr2039 chr2040 chr2041 chr2042 chr2043 chr2044 chr2045 chr2046 chr2047 chr2048 chr2049 chr2050 chr2051 chr2052 chr2053 chr2054 chr2055 chr2056 chr2057 chr2058 chr2059 chr2060 chr2061 chr2062 chr2063 chr2064 chr2065 chr2066 chr2067 chr2068 chr2069 chr2070 chr2071 chr2072 chr2073 chr2074 chr2075 chr2076 chr2077 chr2078 chr2079 chr2080 chr2081 chr2082 chr2083 chr2084 chr2085 chr2086 chr2087 chr2088 chr2089 chr2090 chr2091 chr2092 chr2093 chr2094 chr2095 chr2096 chr2097 chr2098 chr2099 chr2100 chr2101 chr2102 chr2103 chr2104 chr2105 chr2106 chr2107 chr2108 chr2109 chr2110 chr2111 chr2112 chr2113 chr2114 chr2115 chr2116 chr2117 chr2118 chr2119 chr2120 chr2121 chr2122 chr2123 chr2124 chr2125 chr2126 chr2127 chr2128 chr2129 chr2130 chr2131 chr2132 chr2133 chr2134 chr2135 chr2136 chr2137 chr2138 chr2139 chr2140 chr2141 chr2142 chr2143 chr2144 chr2145 chr2146 chr2147 chr2148 chr2149 chr2150 chr2151 chr2152 chr2153 chr2154 chr2155 chr2156 chr2157 chr2158 chr2159 chr2160 chr2161 chr2162 chr2163 chr2164 chr2165 chr2166 chr2167 chr2168 chr2169 chr2170 chr2171 chr2172 chr2173 chr2174 chr2175 chr2176 chr2177 chr2178 chr2179 chr2180 chr2181 chr2182 chr2183 chr2184 chr2185 chr2186 chr2187 chr2188 chr2189 chr2190 chr2191 chr2192 chr2193 chr2194 chr2195 chr2196 chr2197 chr2198 chr2199 chr2200 chr2201 chr2202 chr2203 chr2204 chr2205 chr2206 chr2207 chr2208 chr2209 chr2210 chr2211 chr2212 chr2213 chr2214 chr2215 chr2216 chr2217 chr2218 chr2219 chr2220 chr2221 chr2222 chr2223 chr2224 chr2225 chr2226 chr2227 chr2228 chr2229 chr2230 chr2231 chr2232 chr2233 chr2234 chr2235 chr2236 chr2237 chr2238 chr2239 chr2240 chr2241 chr2242 chr2243 chr2244 chr2245 chr2246 chr2247 chr2248 chr2249 chr2250 chr2251 chr2252 chr2253 chr2254 chr2255 chr2256 chr2257 chr2258 chr2259 chr2260 chr2261 chr2262 chr2263 chr2264 chr2265 chr2266 chr2267 chr2268 chr2269 chr2270 chr2271 chr2272 chr2273 chr2274 chr2275 chr2276 chr2277 chr2278 chr2279 chr2280 chr2281 chr2282 chr2283 chr2284 chr2285 chr2286 chr2287 chr2288 chr2289 chr2290 chr2291 chr2292 chr2293 chr2294 chr2295 chr2296 chr2297 chr2298 chr2299 chr2300 chr2301 chr2302 chr2303 chr2304 chr2305 chr2306 chr2307 chr2308 chr2309 chr2310 chr2311 chr2312 chr2313 chr2314 chr2315
```

[illegible]

The functional annotation table can then be transformed and saved in different data formats using the `prad2*`() functions. Options a

- `pred2bed()` : Format LTR prediction data to BED file format
- `pred2csv()` : Format LTR prediction data to CSV file format
- `pred2fasta()` : Save the sequence of the predicted LTR Transposons in a fasta file
- `pred2gff()` : Format LTR prediction data to GFF3 file format

Developed by Mark George Upton

Site built with [pkgdown \(https://pkgdown.r-lib.org/\)](https://pkgdown.r-lib.org/) 1.6.1.

```
pred2gff (../reference/pred2gff.html){LTR_data = Screenshot1.png
output   = "Screenshot1_functional_LTRpred_gff",
program  = "LTRpred"}
```

### Detailed description of adjustable LTRpred parameters

The *S. cerevisiae* example shown above assumes that users wish to run iTigrid using a default parameter configuration. However, as the figure shown below demonstrates, there are large number of parameters that can be adjusted and altered according to the user's specifications and interests.

Users can adjust and run LTPred with the following parameter options:

[illegible][illegible]

- `penaltiescore` specify the match score used in the PSSM/PSA Smith-Waterman alignment. Default is `penaltiescore = 1`.
- `positivescore` specify the mismatch score used in the PSSM/PSA Smith-Waterman alignment. Default is `positivescore = -10`.
- `phoscorescore` specify the insertion score used in the PSSM/PSA Smith-Waterman alignment. Default is `phoscorescore = -20`.
- `phosidionscore` specify the deletion score used in the PSSM/PSA Smith-Waterman alignment. Default is `phosidionscore = -20`.
- `pfm.in` is a character vector storing the Pfm (2k from <http://pfm.stanford.edu>) that shall be downloaded and used to perform protein domain searches within the sequences between the predicted LTRs.
- `cores` the number of cores that shall be used for multicore processing. In case of `nuc.cores` and `am.cores` are not specified then the value of `cores` is used for those arguments.
- `am.cores` number of cores to be used for multicore processing when running `blast` query (in case `am.cores` = "0" or "1")
- `nuc.cores` number of cores to be used for multicore processing when performing `blast` protein search with `LTRpred`.
- `orf.type` type of predicting open reading frames (see documentation of `ORFfinder`).
- `min.cores` minimum number of codons in the predicted open reading frame.
- `trans.wrap` logical value indicating whether or not predicted open reading frames shall be translated and the corresponding protein sequences stored in the output folder.
- `output.path` a path/folder to store all results returned by `LTRpred`, `LTRpred`, and `LTRpred`. If `output.path` = `NA` (default) then a folder with the name of the input genome file will be generated in the current working directory of R and all results are then stored in this folder.
- `quality.rfilter` shall false positives be filtered out as much as possible? Default is `quality.rfilter` = `TRUE`. See Description for details.
- `u.orf` minimum number of Open Reading Frames that must be found between the LTRs (if `quality.rfilter` = `TRUE`). See Details for further information on quality control.
- `verbose` shall further information be printed on the console or not.

## Metagenome scale annotations

LTRpred allows users to perform annotations not only for single genomes but for multiple genomes (metagenomes) using only one pipeline function named `LTRpred.meta()`.

Please be aware that LTRpred annotations for multiple genomes is a computationally hard task and requires large server or high performance computer access. Computations for hundreds of genomes even using tens to hundreds of cores might take several weeks to terminate. Please make sure that you have the right computational infrastructure to run these processes.

Users can download the bioconductor (<http://bioconductor.org/packages/>) package to automatically retrieve genome assembly files for the species of interest.

```
# specify the scientific names of the species of interest
# that shall first be downloaded and then be used
# to generate LTRpred annotations
species <- c("http://refseq.nlm.nih.gov/assembly/Arabidopsis_thaliana", "Arabidopsis lyrata", "Capsella rubella")
# download the genome assembly files for the species of interest
#
# install.packages("biomart")
biomart::getGenomeSet(http://docs.biomart.org/Biomart/reference/getGenomeSet.html)(db = "refseq", organism = species, path = "store_genome_set")
# run LTRpred.meta() on the 3 species with 3 cores
LTRpred::LTRpred.meta(c(reference=LTRpred.meta.html|genome.folder = "store_genome_set",
                        output.folder = "LTRpred_meta_results",
                        cores = 3))
```