

Welcome !

Principles and history of RepeatExplorer

2007

First paper
on repeat
clustering
from NGS
data

BMC Genomics



Research article

Open Access

**Repetitive DNA in the pea (*Pisum sativum* L.) genome:
comprehensive characterization using 454 sequencing and
comparison to soybean and *Medicago truncatula***

Jiří Macas*, Pavel Neumann and Alice Navrátilová

Address: Biology Centre ASCR, Institute of Plant Molecular Biology, Branišovská 31, Žeské Budějovice, CZ-37005, Czech Republic

Email: Jiří Macas* - macas@umbr.cas.cz; Pavel Neumann - neumann@umbr.cas.cz; Alice Navrátilová - navratil@umbr.cas.cz

* Corresponding author

Published: 21 November 2007

Received: 13 August 2007

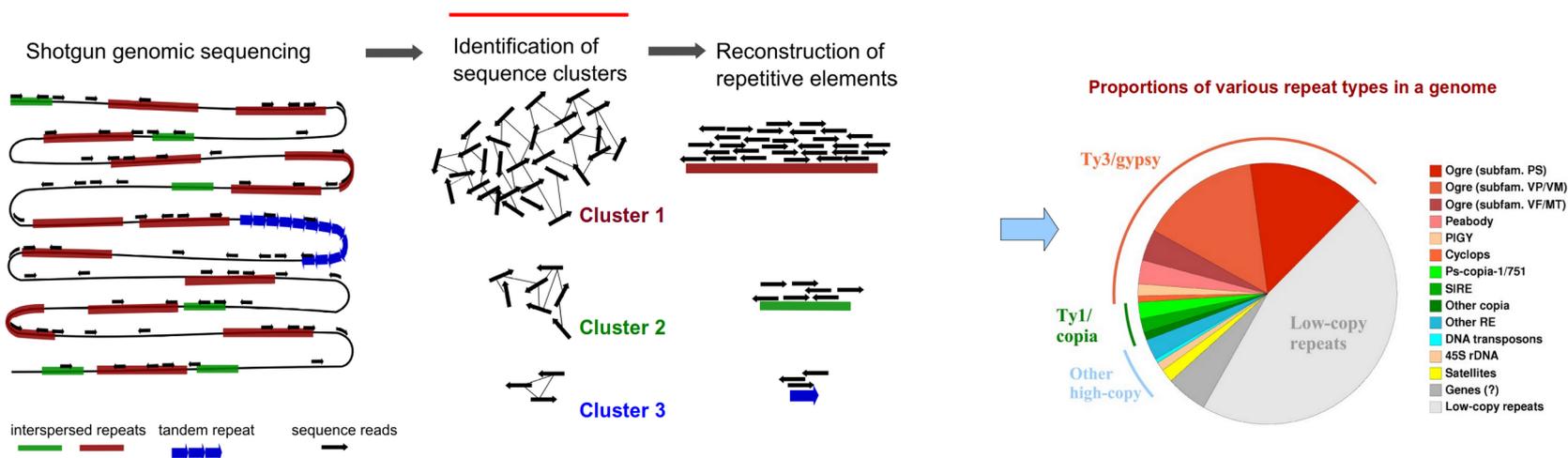
BMC Genomics 2007, 8:427 doi:10.1186/1471-2164-8-427

Accepted: 21 November 2007

Principles and history of RepeatExplorer

2007

First paper on repeat clustering from NGS data

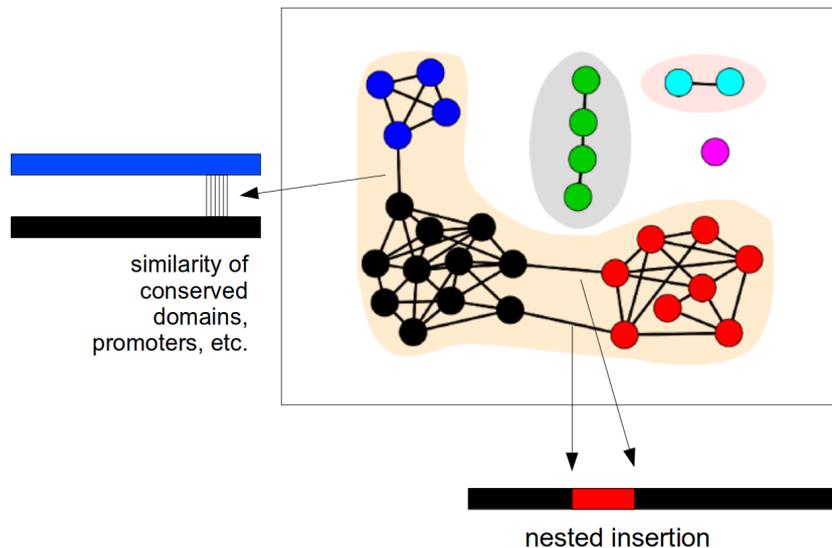


CLUSTER = a set of frequently overlapping reads = REPEAT FAMILY

Principles and history of RepeatExplorer

2007

First paper on repeat clustering from NGS data



Chimeric clusters !

Single linkage clustering => *connected components*

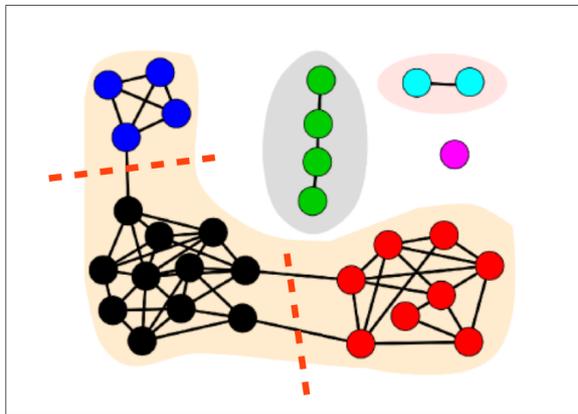
TGICL
(TIGR Gene Indices clustering tool)
Pertea et al., 2003

Principles and history of RepeatExplorer

2007 ... 2010

First paper on repeat clustering from NGS data

Introduction of **graph-based clustering** (Novak et al. 2010)



Graph-based clustering

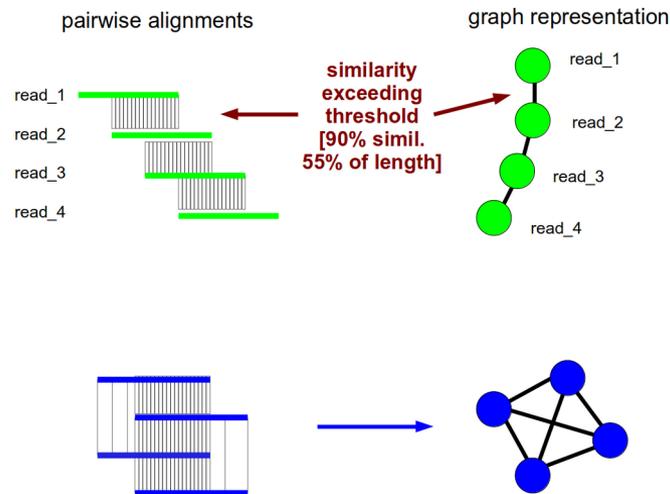
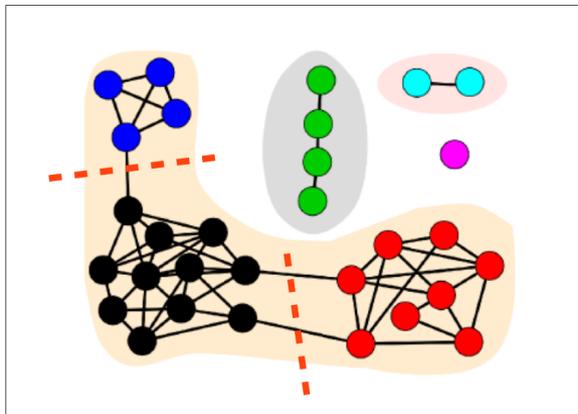
- Sequence overlaps between the reads are transformed to a graph where the **reads** are represented as **nodes** and their **similarities** as **edges** connecting the nodes
- Graph structure is examined to detect **communities of frequently connected nodes** which are **split to separate clusters**

Principles and history of RepeatExplorer

2007 ... 2010

First paper on repeat clustering from NGS data

Introduction of **graph-based clustering** (Novak et al. 2010)



Graph-based clustering

- Sequence overlaps between the reads are transformed to a graph represented as **nodes** and their **similarities** as **edges** connecting them
- Graph structure is examined to detect **communities of frequently connected nodes** which are **split to separate clusters**

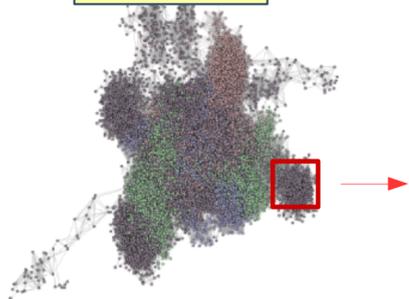


Principles and history of RepeatExplorer

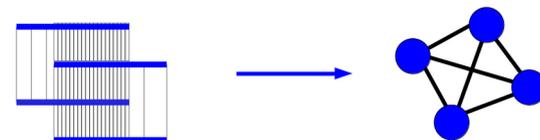
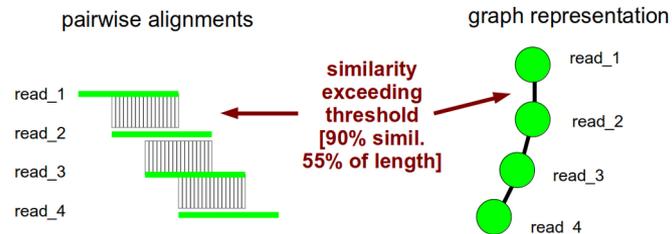
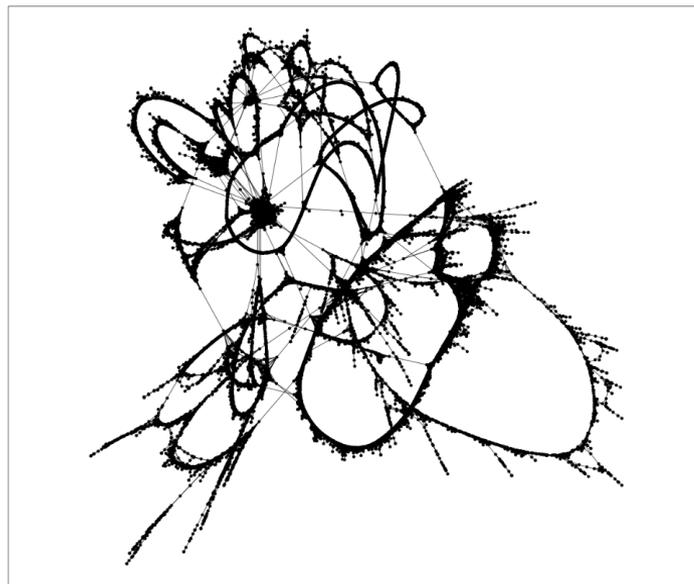
2007 ... 2010

First paper on repeat clustering from NGS data

Introduction of **graph-based clustering** (Novak et al. 2010)



Virtual graphs used to analyze real data contain **up to millions of nodes** (reads)

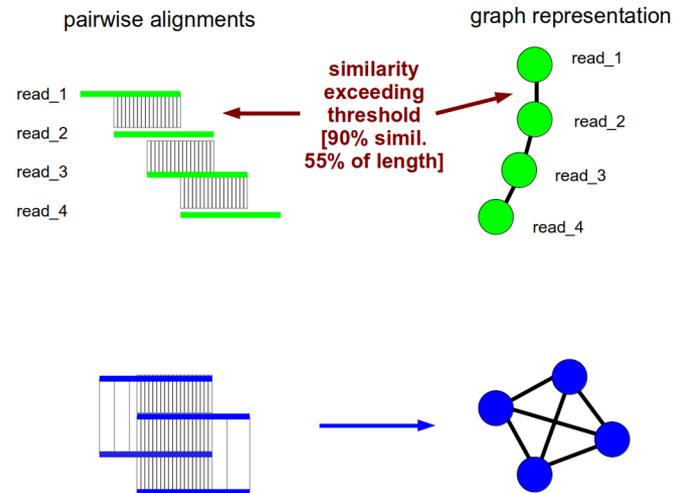
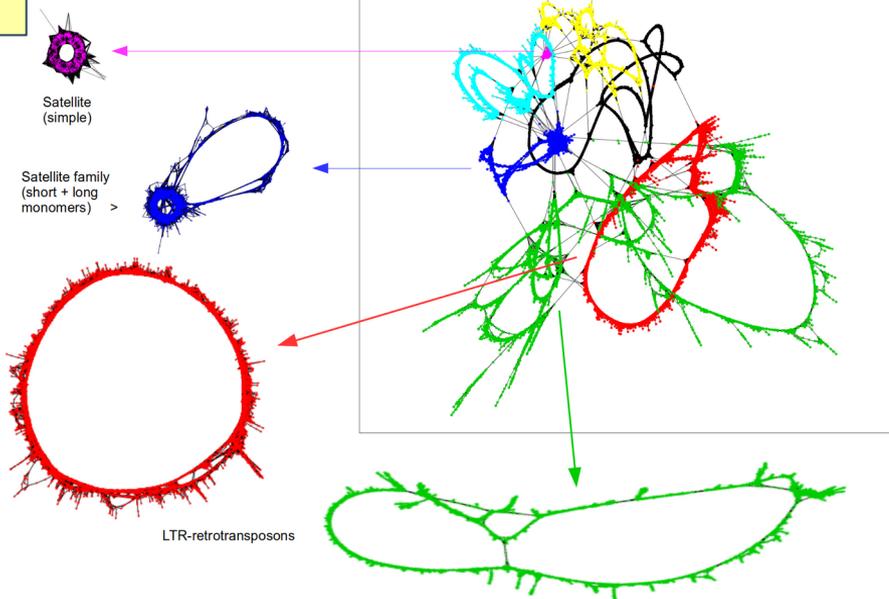


Principles and history of RepeatExplorer

2007 ... 2010

First paper on repeat clustering from NGS data

Introduction of **graph-based clustering** (Novak et al. 2010)



Principles and history of RepeatExplorer

2007

...

2010

...

2013

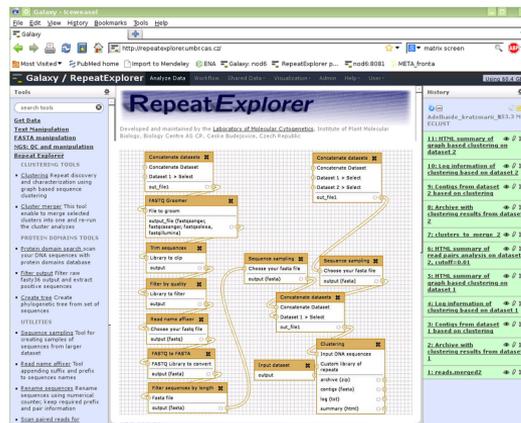
First paper on repeat clustering from NGS data

Introduction of **graph-based clustering** (Novak et al. 2010)

Repeat Explorer in Galaxy (Novak et al. 2013)

command-line version

Public web-based server



Principles and history of RepeatExplorer



2007
First paper on repeat clustering from NGS data

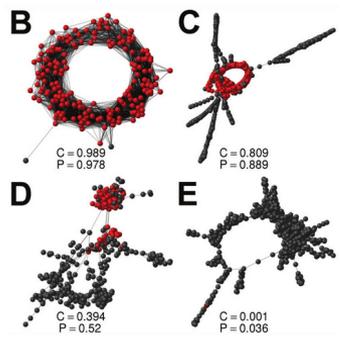
2010
Introduction of **graph-based clustering** (Novak et al. 2010)

2013
Repeat Explorer in Galaxy (Novak et al. 2013)



ELIXIR \$\$\$

2017
TAREAN

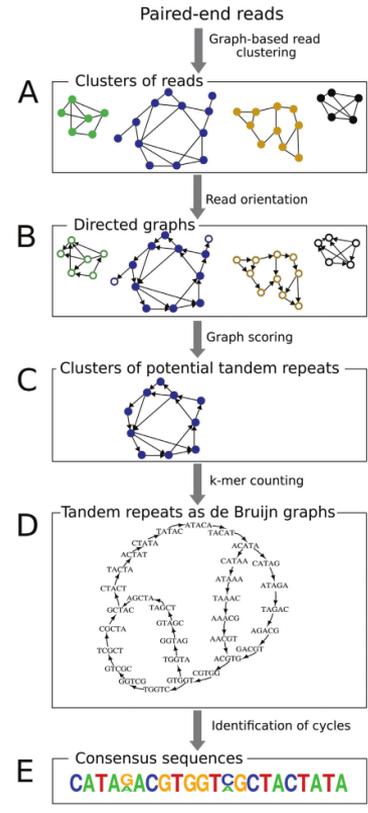


*Nucleic Acids Research, 2017 1
doi: 10.1093/nar/lkx257*

TAREAN: a computational tool for identification and characterization of satellite DNA from unassembled short reads

Petr Novák, Laura Ávila Robledillo, Andrea Koblížková, Iva Vrbová, Pavel Neumann and Jiří Macas*

Institute of Plant Molecular Biology, Biology Centre CAS, České Budějovice CZ-37005, Czech Republic



Principles and history of RepeatExplorer



First paper on repeat clustering from NGS data

Introduction of **graph-based clustering** (Novak et al. 2010)

Repeat Explorer in Galaxy (Novak et al. 2013)

Public web-based server



ELIXIR \$\$\$

TAREAN

ELIXIR / CERIT Galaxy server + data storage

RepeatExplorer
Discover repeats in your next generation sequencing data.

RepeatExplorer includes utilities for Graph-based clustering and characterization of repetitive sequences in next-generation sequencing data and tools for the detection of transposable element protein coding domains.

The user of Galaxy based RepeatExplorer is obliged to use the following acknowledgement formula in all your publications created with the support of RepeatExplorer: Computational resources were provided by the ELIXIR-CZ project (LM2015047), part of the international ELIXIR infrastructure.

If you need help with RepeatExplorer or you want to report a problem and our wiki is not able to give you answers, please contact server_administrator.

RepeatExplorer

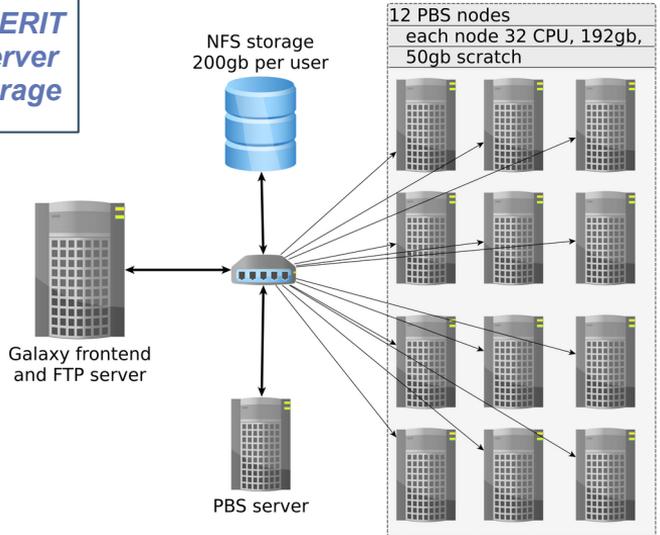
Go to Galaxy RepeatExplorer portal.

Wiki

Do you have questions how to upload data via ftp, etc. RepeatExplorer, etc.? Visit our Galaxy wiki.

Registration

Please, read our [registration manual](#).



Principles and history of RepeatExplorer



First paper on repeat clustering from NGS data

Introduction of **graph-based clustering** (Novak et al. 2010)

Repeat Explorer in Galaxy (Novak et al. 2013)

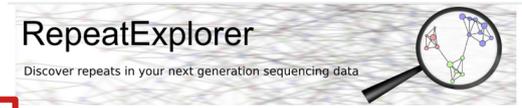
Public web-based server



ELIXIR \$\$\$

TAREAN

ELIXIR / CERIT Galaxy server + data storage



RepeatExplorer includes utilities for Graph-based clustering and characterization of repetitive sequences in next-generation sequencing data and tools for the detection of transposable element protein coding domains.

The user of Galaxy based RepeatExplorer is obliged to use the following acknowledgement formula in all your publications created with the support of RepeatExplorer: **Computational resources were provided by the ELIXIR-CZ project (LM2015047), part of the international ELIXIR infrastructure.**

If you need help with RepeatExplorer or you want to report a problem and our wiki is not able to give you answers, please contact [server_administrator](#).

RepeatExplorer

Go to Galaxy RepeatExplorer portal.

Wiki

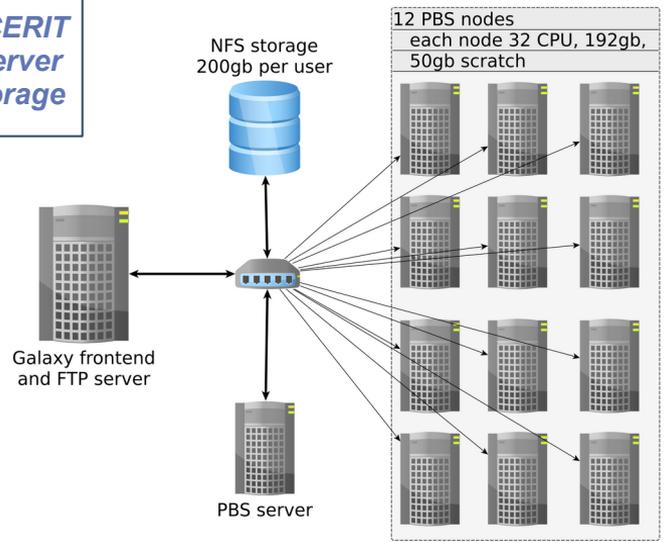
Do you have questions how to upload data via ftp, etc. RepeatExplorer, etc.?
Visit our Galaxy wiki.

Registration

Please, read our [registration manual](#).

The user of Galaxy based RepeatExplorer is obliged to use the following acknowledgement formula in all your publications created with the support of RepeatExplorer: **Computational resources were provided by the ELIXIR-CZ project (LM2015047), part of the international ELIXIR infrastructure.**

Please, acknowledge ELIXIR in your publications !



Principles and history of RepeatExplorer



2007
First paper on repeat clustering from NGS data

2010
Introduction of **graph-based clustering** (Novak et al. 2010)

2013
Repeat Explorer in Galaxy (Novak et al. 2013)



ELIXIR \$\$\$

2017
TAREAN

ELIXIR / CERIT Galaxy server + data storage

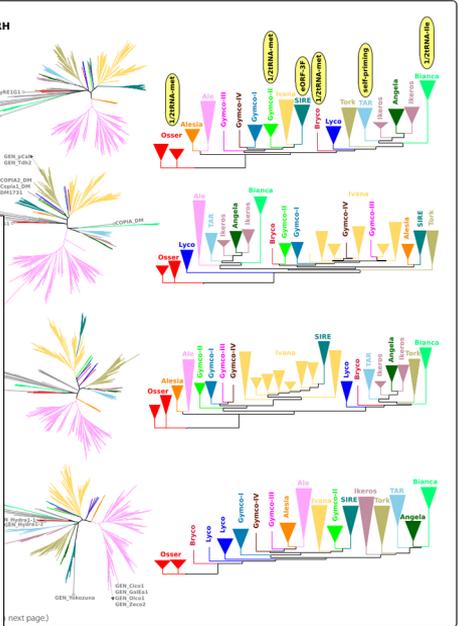
2019
REXdb database (Neumann et al. 2019)

Neumann et al. *Mobile DNA* (2019) 10:1
<https://doi.org/10.1186/s13100-018-0144-1>

RESEARCH **Open Access**

Systematic survey of plant LTR-retrotransposons elucidates phylogenetic relationships of their polyprotein domains and provides a reference for element classification

Pavel Neumann^{*}, Petr Novák, Nina Hošťáková and Jiří Macas



Principles and history of RepeatExplorer



First paper on repeat clustering from NGS data

Introduction of **graph-based clustering** (Novak et al. 2010)

Repeat Explorer in Galaxy (Novak et al. 2013)



TAREAN

REXdb (Neumann et al. 2019)

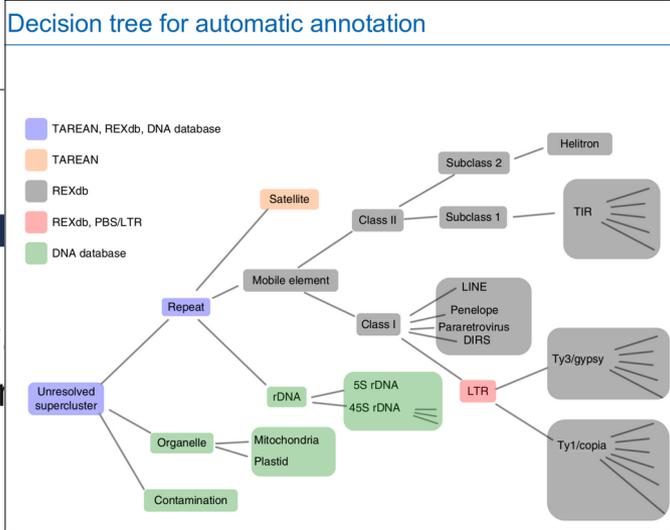
Repeat Explorer ver. 2

Neumann et al. *Mobile DNA* (2019) 10:1
<https://doi.org/10.1186/s13100-018-0144-1>

RESEARCH

Systematic survey of plant LTR-retrotransposons elucidates phylogenetic relationships of their polyprotein domain and provides a reference for element classification

Pavel Neumann*, Petr Novák, Nina Hošťáková and Jiří Macas



Principles and history of RepeatExplorer

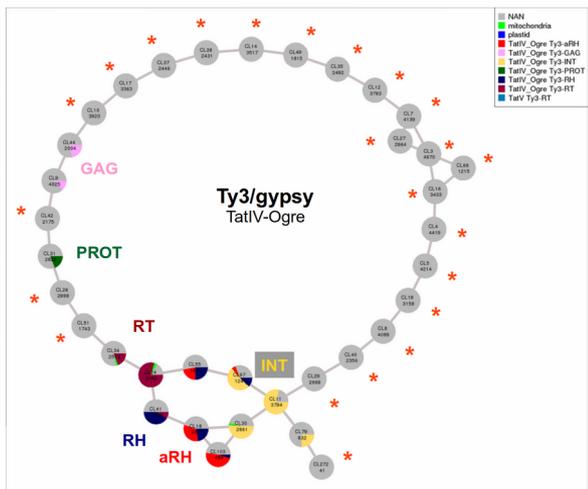


TAREAN

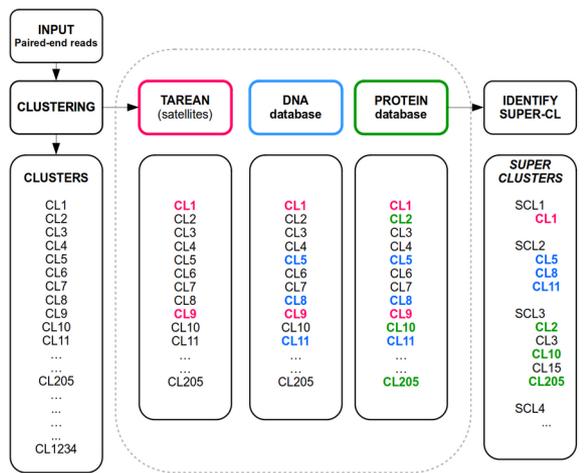
REXdb
(Neumann et al. 2019)

Repeat Explorer ver. 2

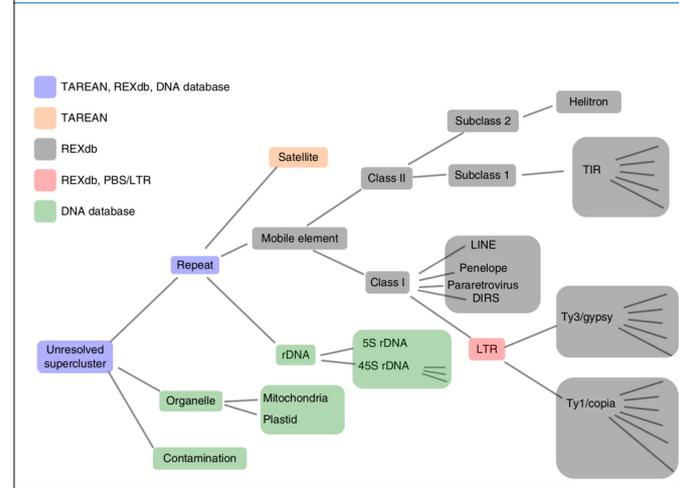
Superclusters provide more complete annotation



RepeatExplorer 2 – automatic detection of superclusters

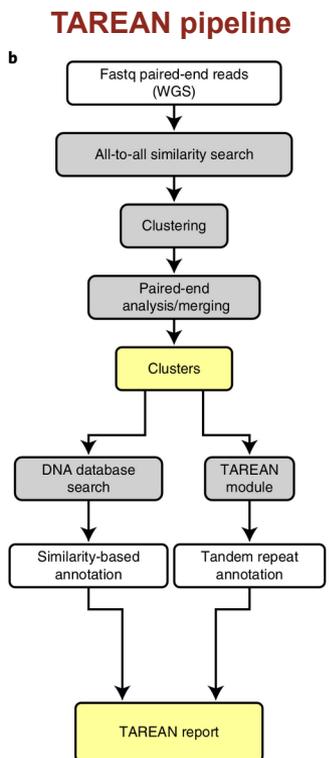
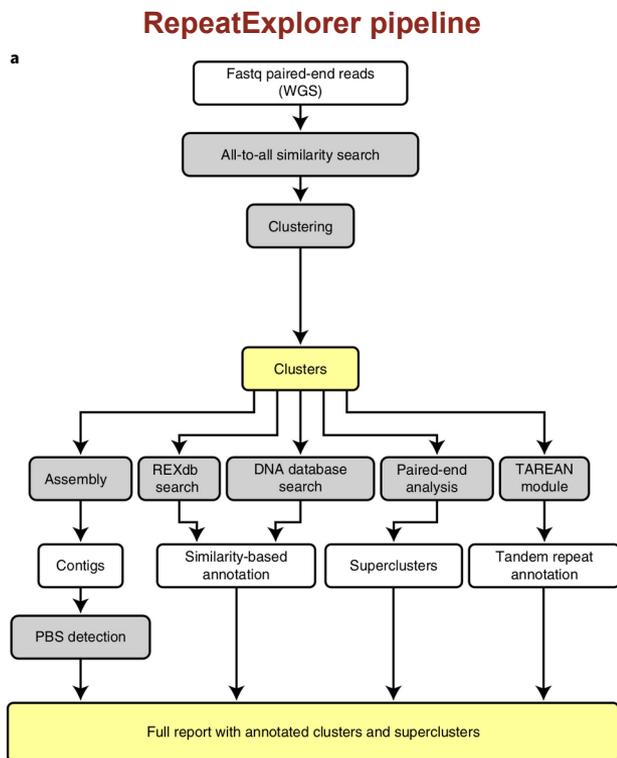


Decision tree for automatic annotation



Principles and history of RepeatExplorer

2007 ... 2010 ... 2013 2014 ... 2016 2017 2018 2019 2020 2023



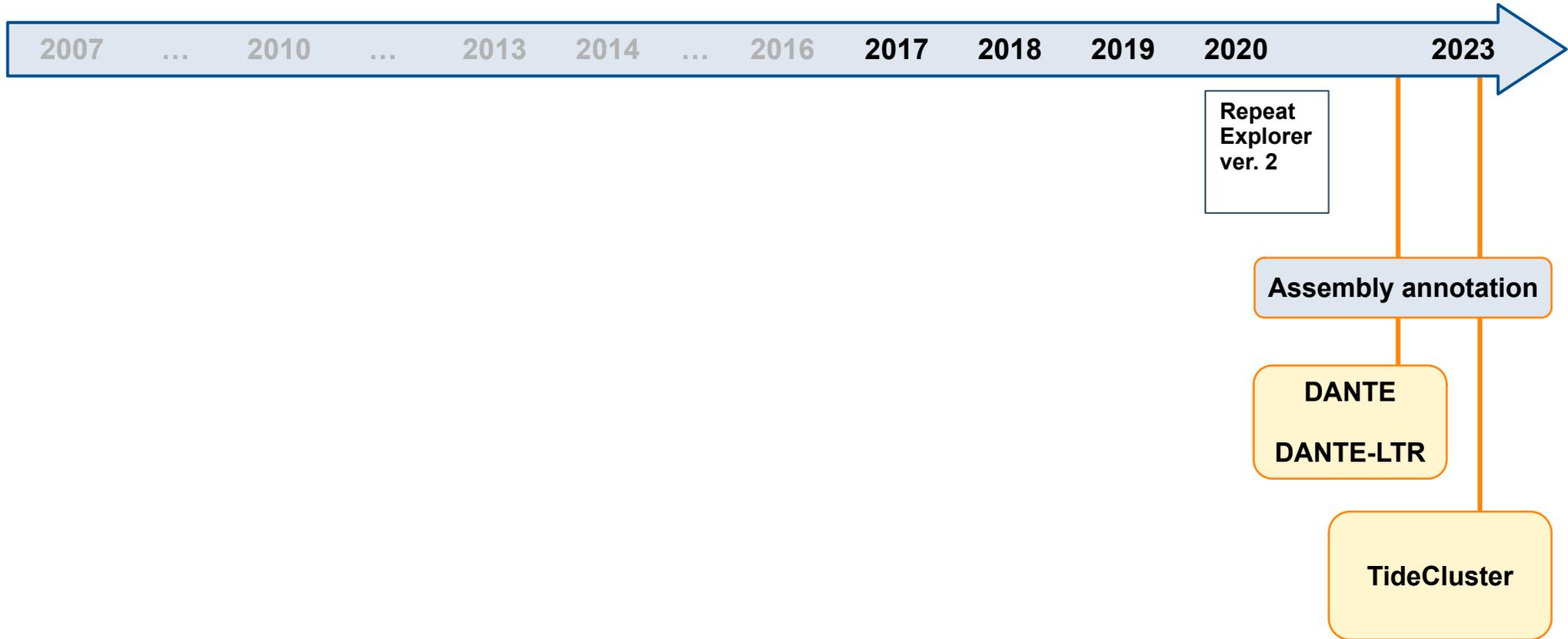
TAREAN

REXdb
(Neumann et al. 2019)

Repeat Explorer ver. 2

- Additional tools:**
- ChIP-seq Mapper
 - Long reads
 - DANTE

Principles and history of RepeatExplorer

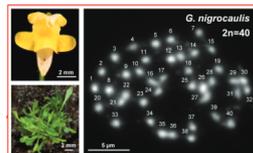
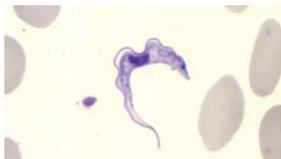
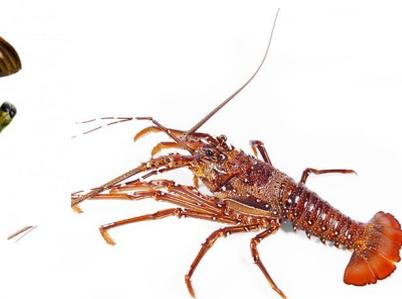


Applications

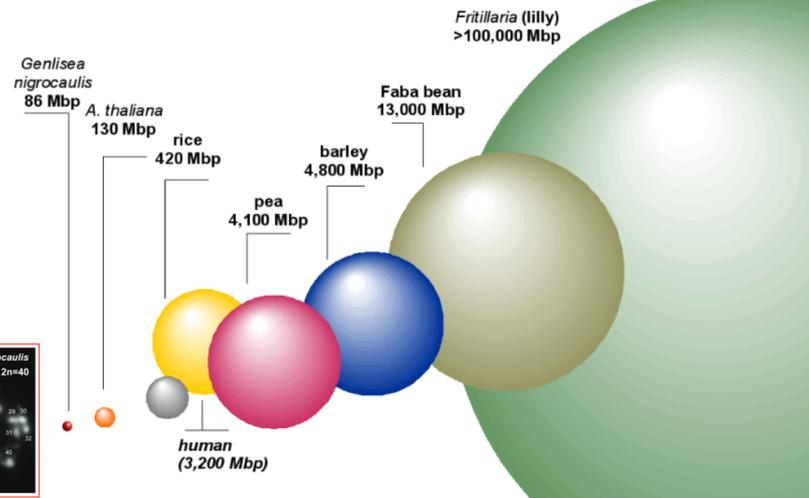
Applications

- Repeat composition

- single species



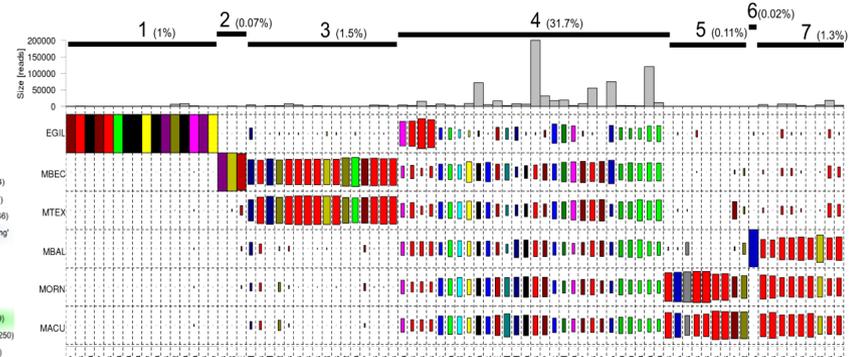
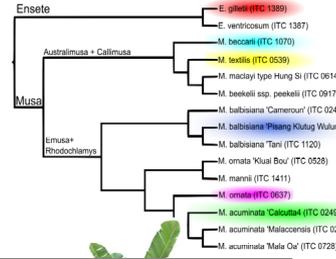
Plant species varying ~2,000-fold in their genome sizes



Applications

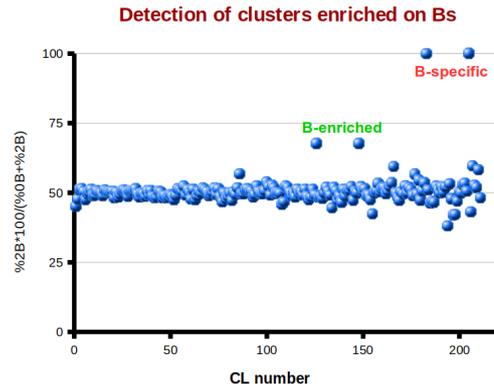
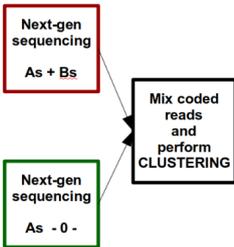
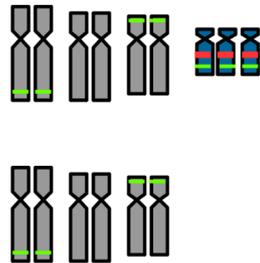
Repeat composition

- single species
- comparative analysis



Identification of chromosome B-specific repeats

Comparative analysis of B+/- plants



Genomic repeat abundances contain phylogenetic signal

Dodsworth et al. (2015) *Syst. Biol.* 64(1): 112-126

- Maximus
 - Angela
 - Tork
 - Ale
 - Ivana
 - TAR
 - Reina
 - Tekay
 - Galadriel
 - Tat
 - Chromovirus/unspecified
 - TRIM
 - pararetrovirus
 - LINE
 - DNA transposon
 - rDNA
 - TR
 - unknown
 - CRM
- (Novak et al., 2014)

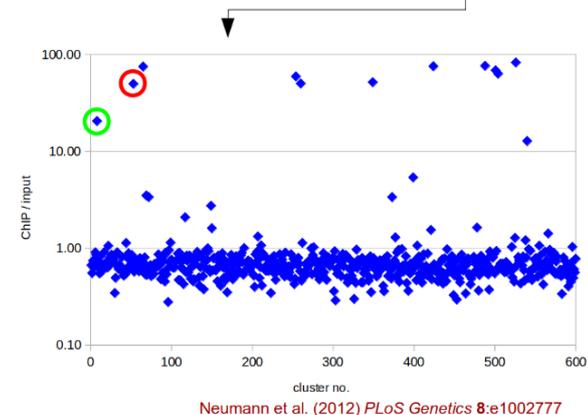
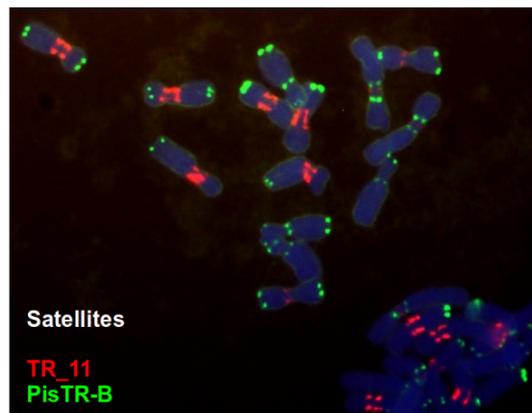
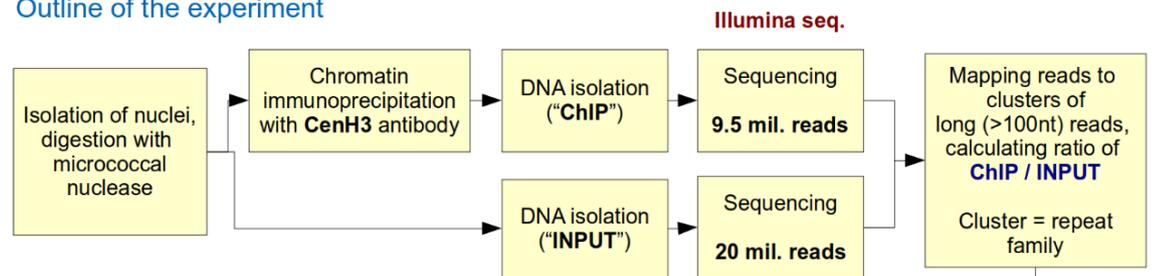
Applications

- Repeat composition
 - single species
 - comparative analysis

- Repeat clusters as a reference
 - ChIP-seq

Identification of centromeric repeats by ChIP-seq

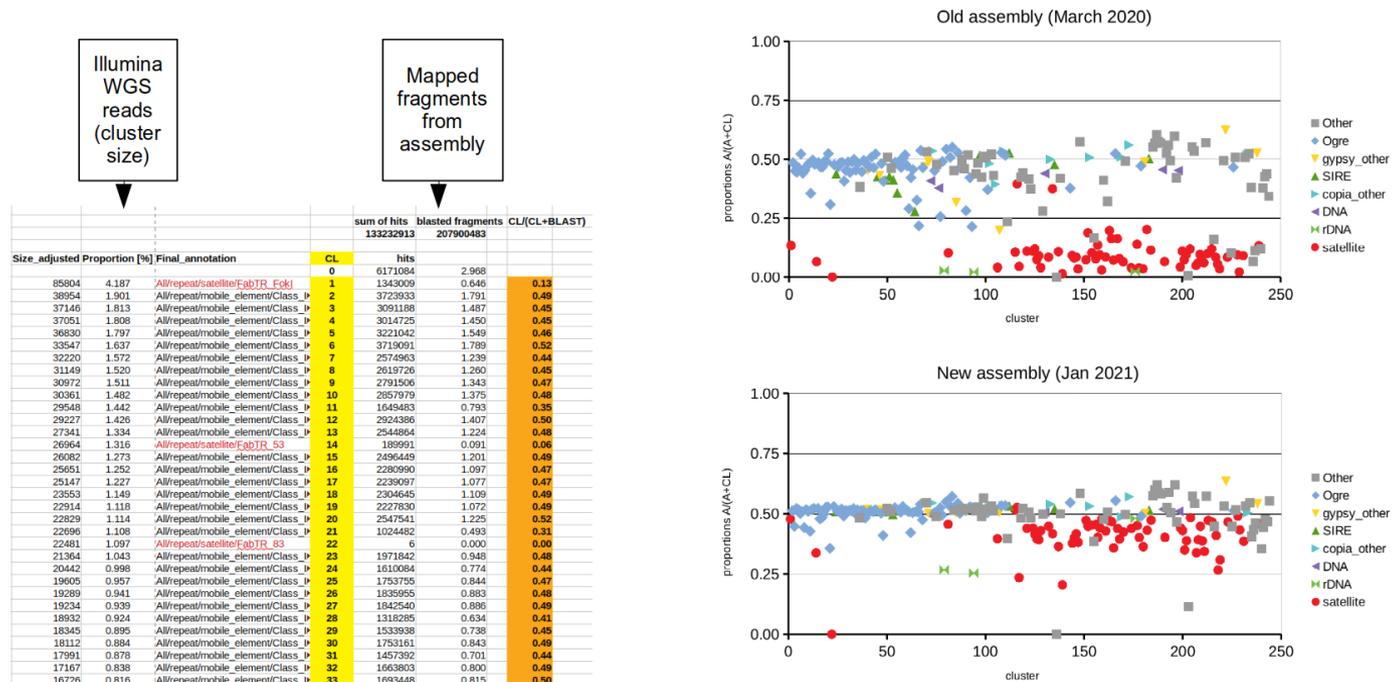
Outline of the experiment



Applications

- Repeat composition
 - single species
 - comparative analysis
- Repeat clusters as a reference
 - ChIP-seq
 - assembly

Assessing completeness of genome assemblies



Applications

Repeat composition

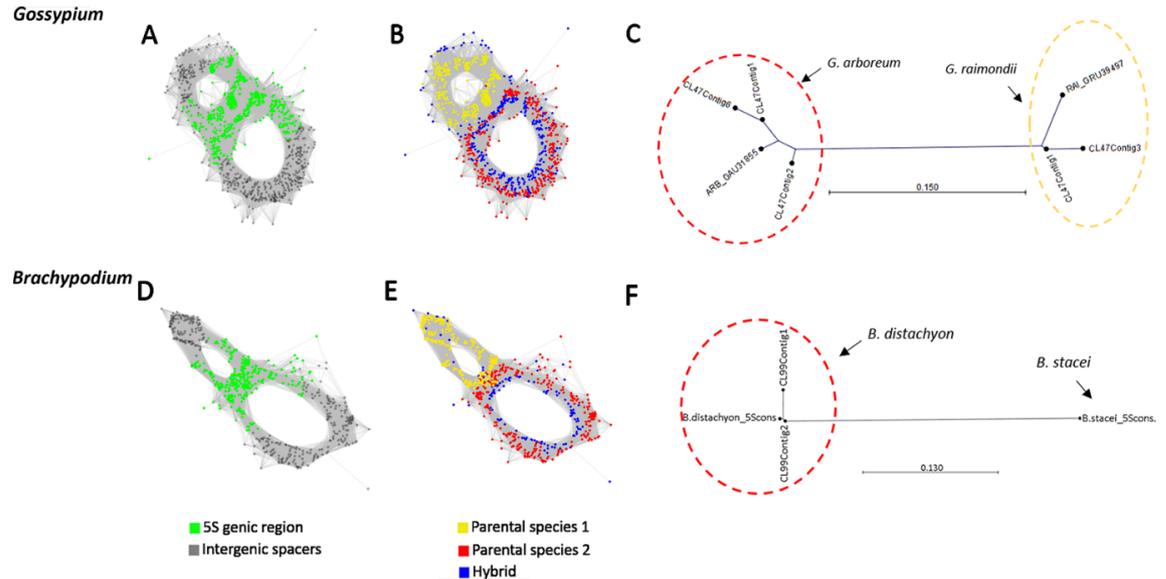
- single species
- comparative analysis

Repeat clusters as a reference

- ChIP-seq
- assembly
- **graph shapes**

The Utility of Graph Clustering of 5S Ribosomal DNA Homoeologs in Plant Allopolyploids, Homoploid Hybrids, and Cryptic Introgressants

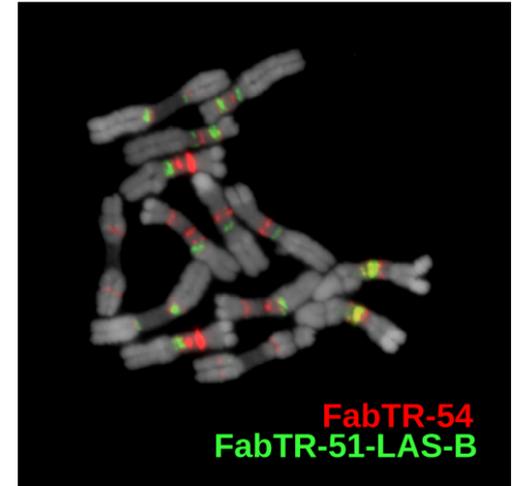
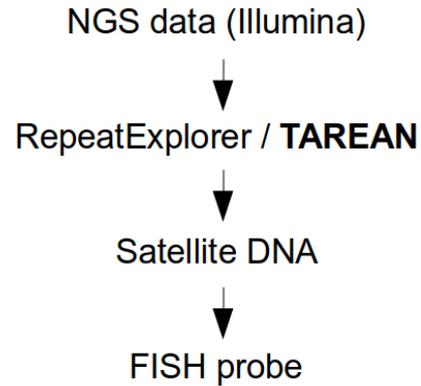
Sonia Garcia^{1,2}, Jonathan F. Wendel³, Natalia Borowska-Zuchowska⁴, Malika Ainouche⁵, Alena Kuderova² and Ales Kovarik^{2*}



Applications

- Repeat composition
 - single species
 - comparative analysis
- Repeat clusters as a reference
 - ChIP-seq
 - assembly
 - graph shapes*
- Satellite DNA
 - cytogenetic studies

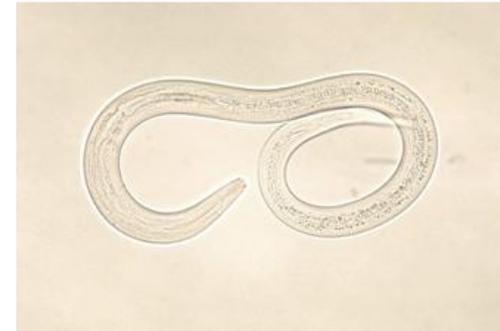
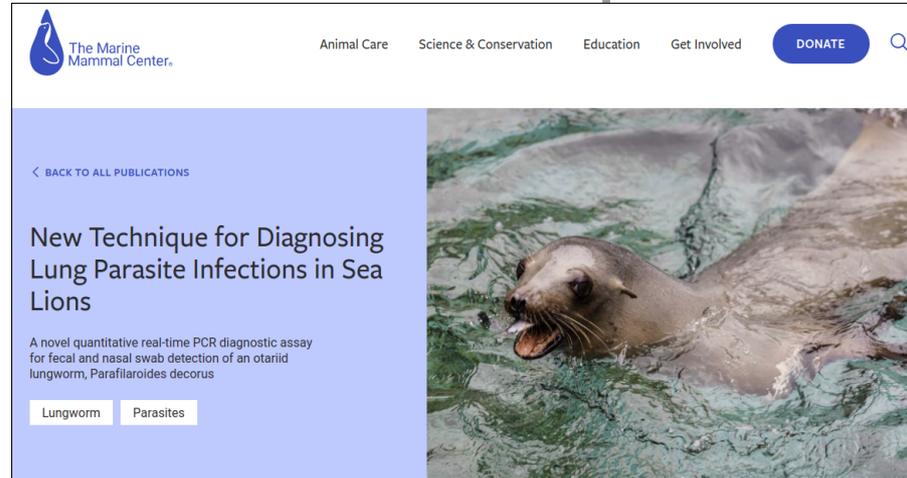
Cytogenetic markers (FISH probes)



Lathyrus sativus
(FISH by L. Avila Robledillo)

Applications

- Repeat composition
 - single species
 - comparative analysis
- Repeat clusters as a reference
 - ChIP-seq
 - assembly
 - graph shapes*
- Satellite DNA
 - cytogenetic studies
 - diagnostic markers



Applications

- Repeat composition

- single species
- comparative analysis

- Repeat clusters as a reference

- ChIP-seq
- assembly
- graph shapes

- Satellite DNA

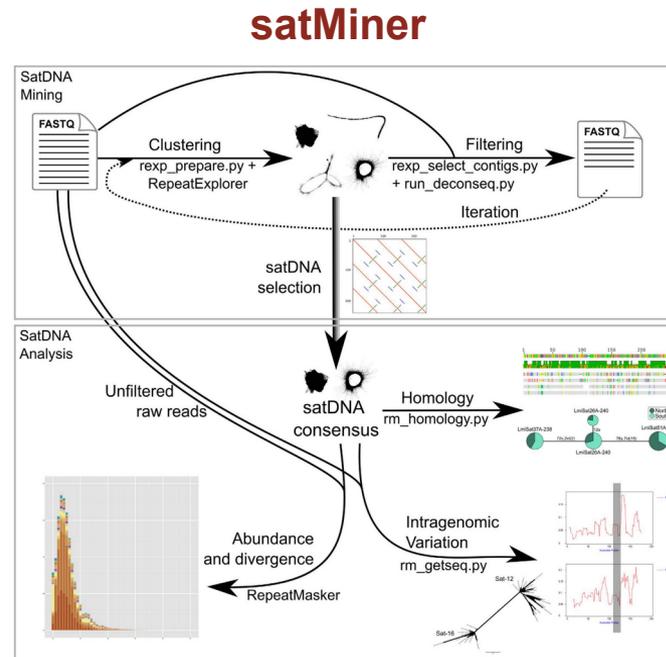
- cytogenetic studies
- diagnostic markers

- RE utilization in other pipelines

SCIENTIFIC REPORTS

OPEN High-throughput analysis of the satellitome illuminates satellite DNA evolution

Received: 14 January 2016 Accepted: 02 June 2016 Francisco J. Ruiz-Ruano, María Dolores López-León, Josefa Cabrero & Juan Pedro M. Camacho



Applications

Repeat composition

- single species
- comparative analysis

Repeat clusters as a reference

- ChIP-seq
- assembly
- graph shapes

Satellite DNA

- cytogenetic studies
- diagnostic markers

RE utilization in other pipelines

Mann et al. *BMC Bioinformatics* (2022) 23:40
<https://doi.org/10.1186/s12859-021-04545-2>

BMC Bioinformatics

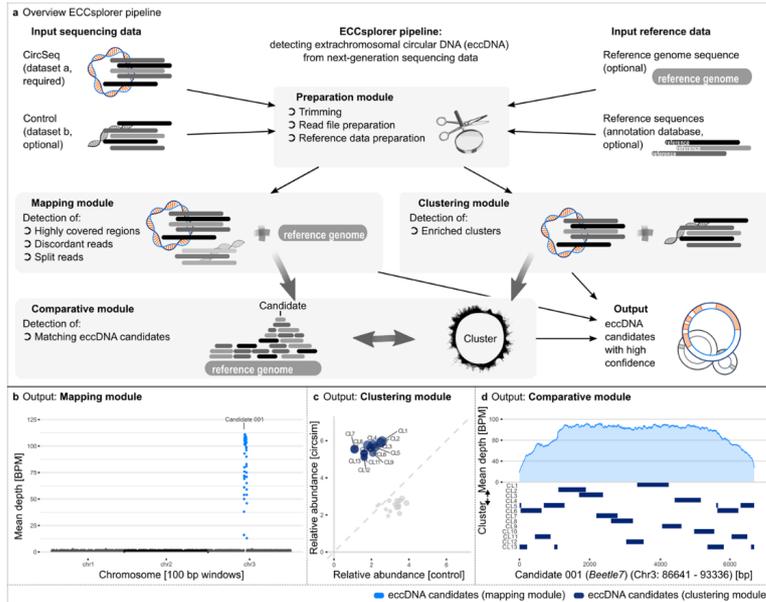
SOFTWARE

Open Access



ECCsplorer: a pipeline to detect extrachromosomal circular DNA (eccDNA) from next-generation sequencing data

Ludwig Mann, Kathrin M. Seibt, Beatrice Weber and Tony Heitkam



Protocols and tutorials

RepeatExplorer principles and example protocols published

Corresponding video tutorials available from  YouTube



Global analysis of repetitive DNA from unassembled sequence reads using RepeatExplorer2

Petr Novák, Pavel Neumann and Jiří Macas

RepeatExplorer2 is a novel version of a computational pipeline that uses graph-based sequencing reads for characterization of repetitive DNA in eukaryotes. The pipeline identifies in any genome by using relatively small quantities of short sequence reads. It performs automatic annotation and quantification of the identified repeats. Galaxy platform, which provides a user-friendly web interface for script execution. Compared to the original version of the pipeline, RepeatExplorer2 provides auto elements, identification of tandem repeats and enhanced visualization of analysis results. RepeatExplorer2 workflow and provide procedures for its application to (i) de novo species, (ii) comparative repeat analysis in a set of species, (iii) development of a pipeline to identify centromeric repeats based on ChIP-seq data, and (iv) identification of centromeric repeats based on ChIP-seq data. RepeatExplorer2 is available at <https://repeatexplorer-eltir.cesb.cas.cz/>

Introduction

Complex eukaryotic genomes, including those of higher plants, are characterized by various types of repetitive sequences. These sequence elements (retroelements and DNA transposons) that are distributed in arrays of tandemly repeated satellite DNA that constitute the repetitive genome regions. Differential amplification and retention of repetitive sequences facilitate rapid genome restructuring and evolution. Consequently, repetitive DNA attracted the interest of genome biologists, which in turn prompted the development of high-throughput sequencing technologies, which can generate large volumes of genomic data rapidly. Because the lack of such data had been the bottleneck in the investigations of repetitive DNA in non-model species, the demand for computational tools specifically tailored for repetitive DNA analysis has increased. Several principally different approaches have been investigated: the reference databases of previously characterized repeats, algorithms that detect repeats based on structural features of the analyzed sequences (reviewed in refs. 1, 2). The unique approach known as graph-based repeat clustering, constitutes the

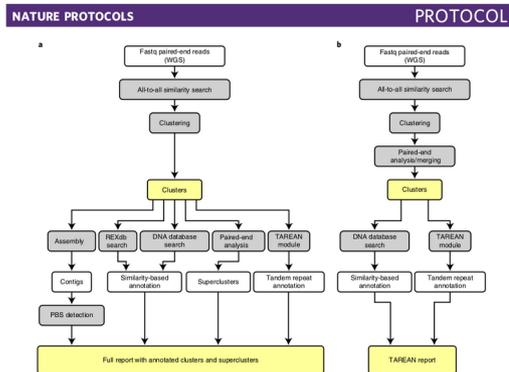


Fig. 1 | Schematic representation of RepeatExplorer (a) and TAREAN (b) pipelines. Analysis modules are represented by gray boxes, and input and output data are white, with the most important outputs highlighted in yellow. The RepeatExplorer pipeline is used in Procedures 1, 2 and 4; TAREAN is used in Procedure 3. WGS, whole-genome sequencing.

Overview of RepeatExplorer2

Recently, we released RepeatExplorer2, a new version of the pipeline that includes improvements of the existing programs and databases, as well as extended functionality due to inclusion of several novel tools. Although the basic workflow remains the same, a new module performs automated annotation of the clusters based on the similarity hits to the reference databases and utilizing

NATURE PROTOCOLS

Box 6 | RepeatExplorer2 archive structure

Below is a description of the most important files within the output data, extracted from a zip archive. The extracted structure (folders are bold):

```
UnzippedGalaxyArchive
├── libdir
├── seqclust
├── small_clusters_assembly
├── custom_databases
├── reads
├── clustering
│   ├── clusters
│   ├── superclusters
│   └── hitsort.cls
├── logfile.txt
├── index.html
├── cluster_report.html
├── tarean_report.html
├── supercluster_report.html
├── summarized_annotation.html
├── PROFREP_CLASSIFICATION_TEMPLATE.csv
├── CLUSTER_TABLE.csv
├── SUPERCLUSTER_TABLE.csv
├── contigs.fasta
├── TAREAN_consensus_rank_1.fasta
├── TAREAN_consensus_rank_2.fasta
├── TAREAN_consensus_rank_3.fasta
├── TAREAN_consensus_rank_4.fasta
├── summary_histogram.png
├── documentation.html
├── HOW_TO_CITE.html
└── style1.css
```

can be opened in a plain text editor to see how your analysis proceeded. The analysis went, as well as whether any error messages were printed out, can be opened in your web browser and provides a summary of the clustering. The other HTML files are linked to index.html, provide more detailed information for a single cluster.

ing basic information about clusters like the size of the cluster (i.e., the number of reads) and results of automatic annotation. The pipeline also performs summarizing automatic annotation of superclusters.