

Welcome !

Program at a glance

Tuesday (May 24)

Introduction and update on RepeatExplorer

- Principles and history of RepeatExplorer
(J. Macas)
- Diagnostic features of repetitive elements – Part I
(P. Neumann)
- Coffee break*
- Diagnostic features of repetitive elements – Part II
(P. Neumann)
- RepeatExplorer tools for genome annotation
(P. Novák)

Lunch

Practical training I

Basic Protocols 1-4: troubleshooting & advanced data analysis

Dinner at CITYgastro restaurant

Wednesday (May 25)

Short talks

- Nicola Schmidt
- Ludwig Mann
- Zirlane Portugal da Costa

Coffee break

- Pol Fernández Mató
- Yennifer Mata-Sucre
- Nusrat Sultana

Lunch

Practical training II

- using RE output for annotating genome assemblies
- REXdb and DANTE
- structure-based annotation of complete LTR-retrotransposons

Thursday (May 26)

Short talks

- Matej Lexa
- Monika Čechová
- Sophie Maiwald
- Vratislav Peška

Coffee break

- Camila do Nascimento Moreira
- Veit Herklotz
- Alice Krumpolcová

Lunch

Practical training III

- reconstructing phylogenetic relationships
- local installation and running RE tools from a command line
- topics proposed by the participants & individual consultations

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 clustering identification in any genome by using relatively small quantities of short sequence repeats perform automatic annotation and quantification of the identified repeats. Galaxy platform, which provides a user-friendly web interface for the execution. Compared to the original version of this pipeline, RepeatExplorer2 provides automated elements, identification of tandem repeats and enhanced visualization of analysis results. The RepeatExplorer2 workflow and provide procedures for its application to (i) different species, (ii) comparative repeat analysis in a set of species, (iii) development of new experiments and (iv) identification of centromeric repeats based on ChIP-seq data. It is designed to complete. RepeatExplorer2 is available at <https://repeatexplorer-elixir.cert.europa.eu>

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 arrays of tandemly repeated satellite DNA that constitute major components of plant genomes. Plant repetitive DNA is one of the main drivers of genome evolution, leading to, for example, the 2,400 plant species.¹ Repetitive DNA is present even in the smallest genomes, up to >15% of nuclear DNA,² and its proportions increase in animal genomes, reaching values as high as 85%.³ In addition, repetitive sequences facilitate rapid genomic restructuring among plant genomes.

Consequently, repetitive DNA attracted the interest of genome biology, which in turn prompted the development of methods for its quantification. This process was greatly accelerated by the introduction of next-generation sequencing (NGS) technologies,⁴ which can generate large volumes of data rapidly. Because the lack of such data had been the bottleneck in the investigations of repetitive DNA in non-model species, the development of new pipelines for the analysis of NGS data of repetitive DNA has been a priority. Several principally different approaches have been implemented in the reference databases of previously characterized algorithms that detect repeats based on structural features of the analyzed sequences (reviewed in refs.^{5,6}). The unique feature of RepeatExplorer⁷, the core component of the pipeline, constitutes the use of graph-based repeat clustering,⁸ which is able to identify clusters of repeats with different lengths and orientations.

NATURE PROTOCOLS

a

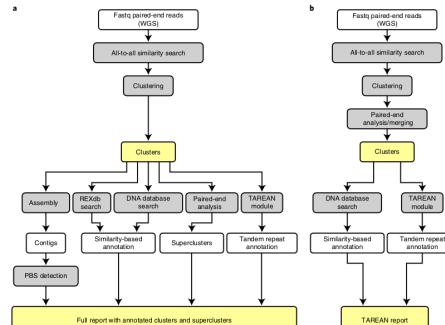


Fig. 1 | Schematic representation of RepeatExplorer (a) and TAREAN (b) pipelines. Analysis modules are represented by gray boxes, and input and output data are white, with the most important outputs highlighted in yellow. The RepeatExplorer pipeline is used in Procedures 1, 2 and 4; TAREAN is used in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

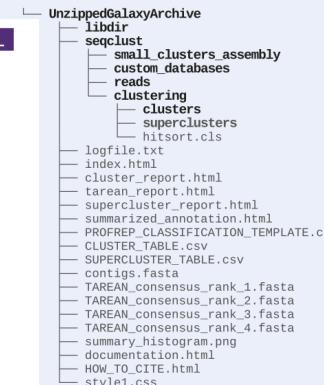
Overview of RepeatExplorer²

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

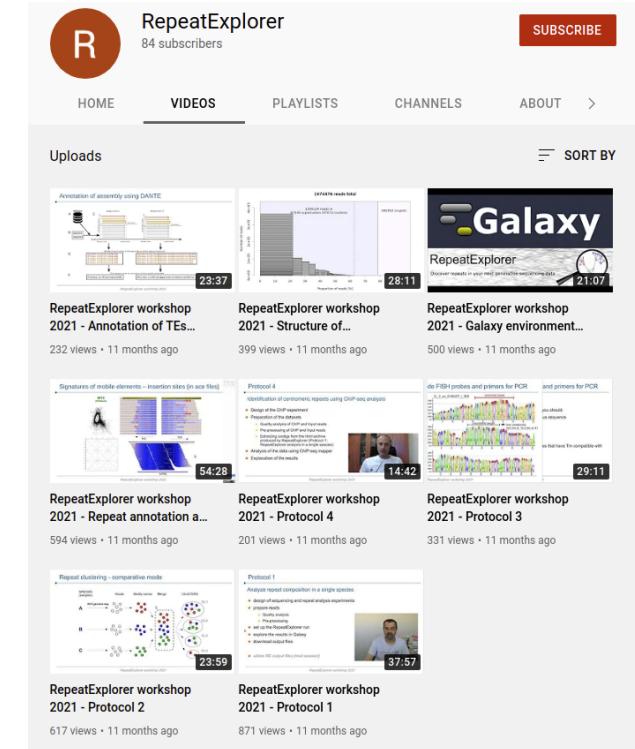
NATURE PROTOCOLS

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):



which 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. You 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 specific cluster. It also provides the size of the cluster, the number of repeats and of automatic annotation. Summarizing automatic annotation of superclusters.



Principles and history of RepeatExplorer

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

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

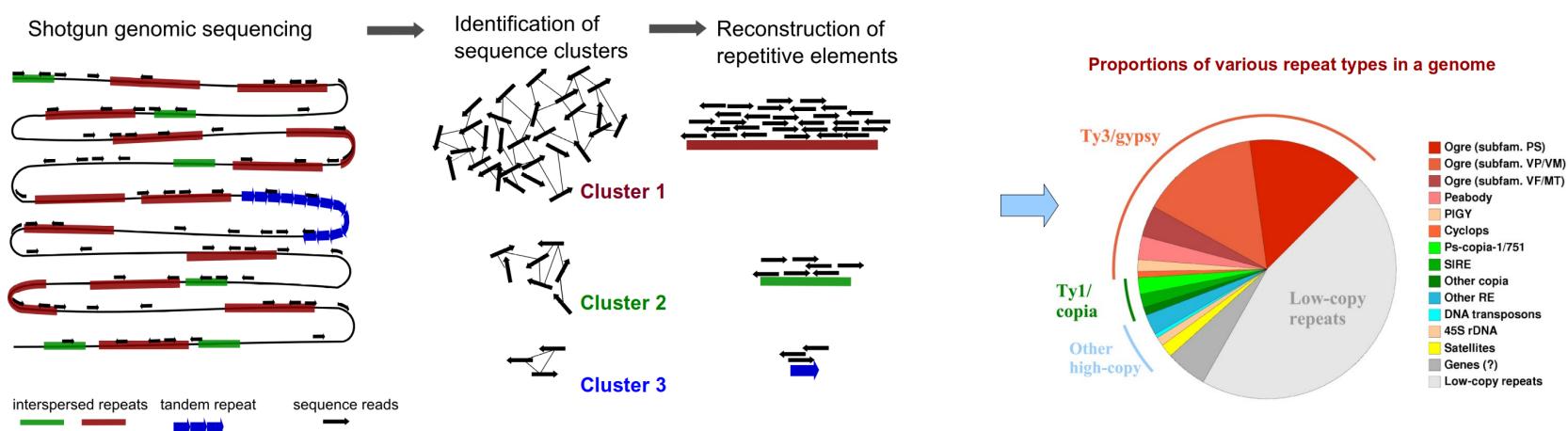
Received: 13 August 2007

Accepted: 21 November 2007

Principles and history of RepeatExplorer

2007

First paper
on repeat
clustering
from NGS
data

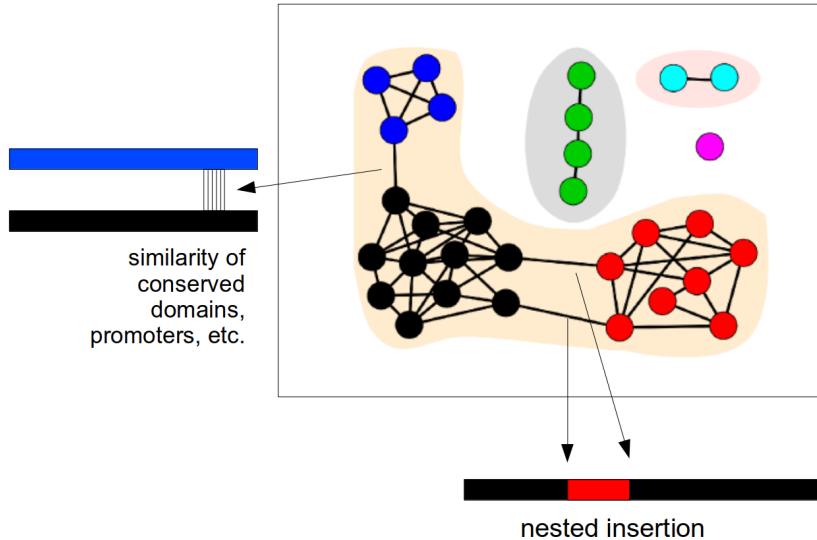


CLUSTER = a set of frequently overlapping reads = REPEAT FAMILY

Principles and history of RepeatExplorer

2007

First paper
on repeat
clustering
from NGS
data



Chimeric clusters !

Single linkage clustering => connected components

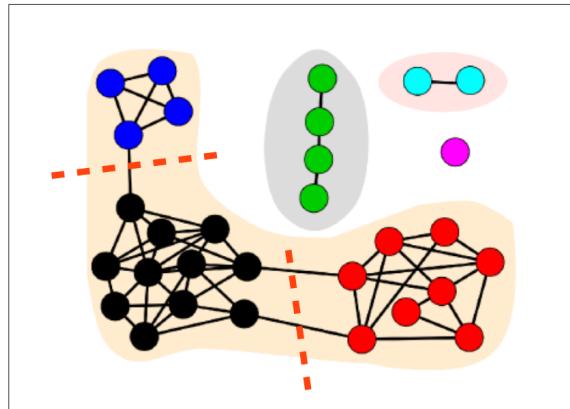
TGICL
(TIGR Gene Indices clustering tool)
Pertea et al., 2003

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

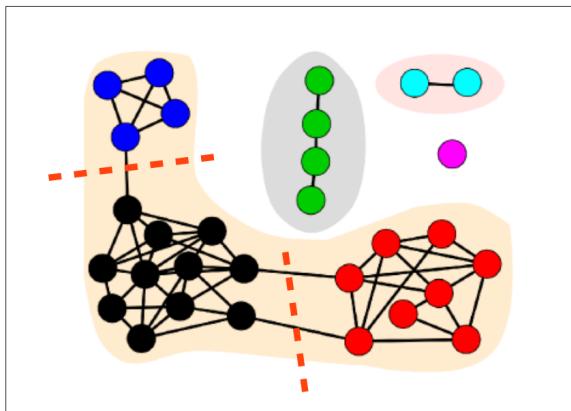
- Sequence overlaps between the reads are transformed to a graph where the **reads** are represented as **nodes** and their **similarities** as **edges** connecting the nodes
- Graph structure is examined to detect **communities of frequently connected nodes** which are **split to separate clusters**

Principles and history of RepeatExplorer

2007 ... 2010

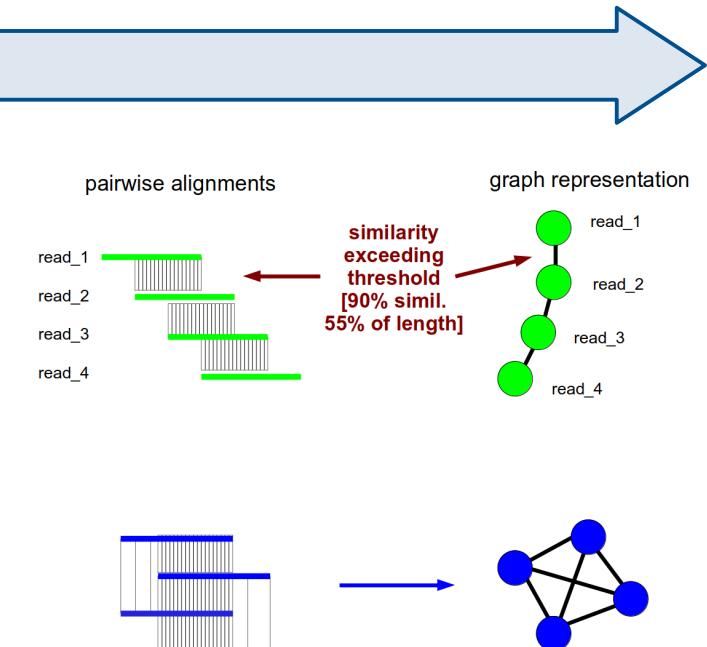
First paper
on repeat
clustering
from NGS
data

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

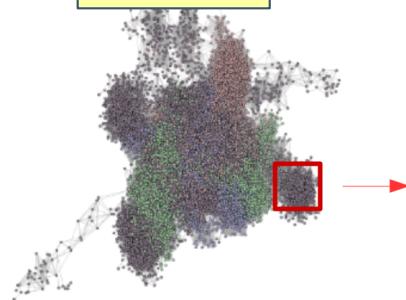


Principles and history of RepeatExplorer

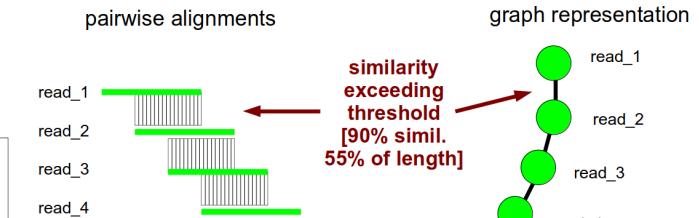
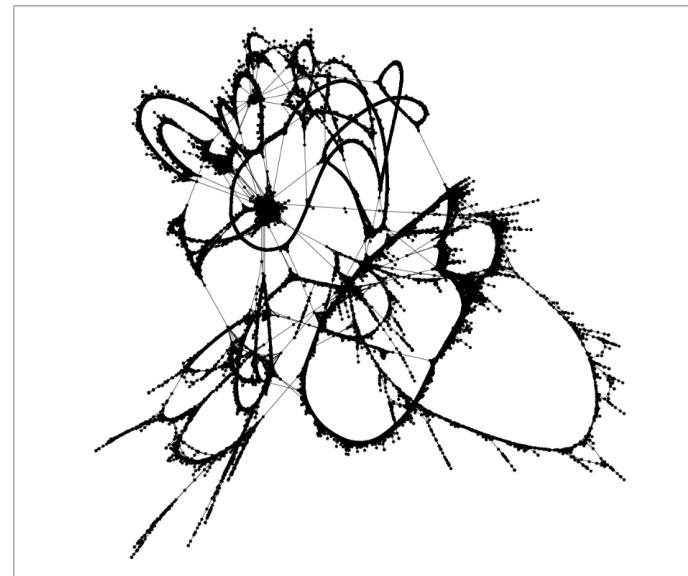
2007 ... 2010

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)

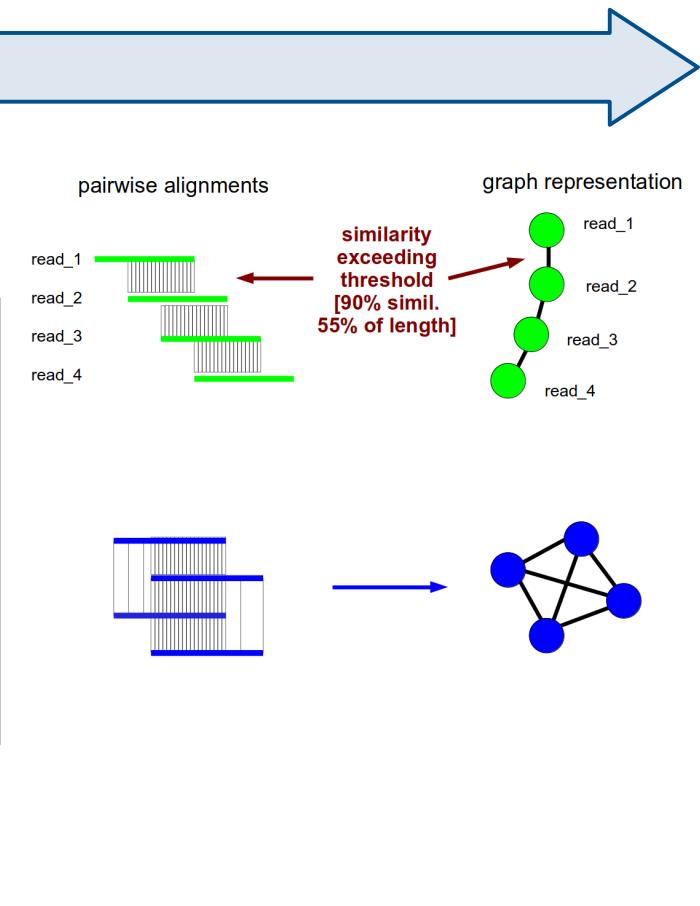
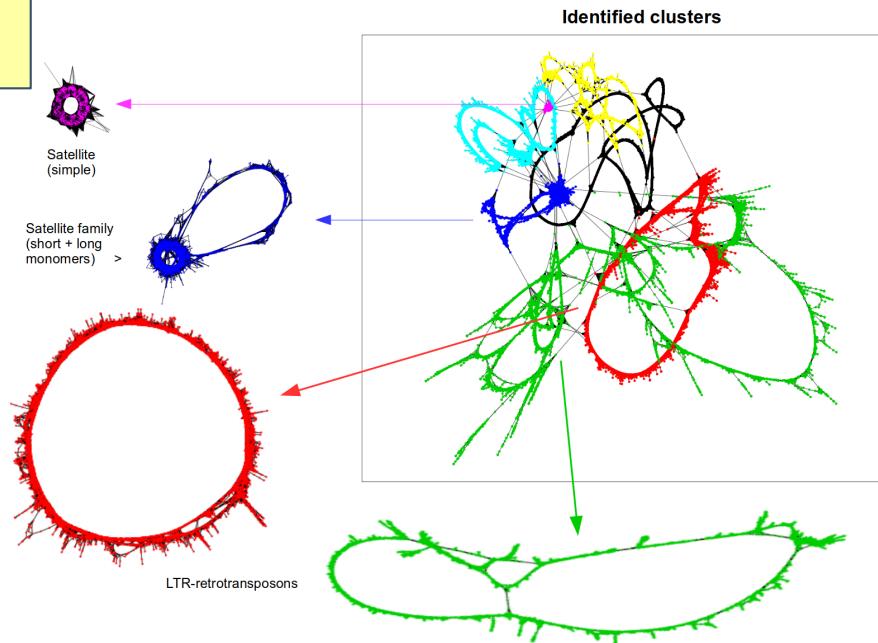


Principles and history of RepeatExplorer

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

First paper
on repeat
clustering
from NGS
data

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

*command-
line
version*

FIRST
WORKSHOP !

Principles and history of RepeatExplorer

2007

2010

2013

First paper on repeat clustering from NGS data

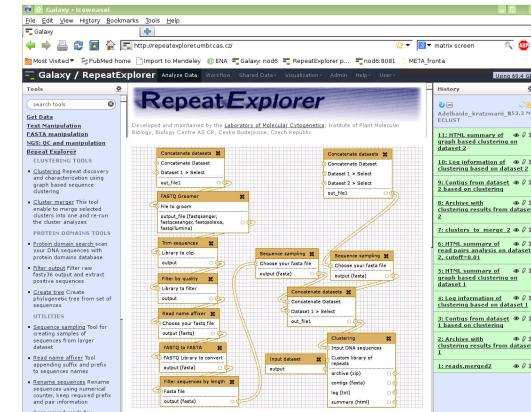
Introduction of graph- based clustering (Novak et al. 2010)

Repeat Explorer in Galaxy (Novak et al. 2013)

command-line version



Public web-based server



Principles and history of RepeatExplorer

2007

...

2010

...

2013

2014

...

2016

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)



ELIXIR \$\$\$

*Public web-
based
server*

Principles and history of RepeatExplorer

2007

...

2010

...

2013

2014

...

2016

2017

First paper
on repeat
clustering
from NGS
data

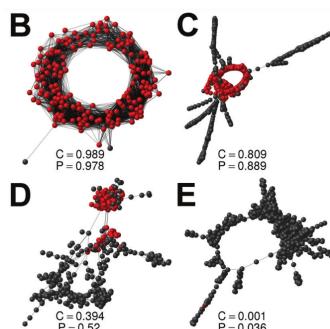
Introduction
of graph-
based
clustering
(Novak et
al. 2010)

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



ELIXIR \$\$\$

TAREAN

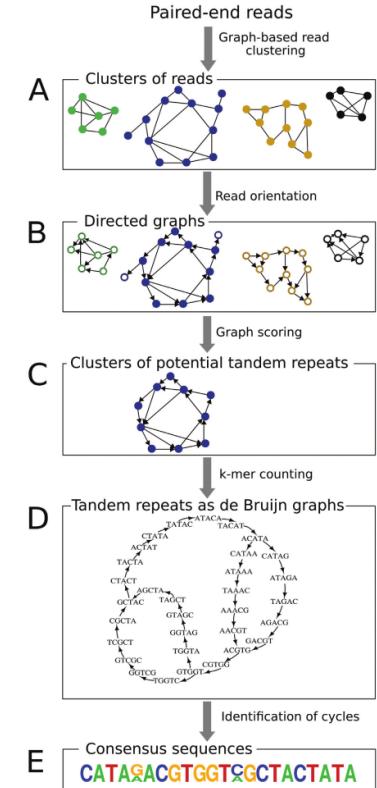


Nucleic Acids Research, 2017 1
doi: 10.1093/nar/gkx257

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

2007

...

2010

...

2013

2014

...

2016

2017

2018

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)



ELIXIR \$\$\$

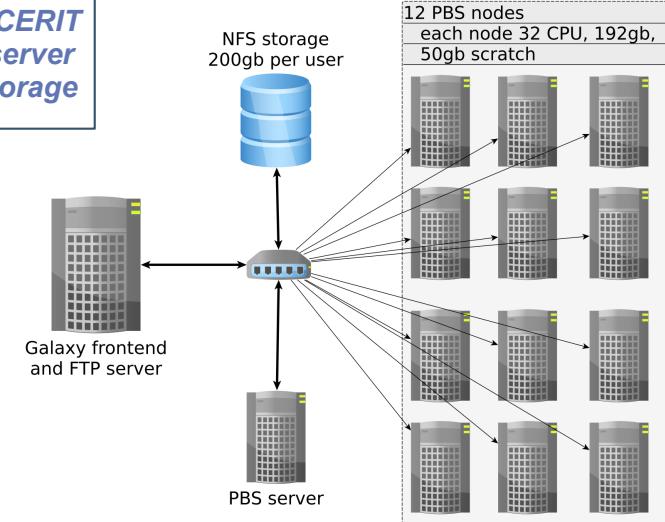
Public web-
based
server

TAREAN



**ELIXIR / CERIT
Galaxy server
+ data storage**

The screenshot shows the RepeatExplorer tool within a Galaxy interface. The main title is "RepeatExplorer" with the subtitle "Discover repeats in your next generation sequencing data". Below this is a magnifying glass icon over a network graph. A detailed description follows: "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." It also mentions the ELIXIR-CZ project support and a contact link for help. At the bottom, there are links for "RepeatExplorer", "Wiki", and "Registration".



Principles and history of RepeatExplorer

2007

...

2010

...

2013

2014

...

2016

2017

2018

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**

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 !



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](#).



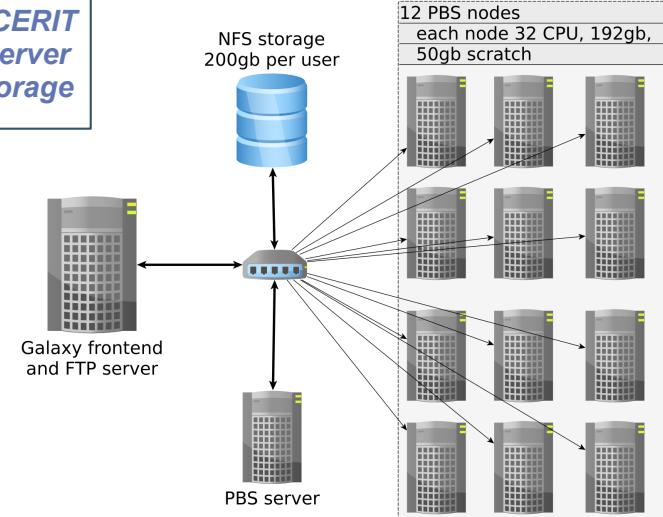
Go to Galaxy RepeatExplorer portal.



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



Please, read our registration manual.



Principles and history of RepeatExplorer

2007

...

2010

...

2013

2014

...

2016

2017

2018

2019

First paper on repeat clustering from NGS data

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

Repeat Explorer in Galaxy
(Novák et al. 2013)



ELIXIR \$\$\$

TAREAN

ELIXIR / CERIT
Galaxy server + data storage

REXdb database
(Neumann et al. 2019)

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

RESEARCH

Mobile DNA

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

2007

...

2010

...

2013

2014

...

2016

2017

2018

2019

2020



First paper on repeat clustering from NGS data

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

Repeat Explorer in Galaxy
(Novák et al. 2013)

TAREAN

REXdb
(Neumann et al. 2019)

Repeat Explorer ver. 2

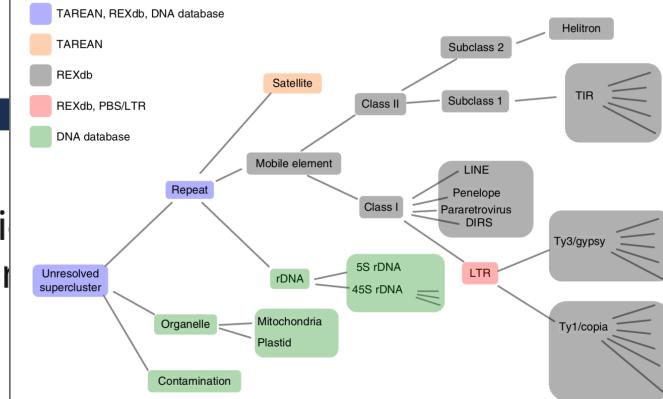
Neumann et al. Mobile DNA (2019) 10:1
<https://doi.org/10.1186/s13100-018-0144-1>

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

Decision tree for automatic annotation



Principles and history of RepeatExplorer

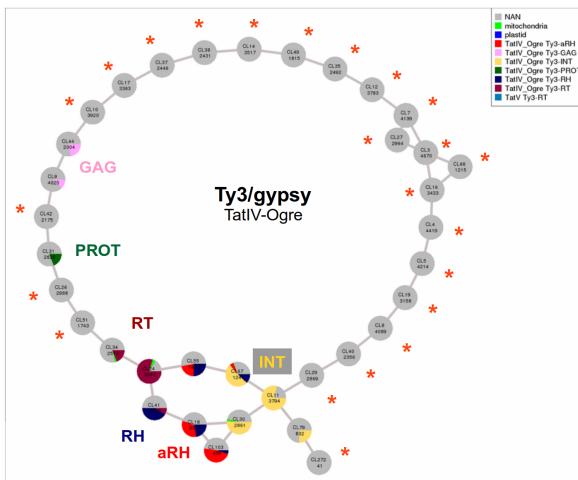
2007 ... 2010 ... 2013 2014 ... 2016 2017 2018 2019 2020

TAREAN

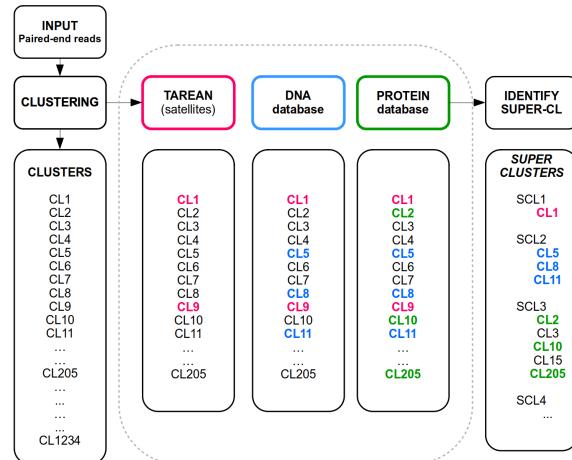
REXdb
(Neumann et al. 2019)

Repeat Explorer ver. 2

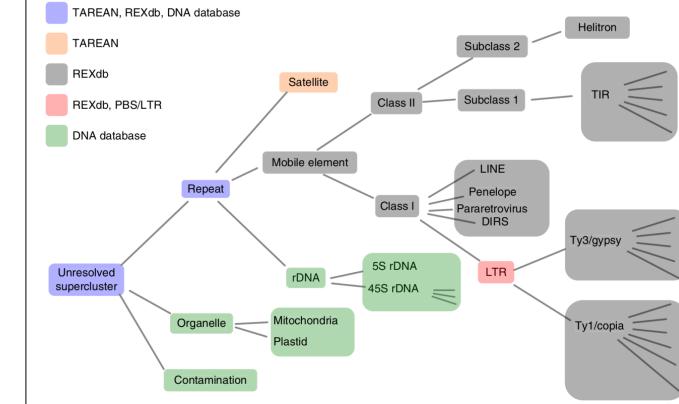
Superclusters provide more complete annotation



RepeatExplorer 2 – automatic detection of superclusters

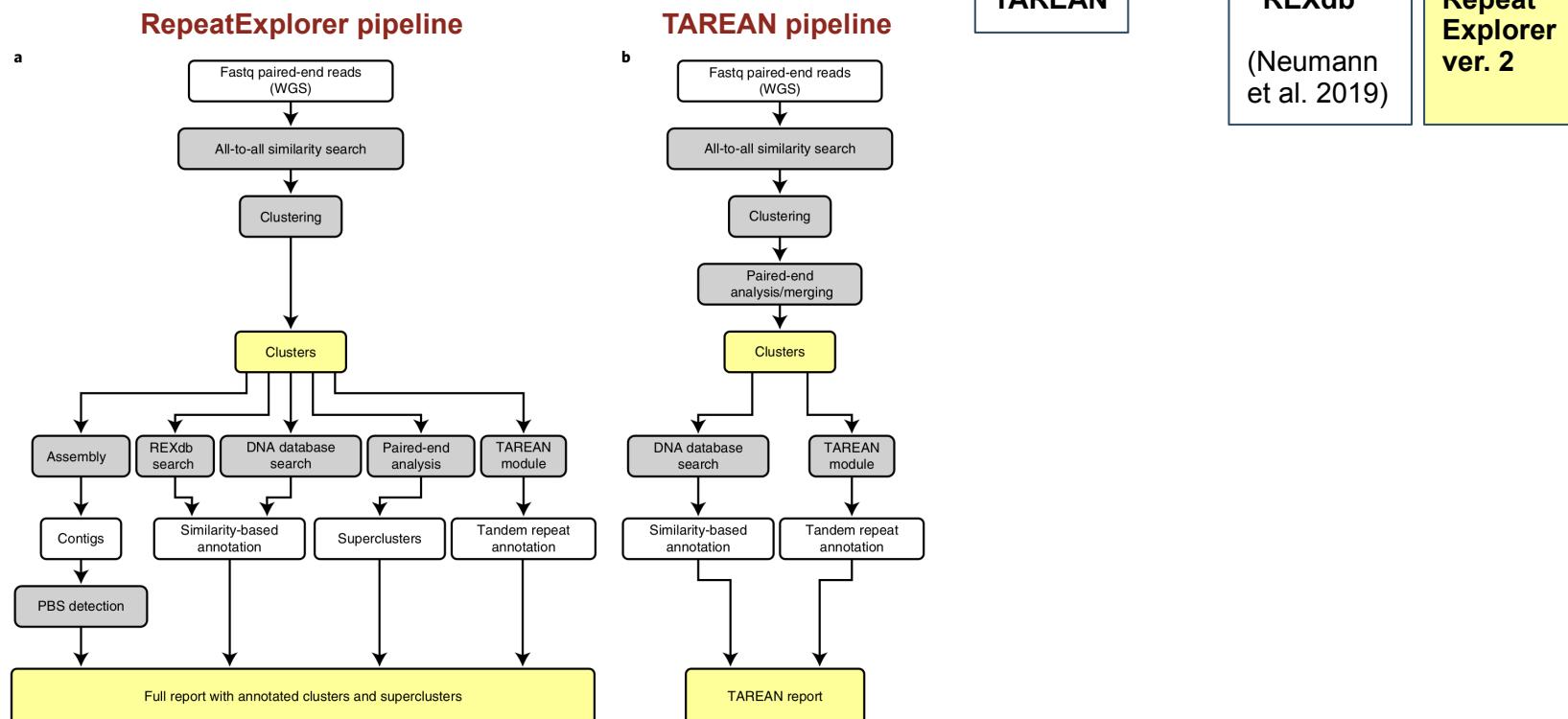


Decision tree for automatic annotation



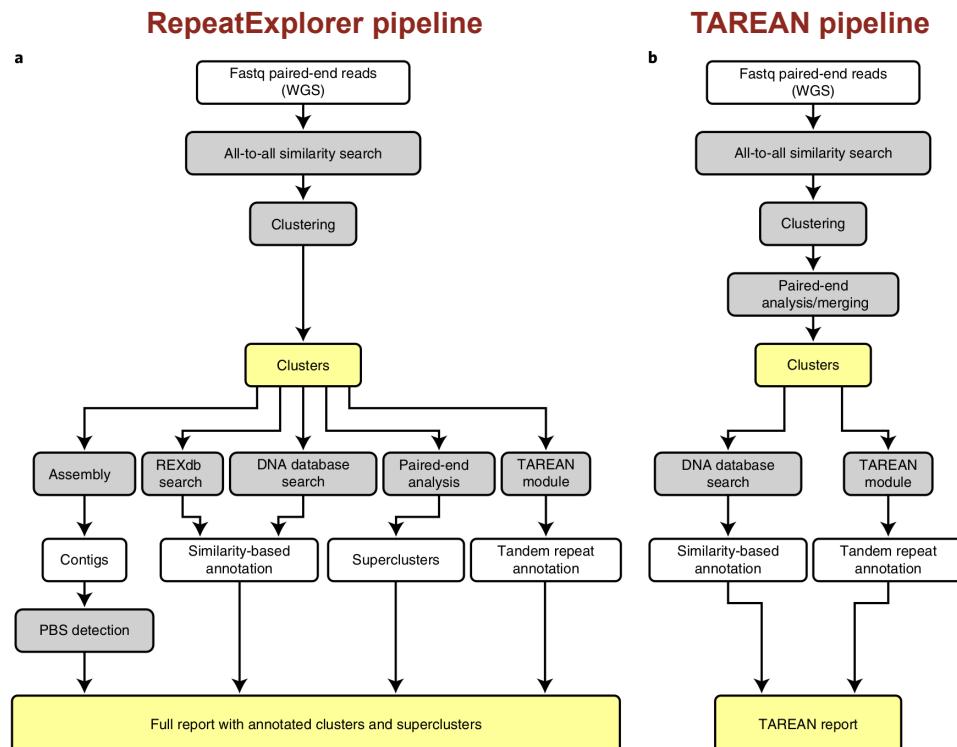
Principles and history of RepeatExplorer

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



Principles and history of RepeatExplorer

2007 ... 2010 ... 2013 2014 ... 2016 2017 2018 2019 2020 2021 2022



TAREAN

REXdb
(Neumann
et al. 2019)

Repeat
Explorer
ver. 2

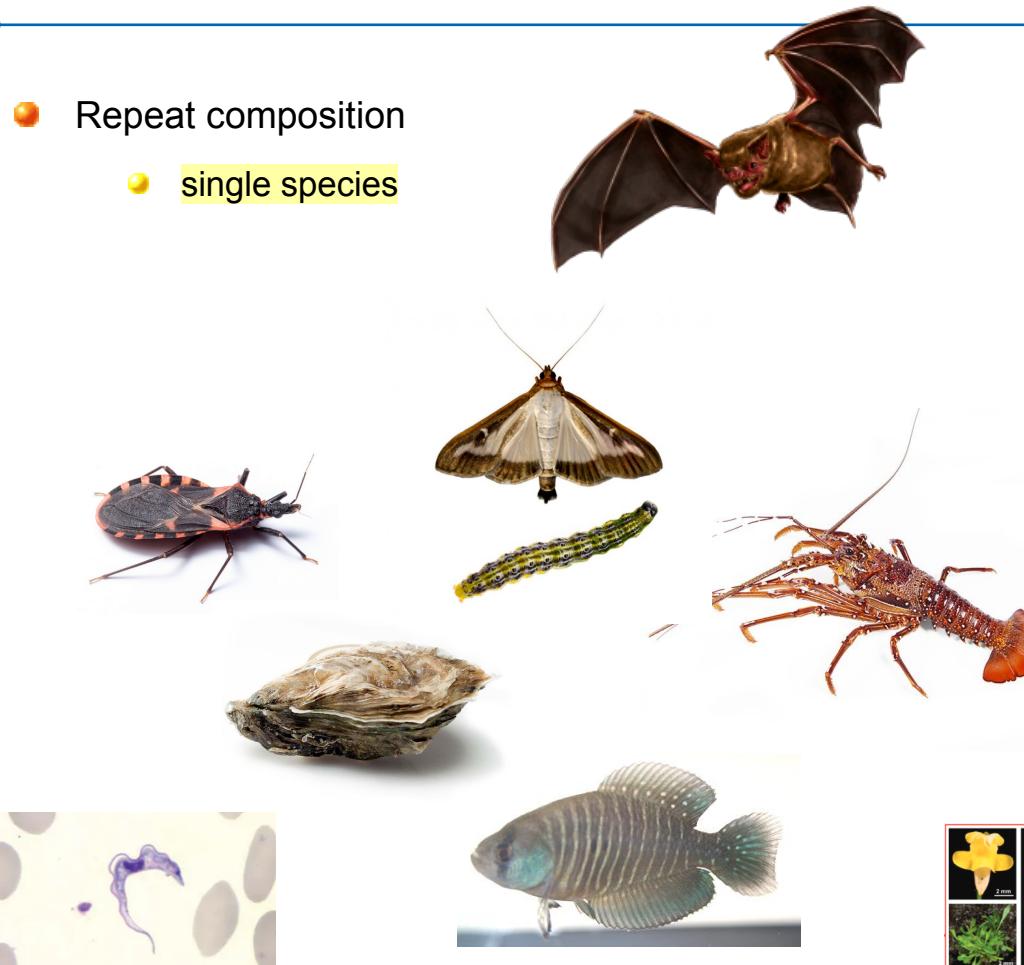
Additional tools:

- ChIP-seq Mapper
 - DANTE
 - Long reads
- Assembly annotation

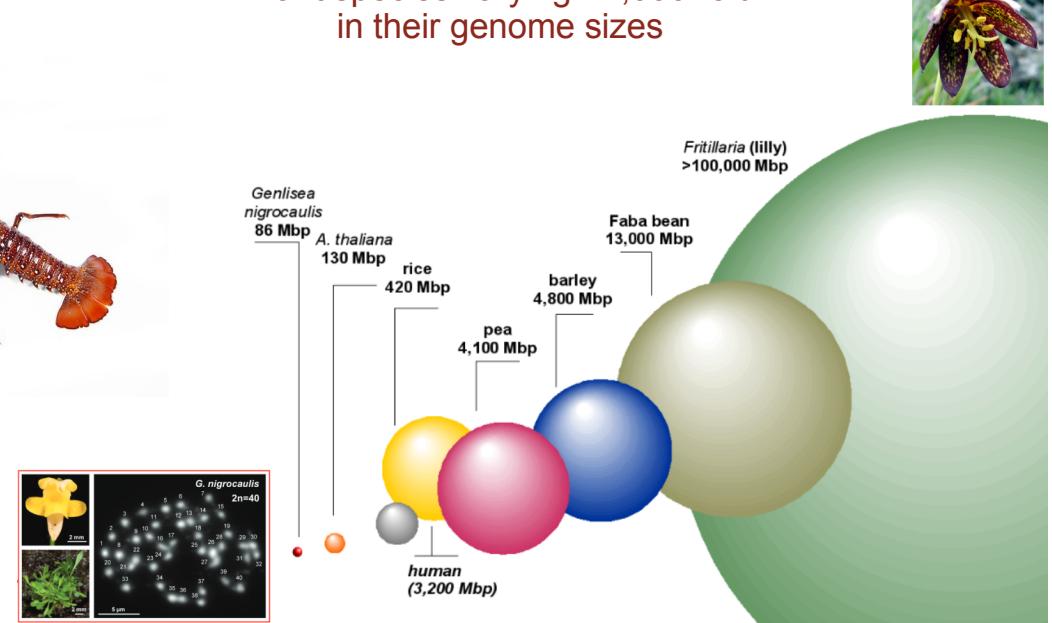
Applications

Applications

- Repeat composition
 - single species



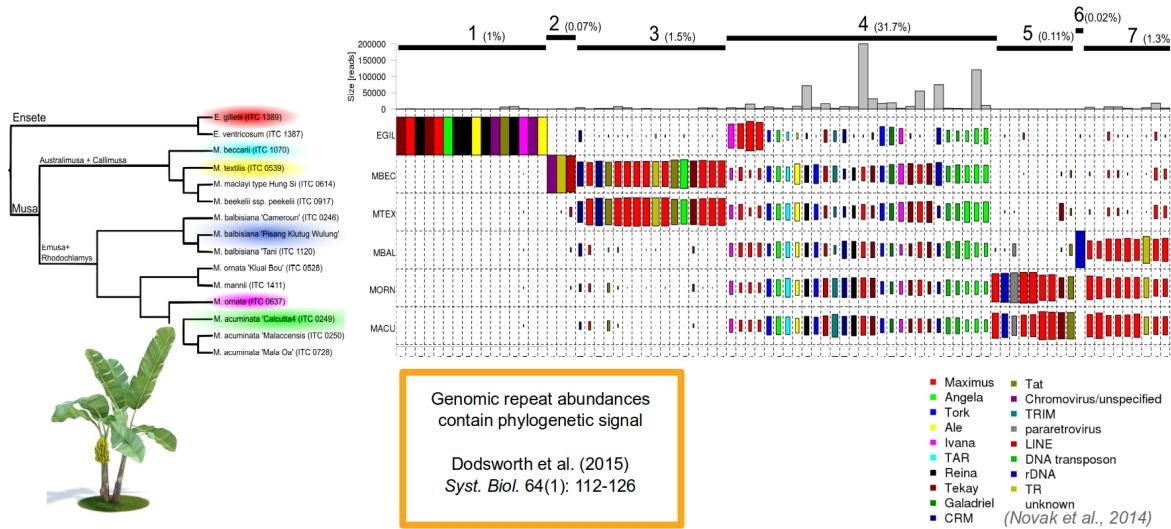
Plant species varying ~2,000-fold
in their genome sizes



Applications

Repeat composition

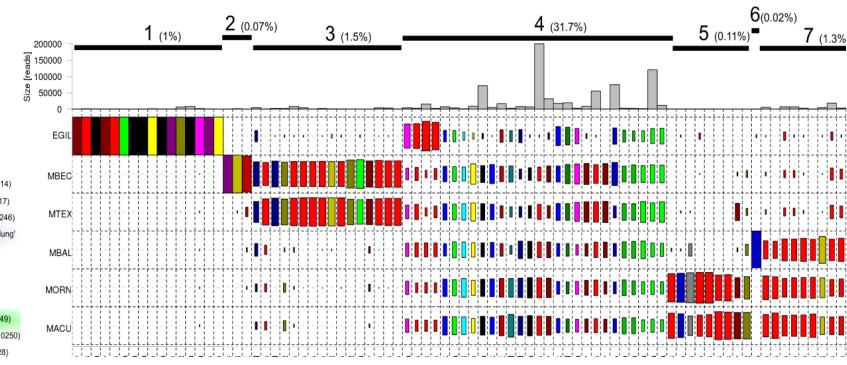
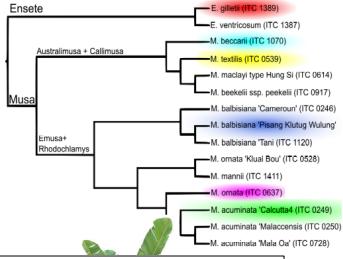
- single species
- comparative analysis



Applications

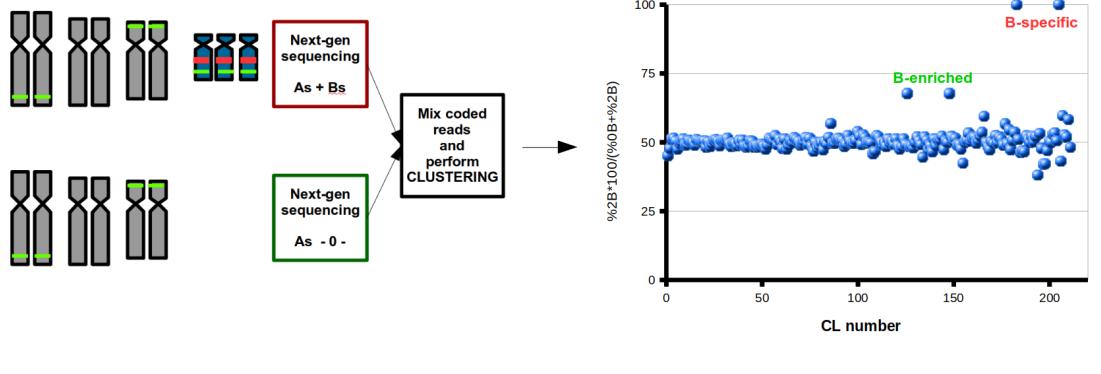
Repeat composition

- 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

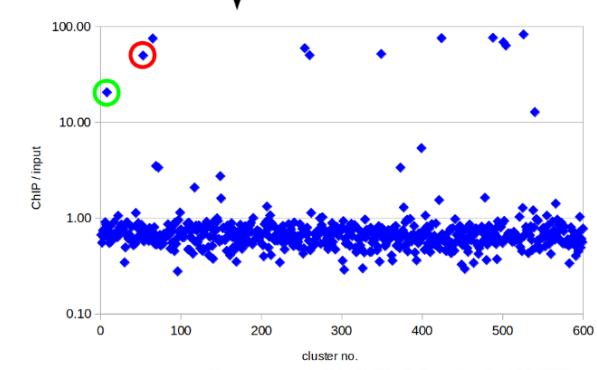
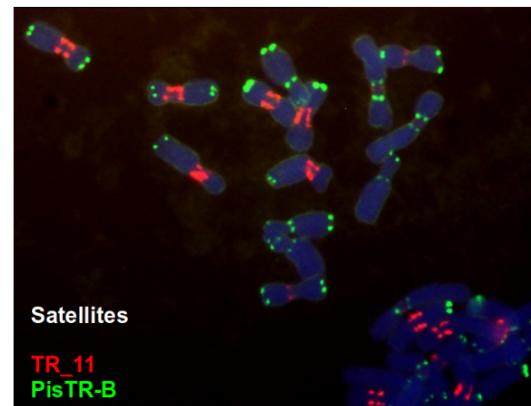
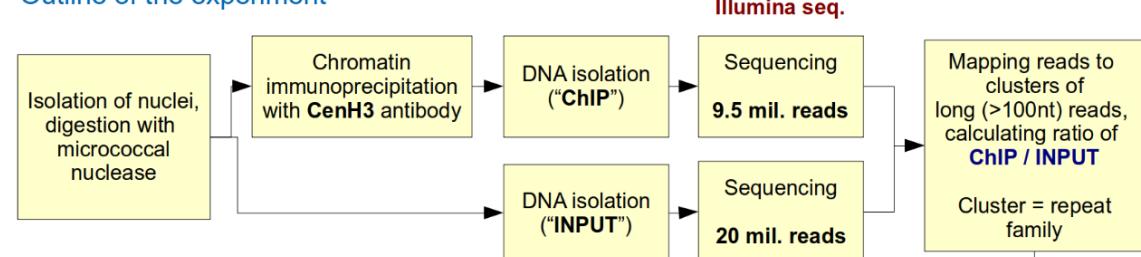
Legend:
Maximus, Angela, Chromovirus/unspecified, Tat, TRIM, pararetrovirus, Ale, LINE, TAR, DNA transposon, Reina, rDNA, Tekay, Galadriel, CRM
(Novak et al., 2014)

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



Neumann et al. (2012) PLoS Genetics 8:e1002777

Applications

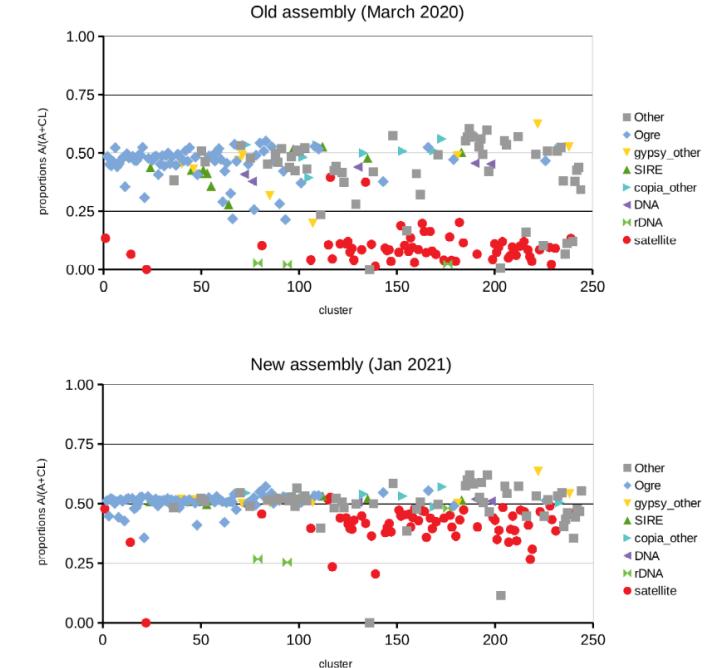
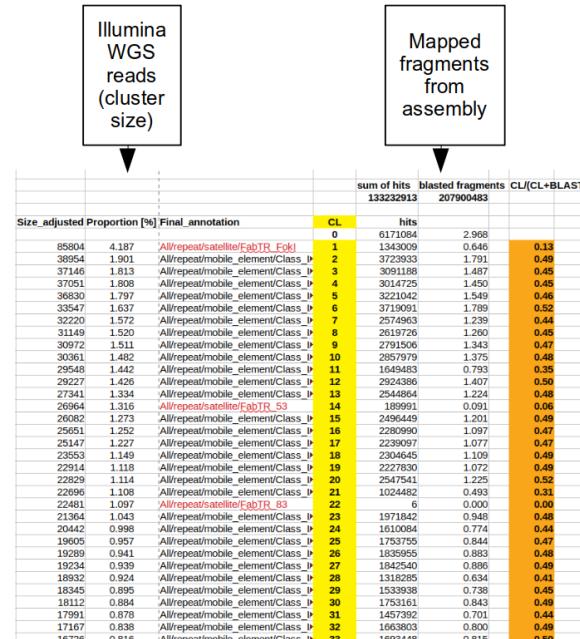
Repeat composition

- 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
 - reference databases for repeat annotation

Applications

- Repeat composition

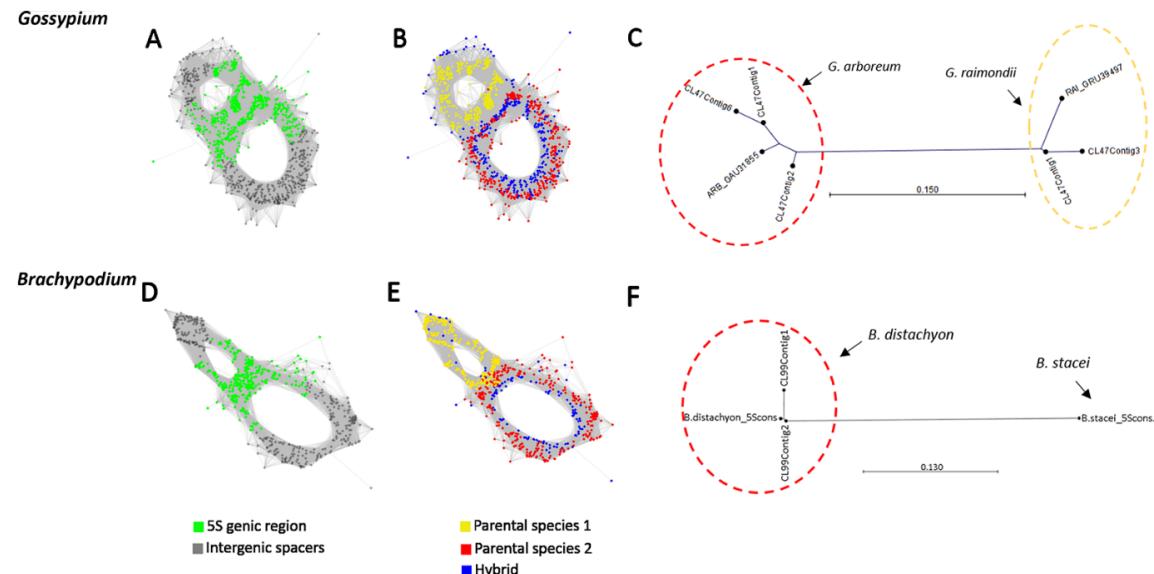
- single species
 - comparative analysis

- Repeat clusters as a reference

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

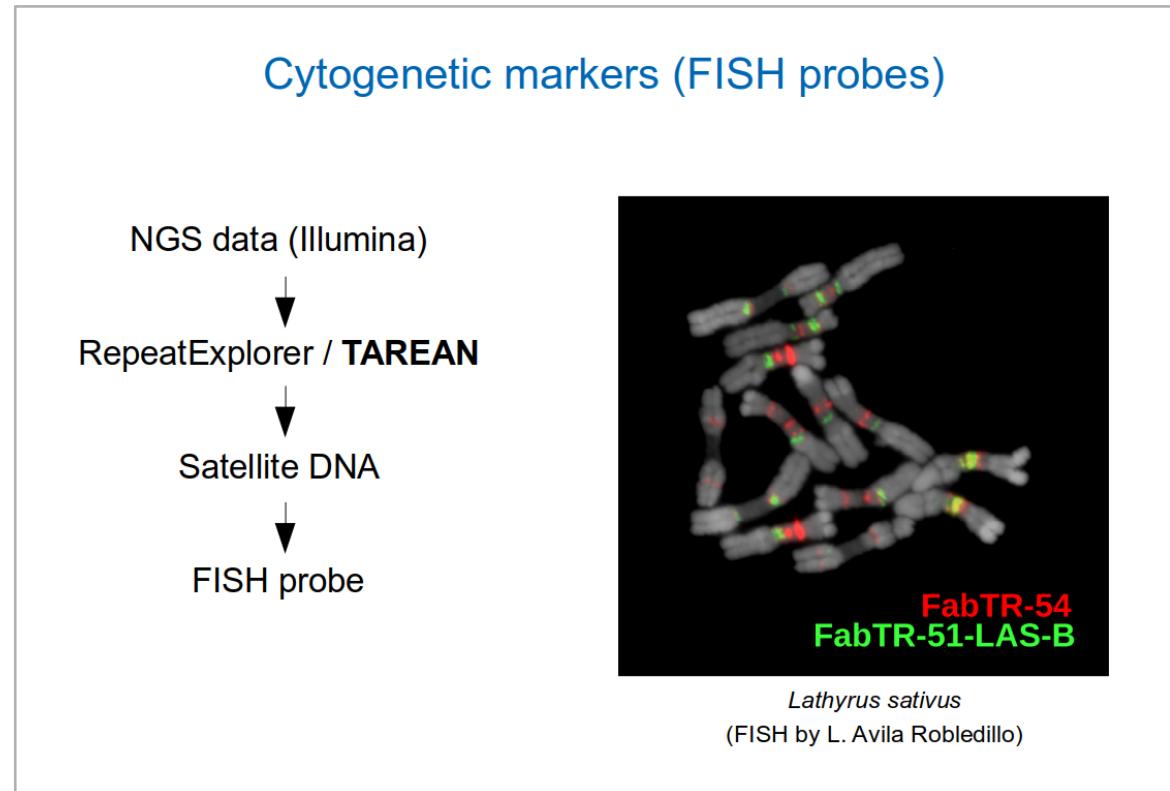
The Utility of Graph Clustering of 5S Ribosomal DNA Homoeologs in Plant Allopolyploids, Homoploid Hybrids, and Cryptic Introgressants

Sònia Garcia^{1,2}, Jonathan F. Wendel³, Natalia Borowska-Zuchowska⁴, Malika Aïnouche⁵,
Alena Kuderová² 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



Applications

- Repeat composition

- single species
 - comparative analysis

- Repeat clusters as a reference

- ChIP-seq
 - assembly
 - graph shapes

- Satellite DNA

- cytogenetic studies
 - diagnostic markers

The screenshot shows a webpage from The Marine Mammal Center. At the top, there's a navigation bar with links for Animal Care, Science & Conservation, Education, Get Involved, a DONATE button, and a search icon. The main content area features a large image of a sea lion swimming in water. To the left of the image, there's a purple sidebar with a back-to-all-publications link and the title "New Technique for Diagnosing Lung Parasite Infections in Sea Lions". Below the title, there's a brief description of the research: "A novel quantitative real-time PCR diagnostic assay for fecal and nasal swab detection of an otariid lungworm, *Parafilaroides decurvis*". At the bottom of the sidebar, there are two buttons: "Lungworm" and "Parasites".



REVIEW

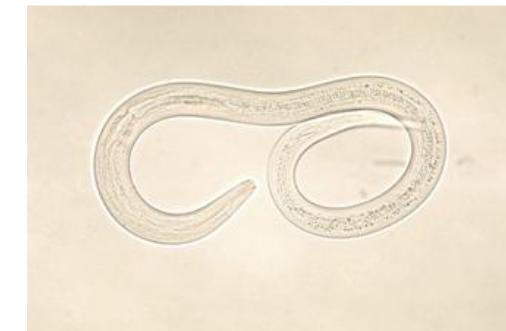
published: 23 September 2019
doi: 10.3389/fgene.2019.00883



A Case for Using Genomics and a Bioinformatics Pipeline to Develop Sensitive and Species-Specific PCR-Based Diagnostics for Soil-Transmitted Helminths

Jessica R. Grant^{1*}, Nils Pilote^{1,2} and Steven A. Williams^{1,2}

¹ Biological Sciences, Smith College, Northampton, MA, United States, ² Molecular and Cellular Biology, University of Massachusetts, Amherst, MA, United States



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

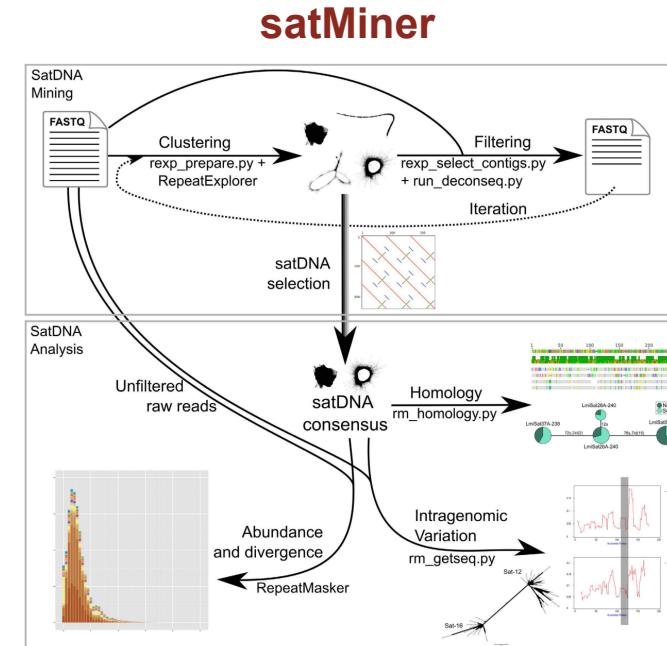
SCIENTIFIC REPORTS

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, Josefina Cabrero & Juan Pedro M. Camacho



Applications

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

BMC Bioinformatics

● Repeat composition

- 🟡 single species
- 🟡 comparative analysis

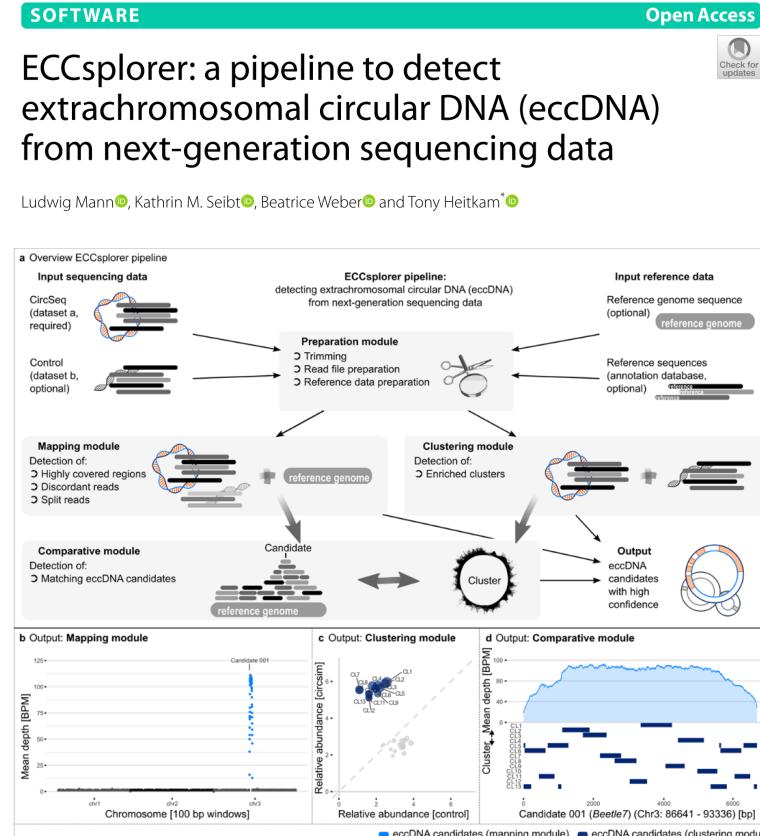
● Repeat clusters as a reference

- 🟡 ChIP-seq
- 🟡 assembly
- 🟡 graph shapes

● Satellite DNA

- 🟡 cytogenetic studies
- 🟡 diagnostic markers

● RE utilization in other pipelines



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



Enjoy the workshop !