

# 6<sup>th</sup> Workshop on the Application of Next Generation Sequencing to Repetitive DNA Analysis in Plants

Welcome !



*Institute of Plant Molecular Biology,  
Biology Centre ASCR  
České Budějovice, Czech Rep.*



# Programme

## *Tuesday (May 23)*

9:30 - Principles and applications of graph-based repeat clustering (J. Macas)

10:10 - RepeatExplorer pipeline version 2.0 (P. Novák)

10:30 - 11:00 *Coffee break*

11:00 - Using RepeatExplorer output for repeat annotation and quantification (J. Macas)

11:20 - Transposon protein databases (P. Neumann)

11:40 - New tool: Detection and annotation of transposon protein domains (N. Hošťáková)

12:00 - 13:30 *Lunch*

13:30 – 14:00 - Steven Dodsworth (Kew, UK) - Phylogenetic signal in repeat abundances: Angiosperm examples from tomatoes to orchids

### 14:00 - (18:00) Practical training I

- design of sequencing and repeat analysis experiments
- introduction to Galaxy environment
- quality control and pre-processing of NGS reads, dealing with various read formats
- setting up clustering analysis
- comparative clustering of multiple samples

**19:00 → : Dinner at “CITYgastro”**

## *Wednesday (May 24)*

9:00 - 12:30

- Gustavo Souza – Using genomic repeat abundance and cytogenomic approaches to infer phylogenetic relationships in *Caesalpinia sensu lato* (Fabaceae)
- Maria Gonzalez – Chromosome evolution of South American and Antarctic species of *Deschampsia* (Poaceae)
- Beatrice Weber – Chromoviruses in the genome of sugar beet *Beta vulgaris*
- Danijela Greguraš – Repeatome dynamics in the earliest evolutionary stages of apomictic plants
- Alevtina Ruban – Why does the genome size differ between roots and shoots in some *Aegilops speltoides* plants?
- Nusrat Sultana – Bioinformatics and molecular characterization of *Vaccinium corymbosum* genome
- Christiaan Henkel – Can we sequence a repeat-rich, 35 Gbp tulip genome?

12:30 - 13:30 *Lunch*

### 13:30 - (18:00) Practical training II

- identification of satellite DNA using TAREAN
- understanding RepeatExplorer output
- cluster annotation and repeat composition of the genome
- comparative clustering of multiple species – data interpretation
- repeat quantification (principles, sensitivity and reproducibility)
- design of hybridization probes based on RE

## *Thursday (May 25)*

9:00 – 12:30

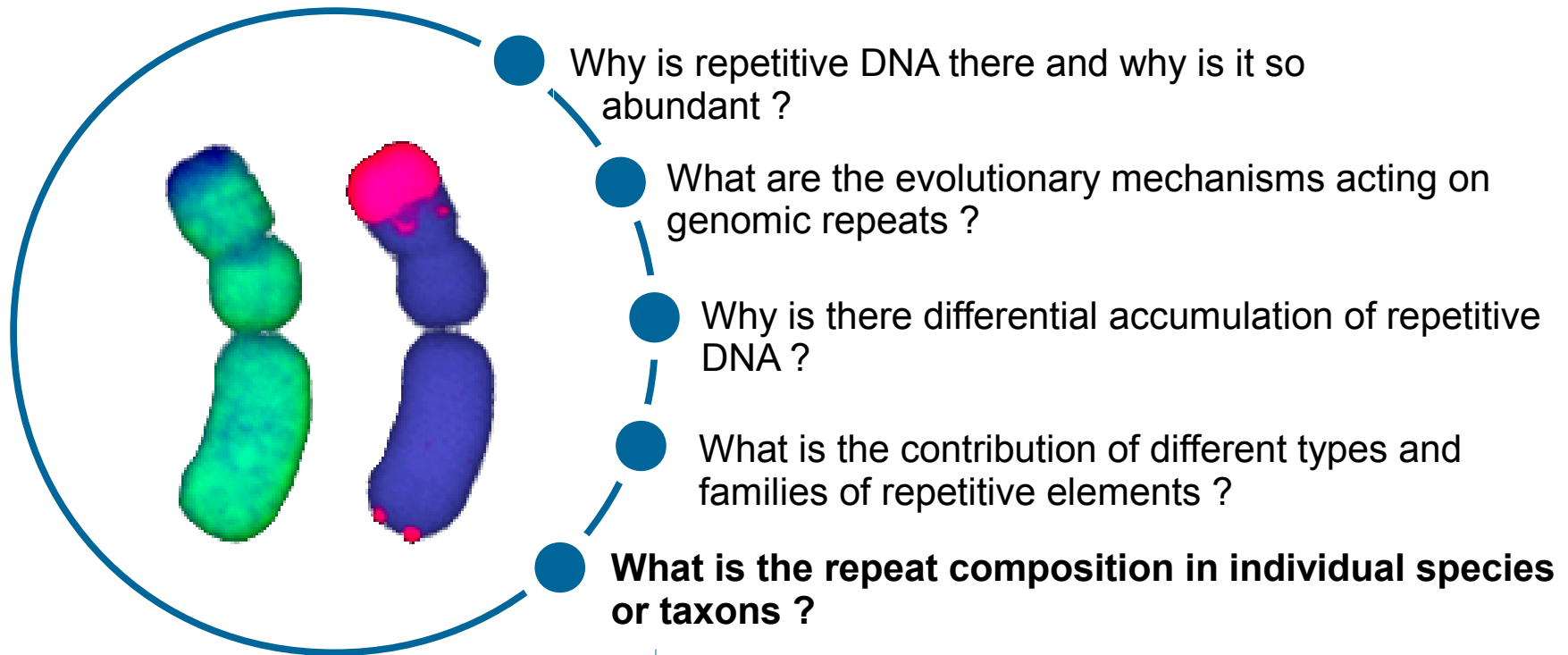
- Amanda Grusz – Genome evolution in the fern family *Pteridaceae*
- Andrew Leitch – The placement of Gnetales amongst seed plants
- Sonia Garcia, Daniel Viales – Concerted evolution under the microscope: rDNA arrangements in three Asteraceae genera
- Ales Kovarik – Higher-order repeat structure of 5S rRNA genes of *Esox lucius* (fish) determined from long PacBio reads
- Tanja Vojvoda Zeljko – Characterization of transposable elements containing internal tandem repeats in the genome of the Pacific oyster *Crassostrea gigas*
- Rodolpho Menezes – Cytogenetics meets phylogeography and phylogenomics: exploring the evolutionary history of Neotropical swarm-founding social wasps
- Abhijeet Shah – Mobile DNA in Acrididae grasshoppers

12:30 - 13:30 *Lunch*

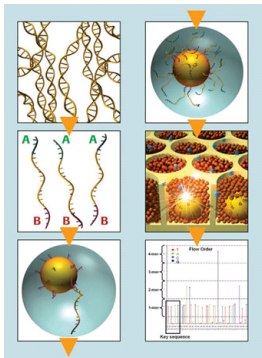
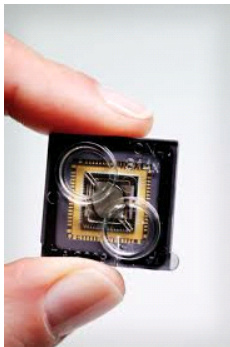
### 13:30 – (18:00) - Practical training III

- combining repeat clustering with ChIP-seq data
- identification and phylogenetic analysis of retrotransposon protein domains
- SeqGrappleR – visualization and annotation of the cluster graphs
- advanced topics, troubleshooting

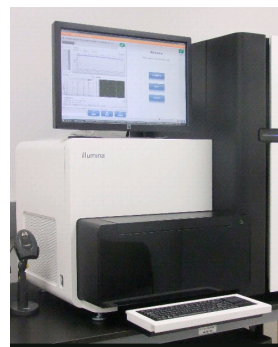
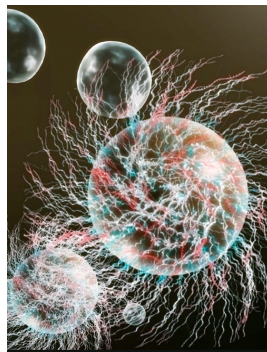
# The challenge of repeat identification in complex genomes



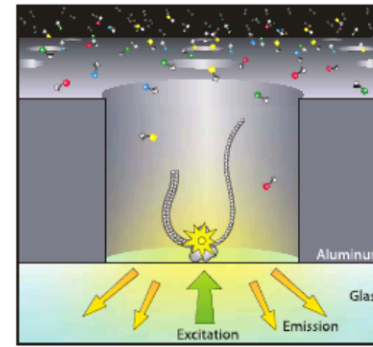
***Next generation sequencing: getting sequence data is no longer a limiting factor***



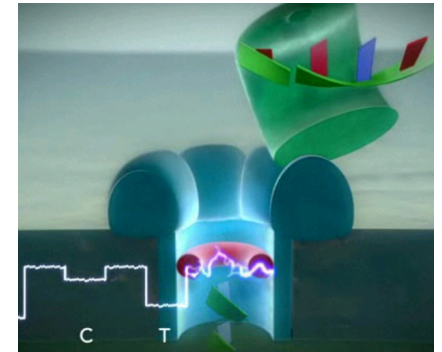
454/Roche



Illumina



Pacific Biosciences



Oxford Nanopore

# Principles and tools for global repeat identification

		Assembled genome or long contigs (BACs etc.)	
		Y E S	N O
Reference database	Y E S	<div>RepeatMasker / RepBase</div> <ul style="list-style-type: none"><li>• <i>first “model” species</i></li><li>• <i>extensive resources (+ wet lab)</i></li></ul>	<ul style="list-style-type: none"><li>• <i>model-related taxa</i></li><li>• <i>often shotgun, “assembled”</i></li></ul>
	N O		

## Repeat identification in sequence data

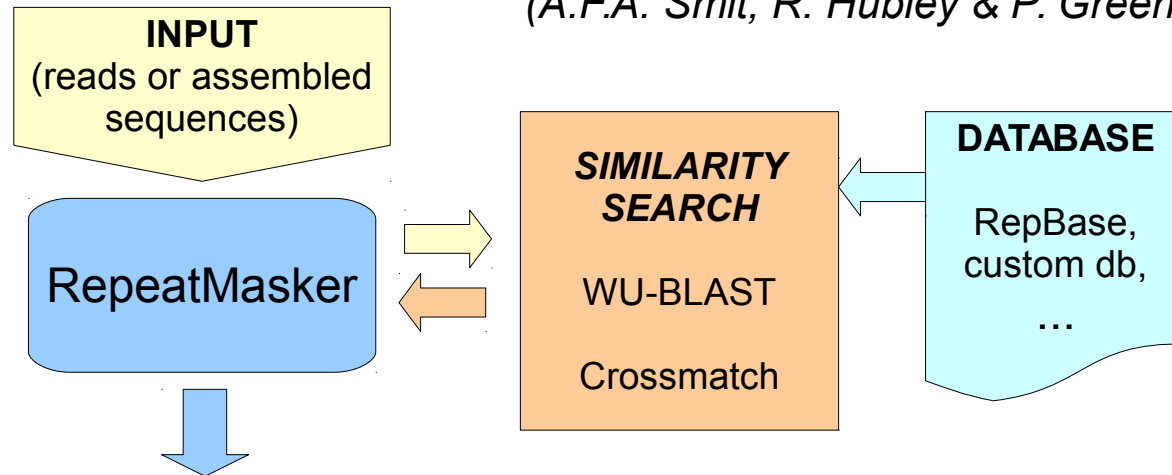
- Pure algorithmic (evaluates “repetitiveness” of substrings in DNA sequence)
- Signature or feature-based (searches for specific structures of biological importance)
- Similarity to known repeats (uses curated databases of known repeats)



## Principles and tools for global repeat identification

# RepeatMasker

(A.F.A. Smit, R. Hubley & P. Green)



## OUTPUT

## table of hits

SW score	perc div.	perc del.	perc ins.	query sequence	position begin	position end	query (left)	matching repeat	repeat class/family
459	20.0	1.1	0.0	LAL_1000603f	2	96	(4) +	Ogre_PS1_from_AY299398	LTR/Gypsy/Ogre
361	27.8	0.0	0.0	LAL_1000670r	1	90	(10) C	Ogre_PA_z_c700	LTR/Gypsy/Ogre_PA
449	23.5	0.0	0.0	LAL_1001438f	1	98	(2) +	PSC454_CL1Contig253	LTR/Gypsy/Ogre
491	23.5	0.0	0.0	LAL_1001438r	3	100	(0) C	PSC454_CL1Contig253	LTR/Gypsy/Ogre
487	24.0	0.0	0.0	LAL_1001478f	1	100	(0) C	PSC454_CL1Contig253	LTR/Gypsy/Ogre
398	28.0	0.0	0.0	LAL_1001478r	1	100	(0) C	PSC454_CL1Contig2773	LTR/Gypsy/Ogre
417	20.6	2.0	2.0	LAL_1002130r	2	100	(0) C	PSC454_CL1Contig253	LTR/Gypsy/Ogre
373	29.0	0.0	0.0	LAL_1002357f	1	100	(0) +	Ogre_PS1_from_AY299398	LTR/Gypsy/Ogre

[illegible]

**< masked  
sequence**

```
summary >
table
```

	number of elements*	length occupied	percentage of sequence
Retroelements	1	12014 bp	96.40 %
SINES:	0	0 bp	0.00 %
Penelope	0	0 bp	0.00 %
LINES:	0	0 bp	0.00 %
CRE/SLACS	0	0 bp	0.00 %
L2/CR1/Rex	0	0 bp	0.00 %
R1/L0A/Jockey	0	0 bp	0.00 %
R2/R4/NeSL	0	0 bp	0.00 %
RTE/Bov-B	0	0 bp	0.00 %
L1/CIN4	0	0 bp	0.00 %
LTR elements:	1	12014 bp	96.40 %
BEL/Pao	0	0 bp	0.00 %
Ty1/Copia	0	0 bp	0.00 %
Gypsy/DIRS1	1	12014 bp	96.40 %
Retroviral	0	0 bp	0.00 %
DNA transposons	0	0 bp	0.00 %
hobo-Activator	0	0 bp	0.00 %
Tcl-IS630-Pogo	0	0 bp	0.00 %
En-Spm	0	0 bp	0.00 %
MuDR-IS905	0	0 bp	0.00 %
PiggyBac	0	0 bp	0.00 %
Tourist/Harbinger	0	0 bp	0.00 %
Other (Mirage, P-element, Transib)	0	0 bp	0.00 %

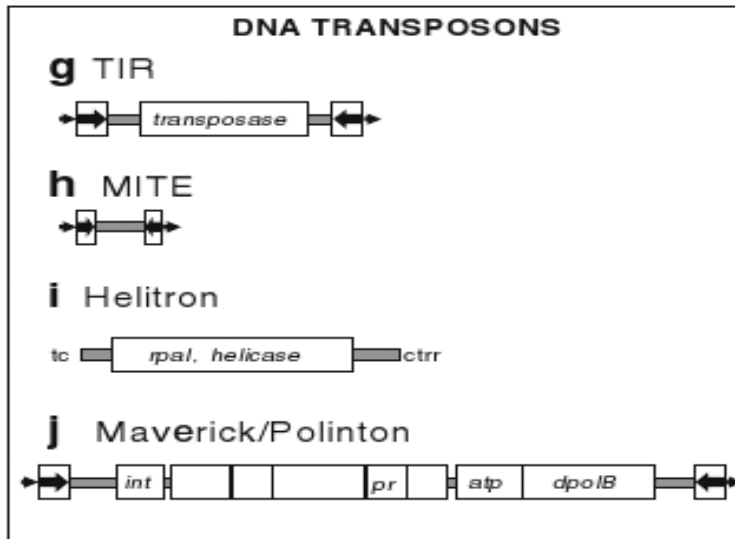
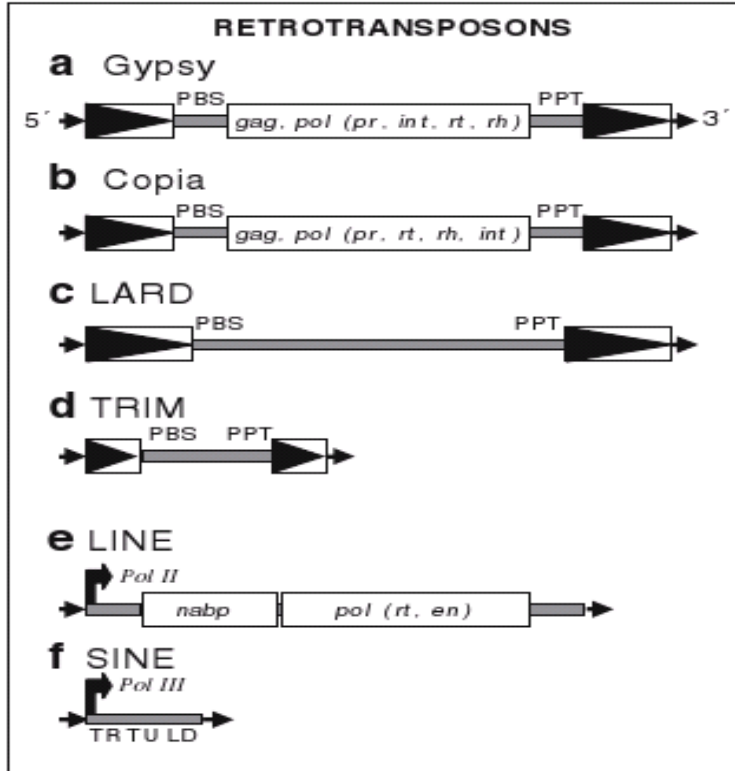
# Principles and tools for global repeat identification

		<b>Assembled genome or long contigs (BACs etc.)</b>	
		<b>Y E S</b>	<b>N O</b>
<b>Reference database</b>	<b>Y E S</b>	<div>RepeatMasker / RepBase</div> <ul style="list-style-type: none"> <li><i>first “model” species</i></li> <li><i>extensive resources (+ wet lab)</i></li> </ul>	<ul style="list-style-type: none"> <li><i>model-related taxa</i></li> <li><i>often shotgun, “assembled”</i></li> </ul>
	<b>N O</b>	<ul style="list-style-type: none"> <li><b>LTR_STRUC, MITE-Hunter , etc.</b></li> </ul>	

## Repeat identification in sequence data

- Pure algorithmic (evaluates “repetitiveness” of substrings in DNA sequence)
- **Signature or feature-based** (searches for specific structures of biological importance)
- Similarity to known repeats (uses curated databases of known repeats)

# Principles and tools for global repeat identification



## Signature / structure-based approaches

**LTR\_STRUC** (McCarthy and McDonald, 2003)

**LTR\_FINDER** (Xu and Wang, 2007)

**SINE-Finder** (Wenke et al. 2011)

**MITE-Hunter** (Han and Wessler, 2010)

**MITE Analysis Toolkit** (Yang and Hall, 2003)

**HelitronFinder** (Du et al. 2008) [HelA helitrons in maize]

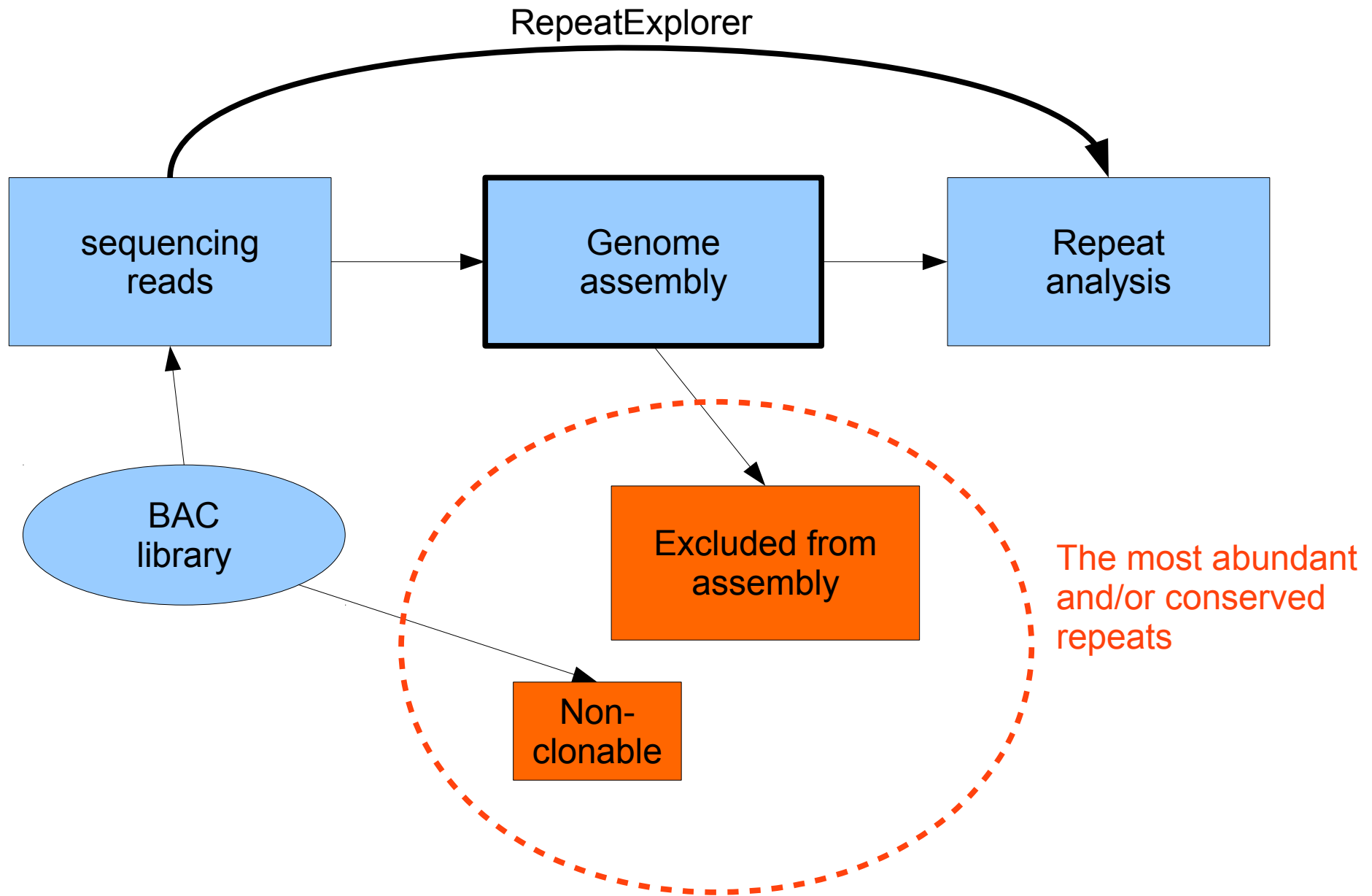
# Principles and tools for global repeat identification

		Assembled genome or long contigs (BACs etc.)	
		YES	NO
Reference database	YES	<div>RepeatMasker / RepBase</div> <ul style="list-style-type: none"><li>• <i>first “model” species</i></li><li>• <i>extensive resources (+ wet lab)</i></li></ul>	<ul style="list-style-type: none"><li>• <i>model-related taxa</i></li><li>• <i>often shotgun, “assembled”</i></li></ul>
	NO	<ul style="list-style-type: none"><li>• LTR_STRUC, MITE-Hunter , etc.</li><li>• <b>REPET</b></li></ul>	<p><i>(using NGS reads)</i></p> <ul style="list-style-type: none"><li>• direct assembly (phrap,...)</li><li>• <b>clustering-based</b></li></ul>

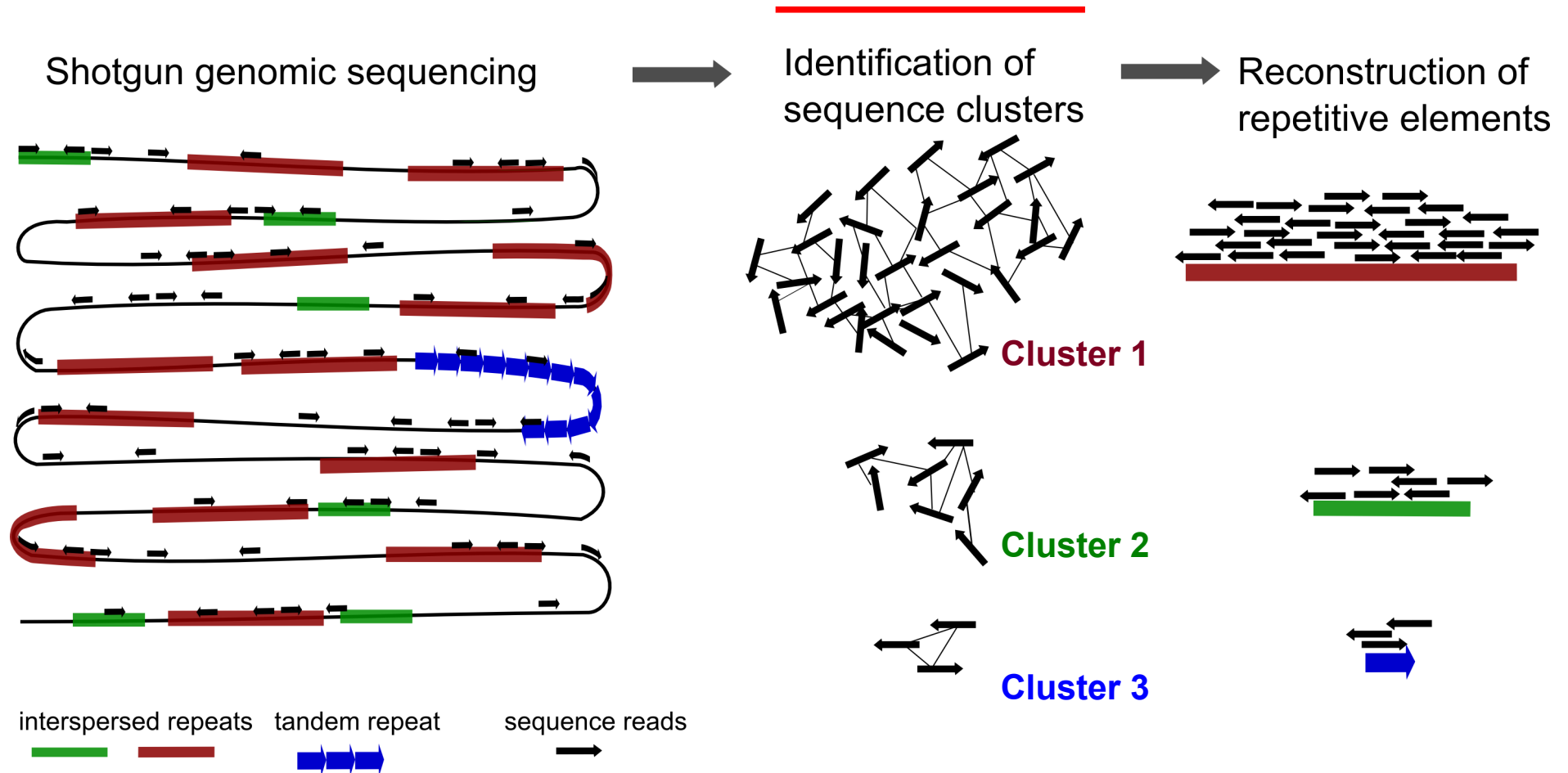
## Repeat identification in sequence data

- **Pure algorithmic (evaluates “repetitiveness” of substrings in DNA sequence)**
- Signature or feature-based (searches for specific structures of biological importance)
- Similarity to known repeats (uses curated databases of known repeats)

# Principles and tools for global repeat identification

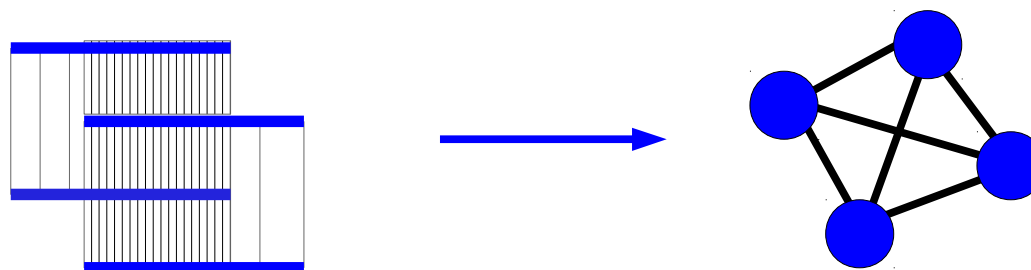
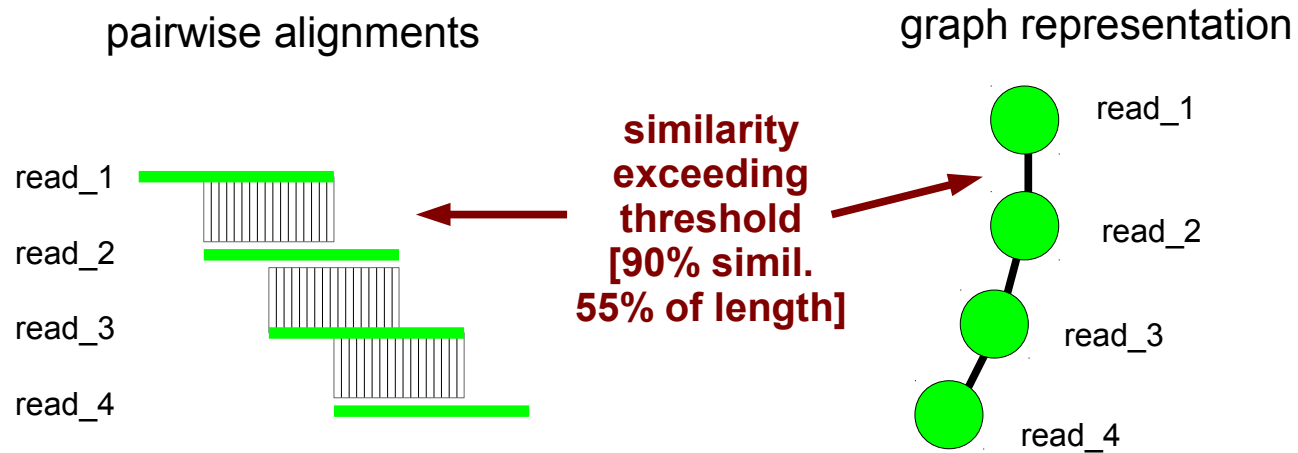


# Clustering-based repeat identification



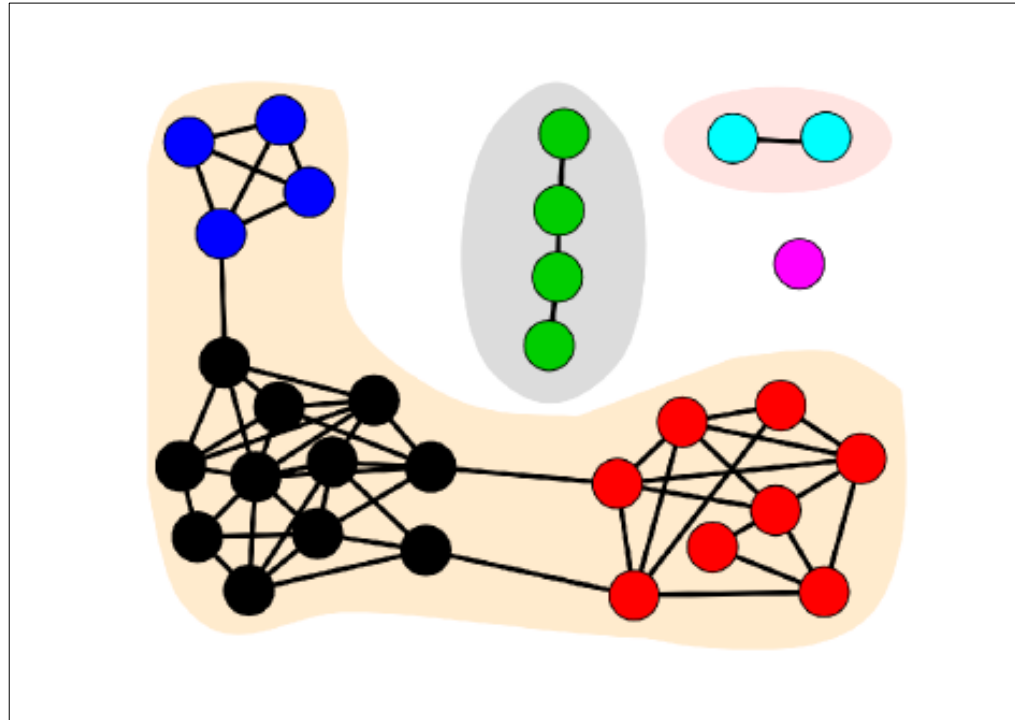
**CLUSTER = a set of frequently overlapping reads = REPEAT FAMILY**

# Clustering-based repeat identification





# Clustering-based repeat identification

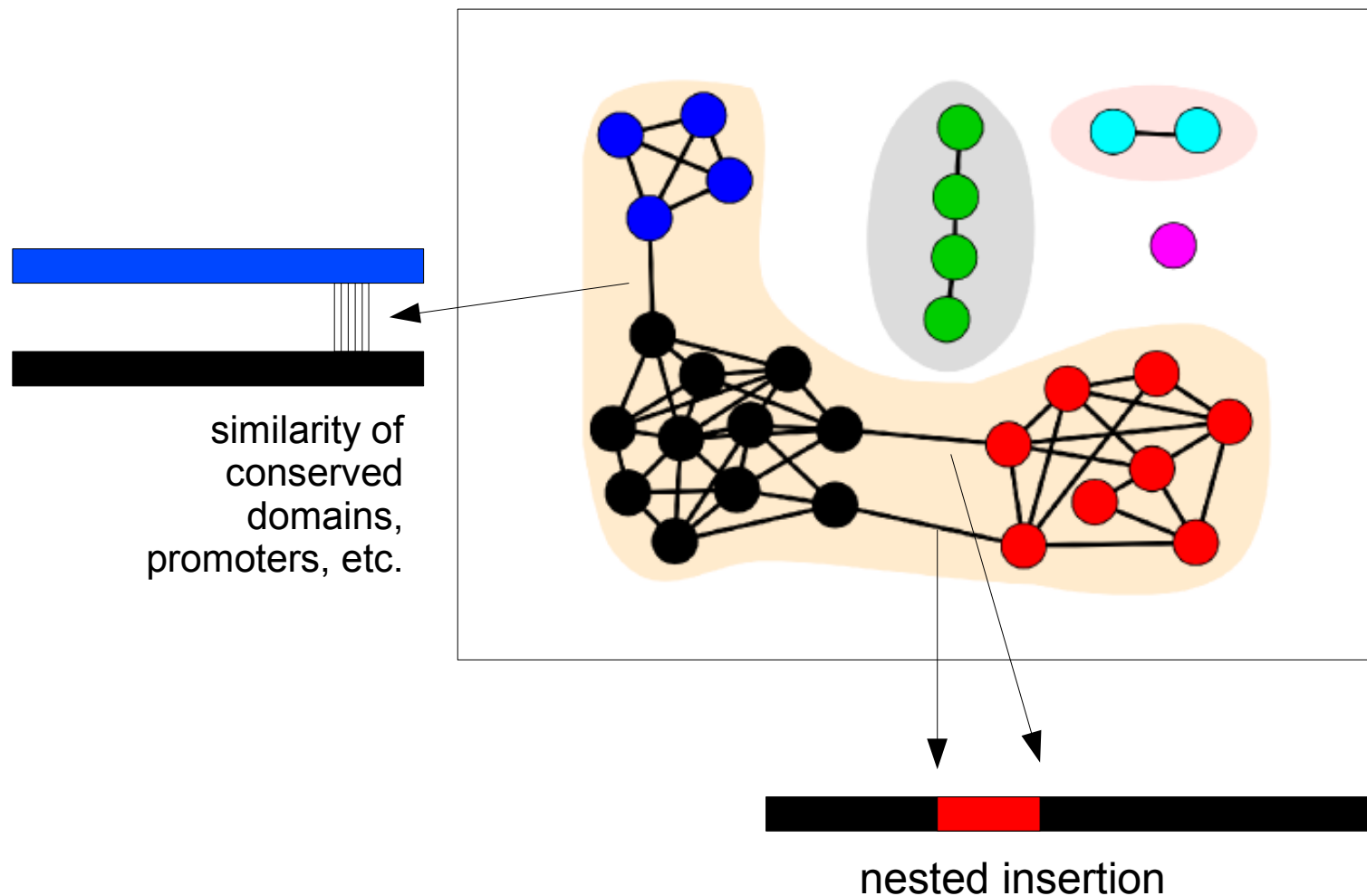


Single linkage clustering => connected components

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

Macas et al. (2007) - Repetitive DNA in the pea (*Pisum sativum* L.) genome: comprehensive characterization using 454 sequencing and comparison to soybean and *Medicago truncatula*.

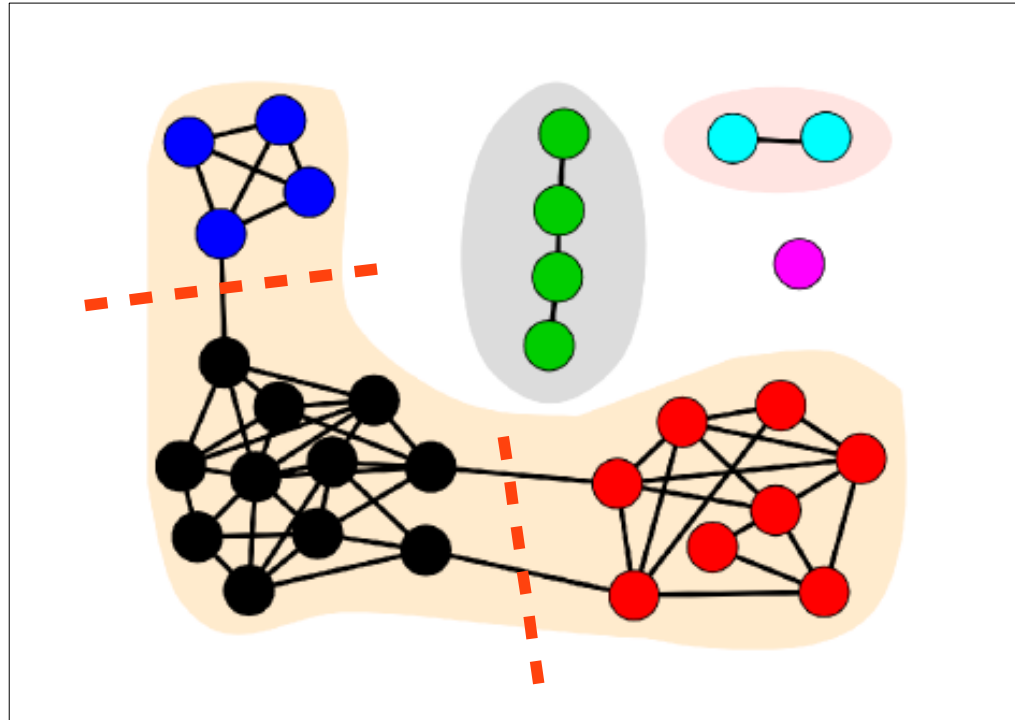
# Clustering-based repeat identification



Single linkage clustering => connected components

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

# Graph-based clustering

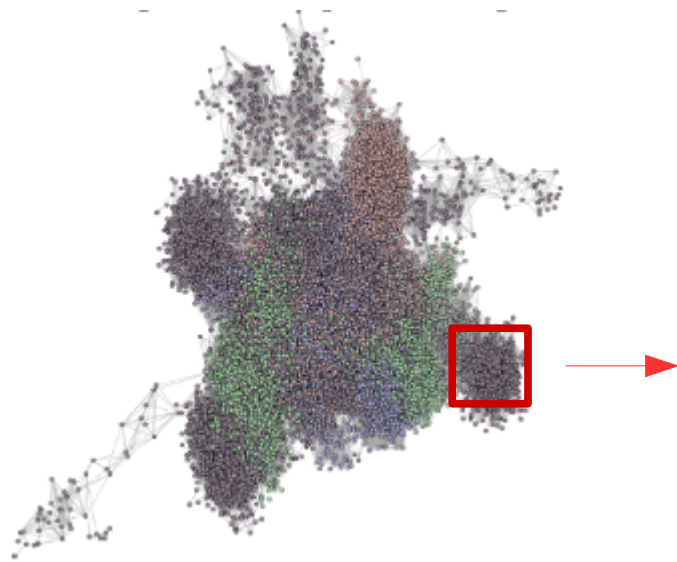


## Graph-based clustering

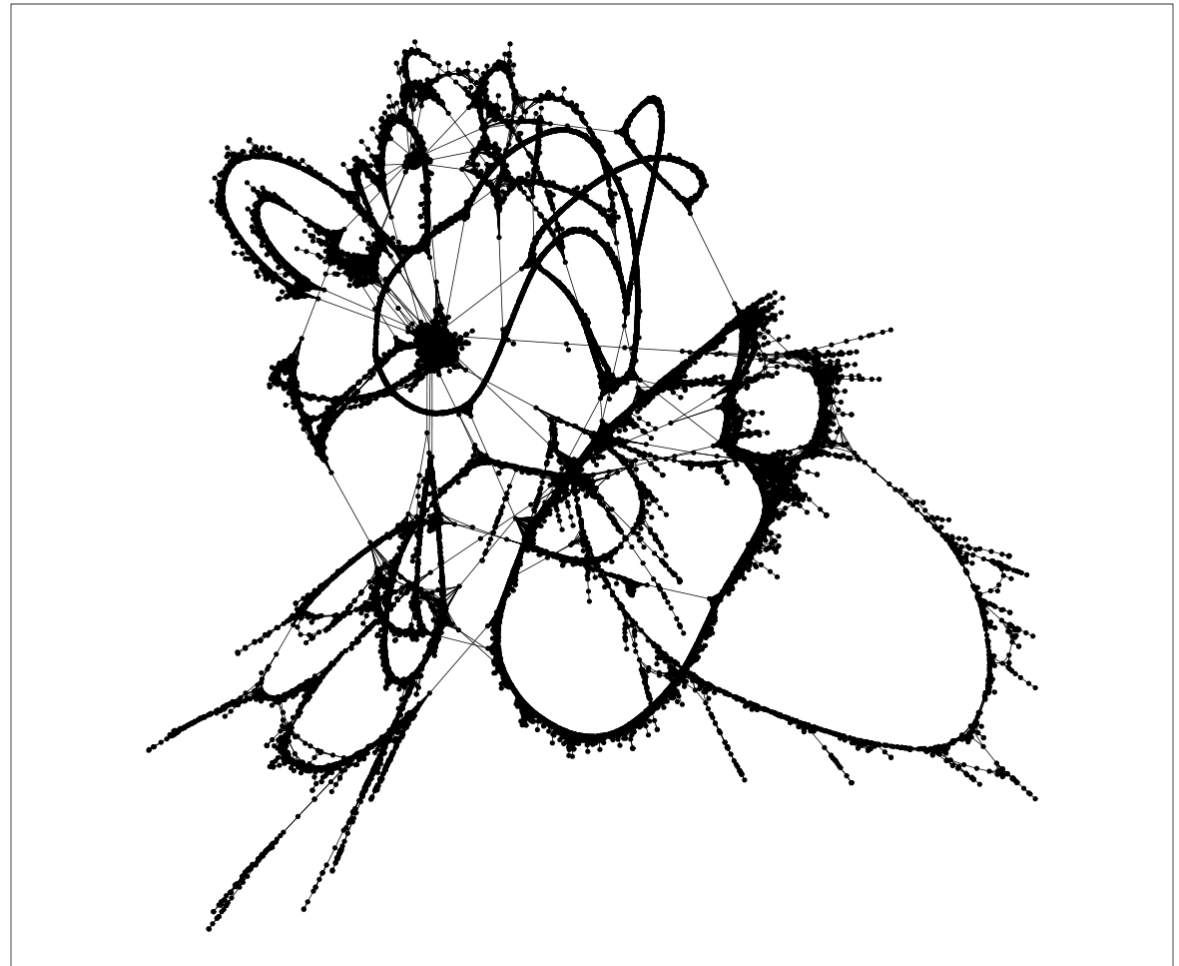
- Sequence overlaps between the reads are transformed to a graph where the **reads** are represented as **nodes** and their **similarities** as **edges** connecting the nodes
- Graph structure is examined to detect *communities of frequently connected nodes* which are split to separate clusters



# Graph-based characterization of repeat clusters



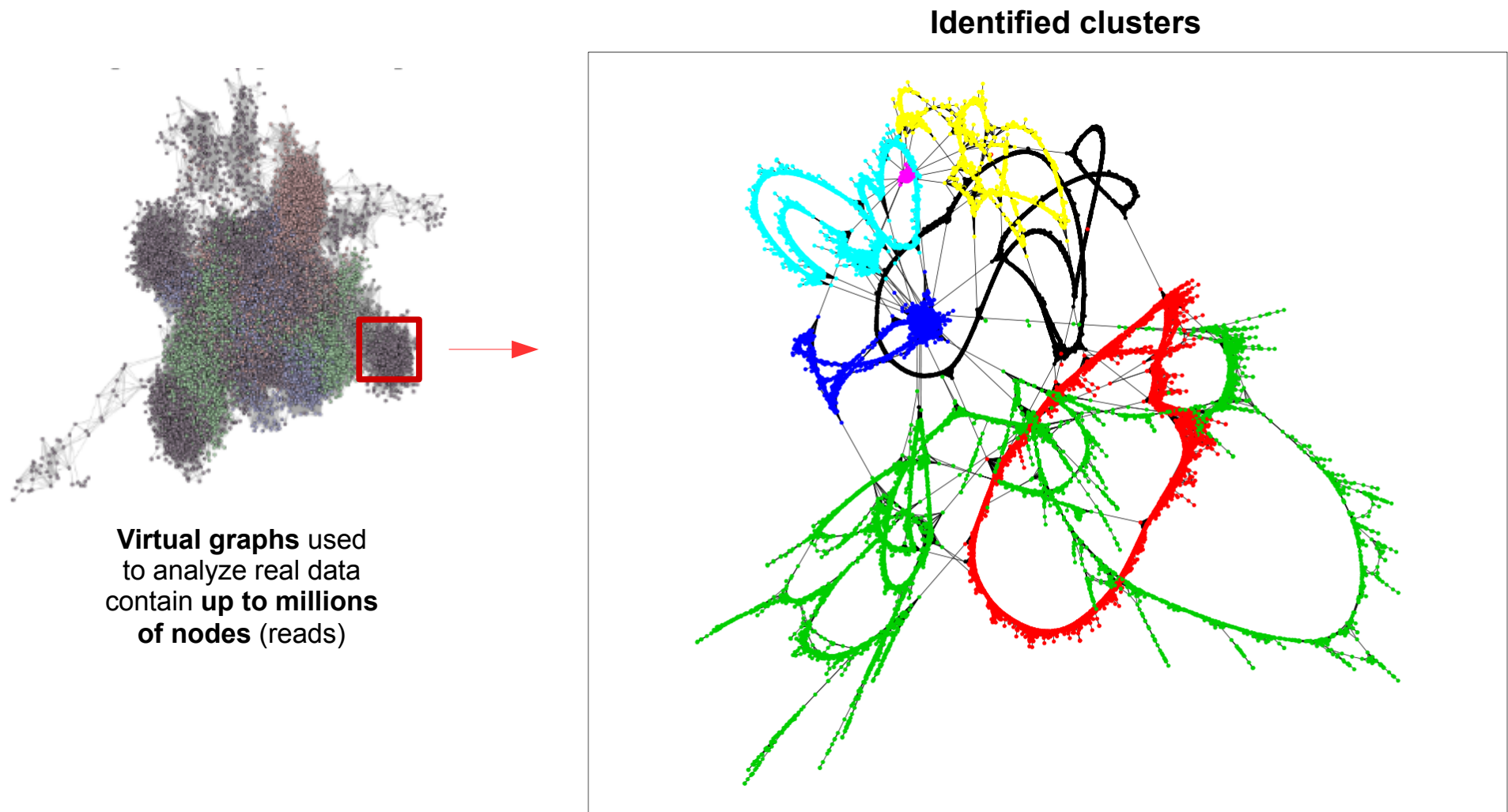
Virtual graphs used to analyze real data contain **up to millions of nodes** (reads)



## 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

# Graph-based characterization of repeat clusters



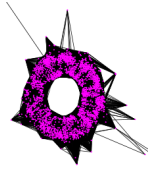
## Graph-based clustering

- 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

# Graph-based characterization of repeat clusters

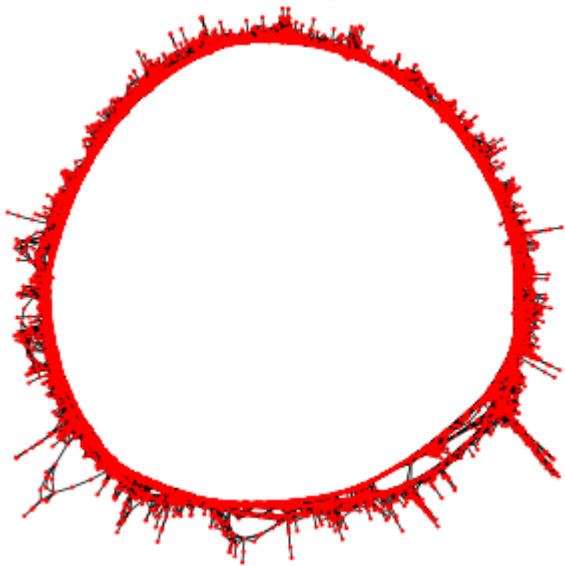
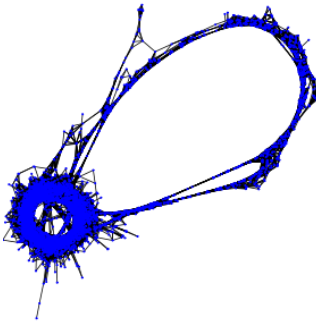
Extraction of individual clusters  
(graph shapes indicate repeat types)

Identified clusters

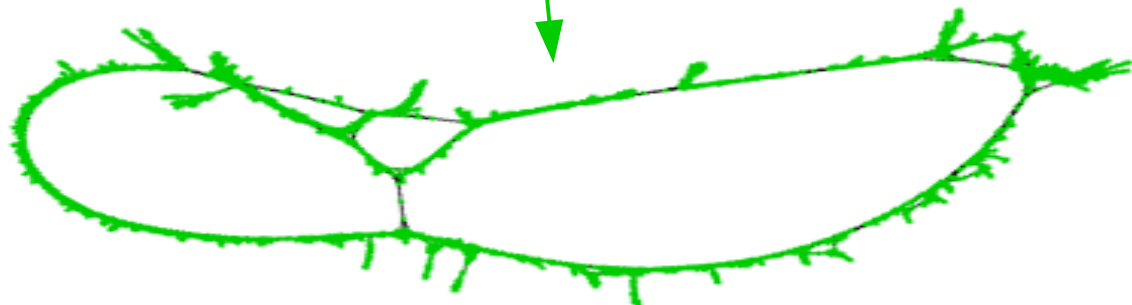
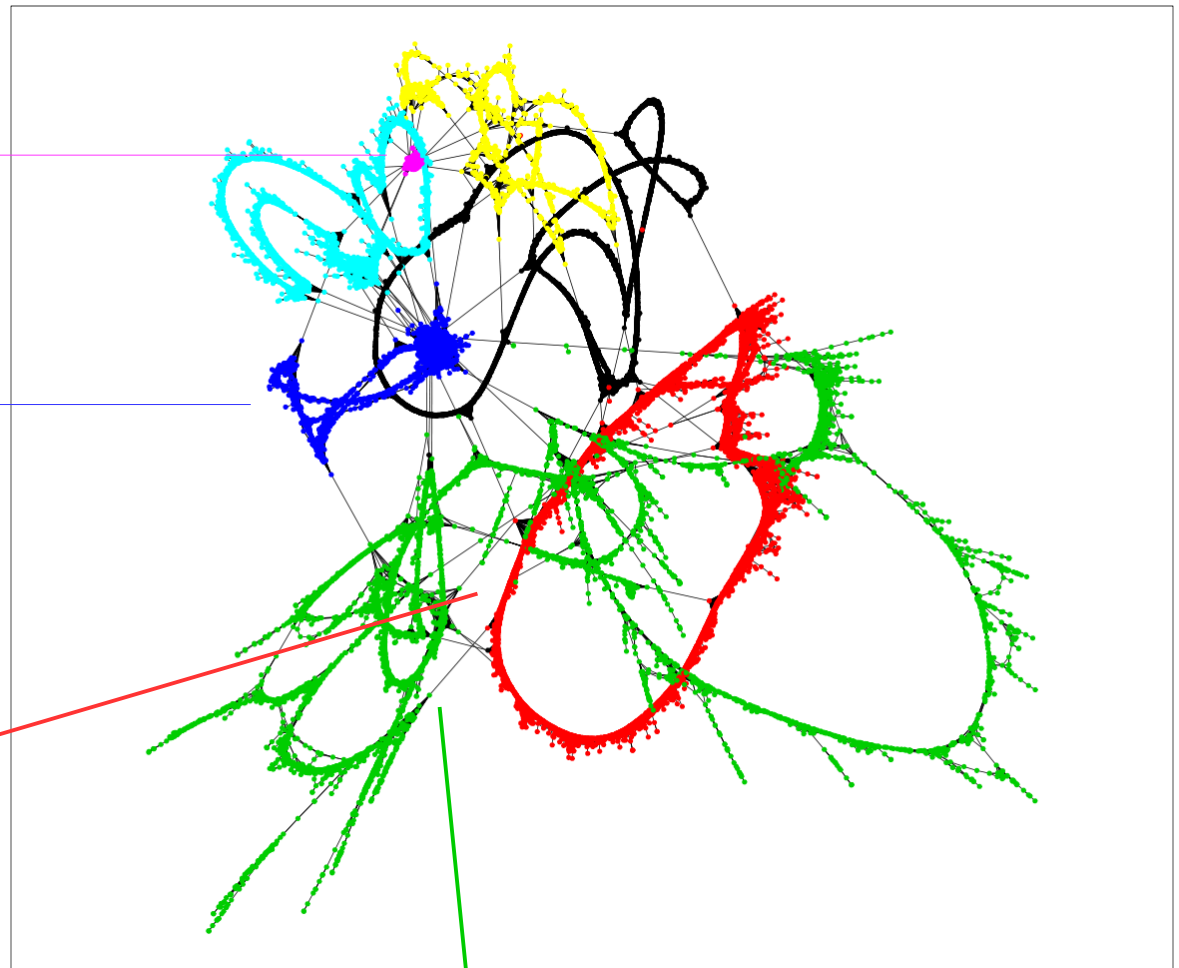


Satellite  
(simple)

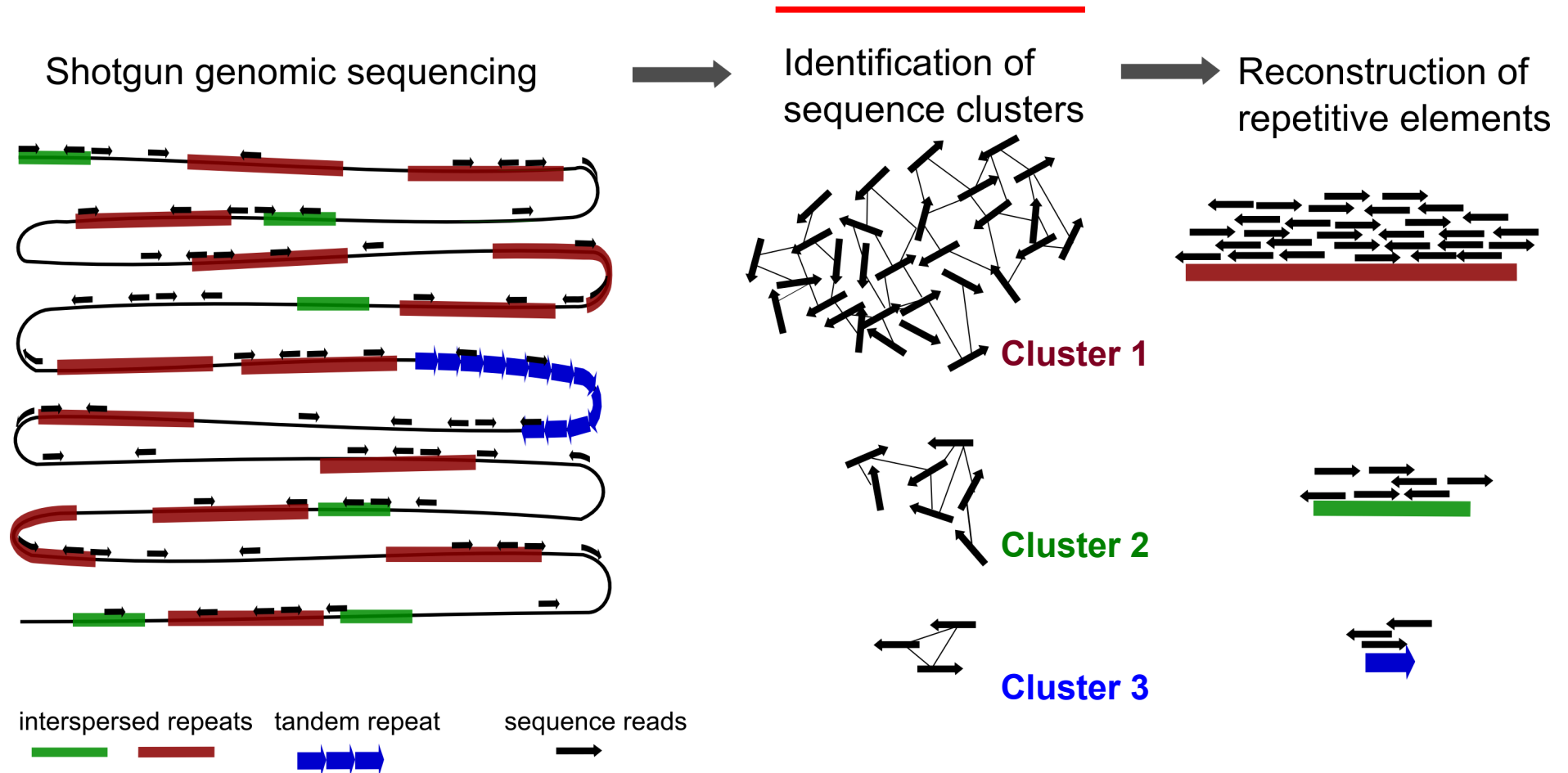
Satellite family  
(short + long monomers) >



LTR-retrotransposons



# Clustering-based repeat identification

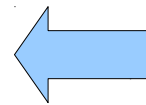


**CLUSTER = a set of frequently overlapping reads = REPEAT FAMILY**



# Repeat characterization

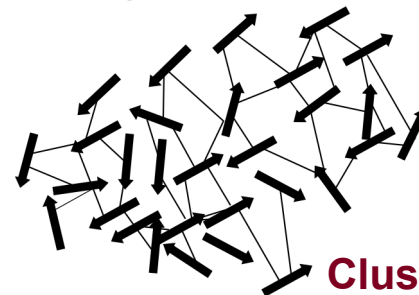
## Cluster annotation and quantification



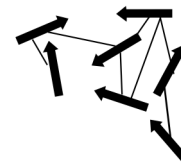
## Identification of sequence clusters

## Reconstruction of repetitive elements

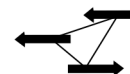
	7192836 (= 100%)	64.1			
CL	reads	genome %	class	type	note
1	304159	4.229	gypsy	Tat	PROT-RT/RH-INT
2	234749	3.264			?, PE->24,18
3	216307	3.007	gypsy	chromo	ALL domains
4	202822	2.820	copla	Maximus	ALL domains
5	149693	2.081	gypsy	Athila	ALL domains
6	145911	2.029	gypsy	Tat	ALL domains
7	143766	1.999	gypsy	chromo	ALL domains
8	142608	1.983	copla	Maximus	ALL domains
9	141836	1.972	LINE		RT
10	123886	1.722	gypsy	chromo	GAG
11	79345	1.103			?,PE->21,95
12	72781	1.012	copla	Angela	ALL domains
13	67096	0.933	gypsy	Tat	ALL domains
14	65455	0.910	gypsy	Athila	GAG, PE->1(!!!),36
15	62334	0.867	gypsy	Tat	ALL domains
16	53845	0.749	copla	Ivana/Oryco	ALL domains
17	49341	0.686			? DNA transp ?? + TR
18	45062	0.626			?,PE->2,63
19	44762	0.622			?,PE->28
20	43332	0.602	tandem		monom ?? ~400 (~1200)
21	42344	0.589	gypsy	chromo	ALL domains
22	40125	0.558	gypsy	Tat	PE->15
23	39923	0.555			?,PE->7,73
24	36353	0.505	gypsy	chromo	(GAG), PE->2,3,28
25	35977	0.500			?,PE->2,81
26	35674	0.496			?
27	34829	0.484	rDNA	5S	
28	34534	0.480	gypsy	chromo	PE->29,19,24
29	34302	0.477	gypsy	chromo	PE->28
30	33114	0.460			?
31	32930	0.458			?



Cluster 1



Cluster 2

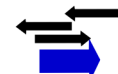


Cluster 3

...

...

(Cluster x)

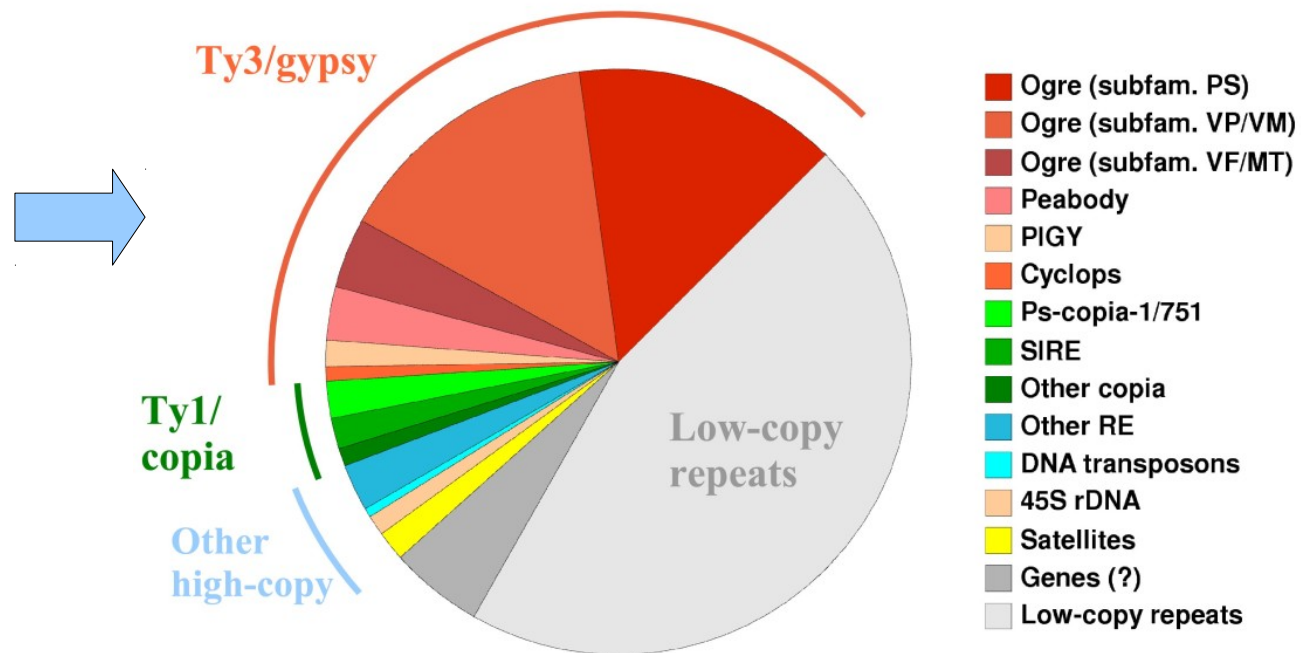


# Repeat characterization

## Cluster annotation and quantification

	7192836 (= 100%)	64.1			
CL	reads	genome %	class	type	
1	304159	4.229	gypsy	Tat	
2	234749	3.264			
3	216307	3.007	gypsy	chromo	
4	202822	2.820	copia	Maximus	
5	149693	2.081	gypsy	Athila	
6	145911	2.029	gypsy	Tat	
7	143766	1.999	gypsy	chromo	
8	142608	1.983	copia	Maximus	
9	141836	1.972	LINE		
10	123886	1.722	gypsy	chromo	
11	79345	1.103			
12	72781	1.012	copia	Angela	
13	67096	0.933	gypsy	Tat	
14	65455	0.910	gypsy	Athila	
15	62334	0.867	gypsy	Tat	
16	53845	0.749	copia	Ivana/Oryco	
17	49341	0.686			
18	45062	0.626			
19	44762	0.622			
20	43332	0.602	tandem		
21	42344	0.589	gypsy	chromo	
22	40125	0.558	gypsy	Tat	
23	39923	0.555			
24	36353	0.505	gypsy	chromo	
25	35977	0.500			
26	35674	0.496			
27	34829	0.484	rDNA	5S	
28	34534	0.480	gypsy	chromo	
29	34302	0.477	gypsy	chromo	
30	33114	0.460			
31	32930	0.458			

## Proportions of various repeat types in a genome



# RepeatExplorer

- *What RepeatExplorer can do for you:*

- Identify all sequences with certain number of repetitions
- Repeat quantification
- Provide models of repeat populations (sequence variability)
- Help in repeat classification/annotation (in plant genomes)

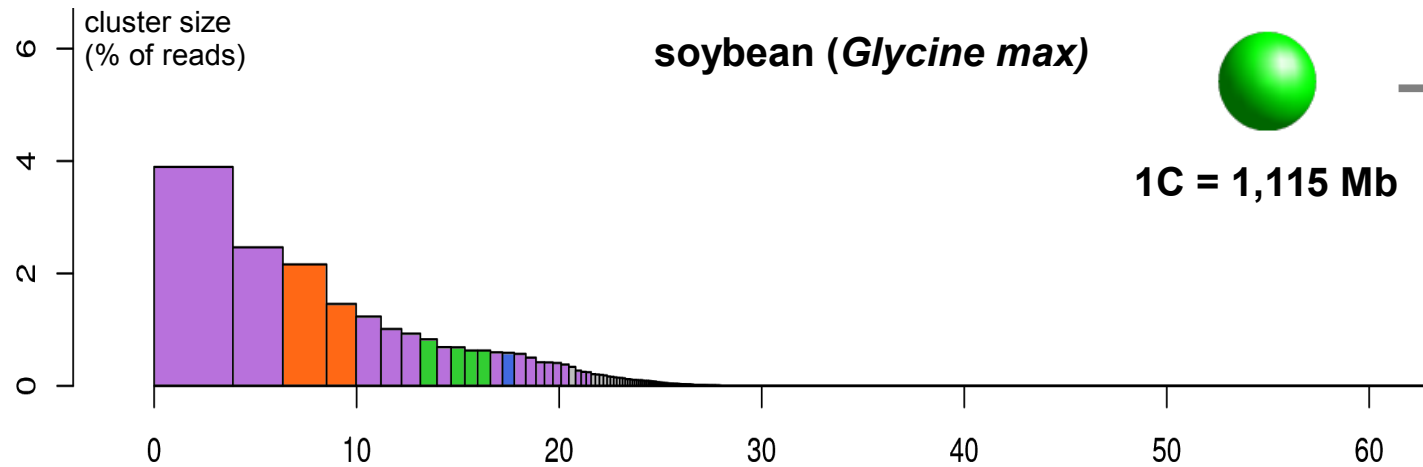
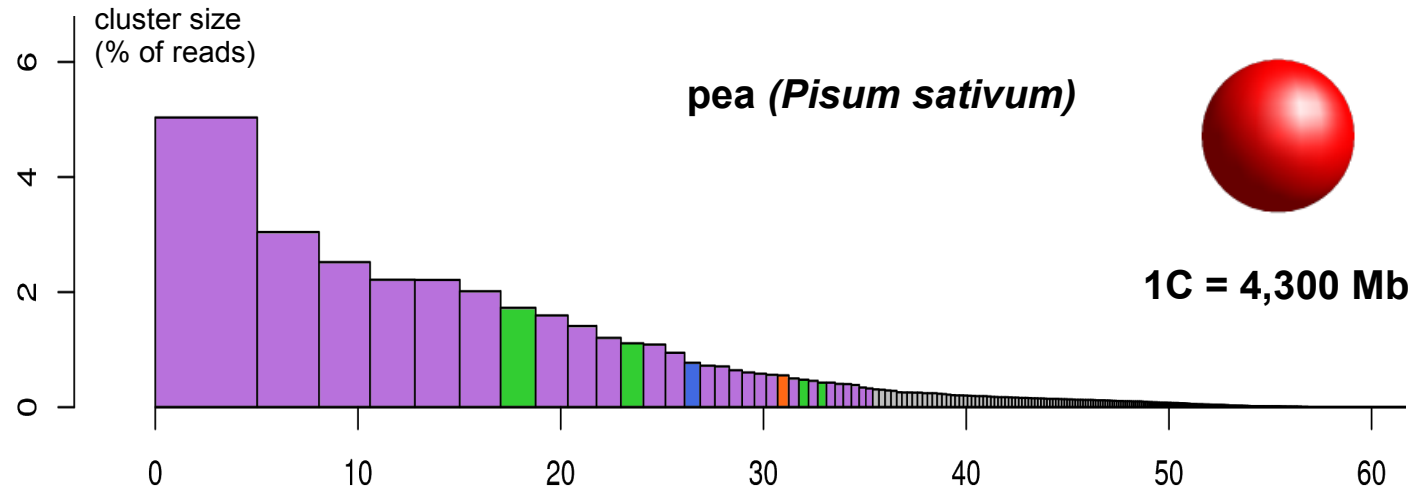
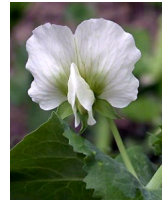
- *What it cannot do:*

- Genome assembly or reconstruction of individual repeat copies
- Analyze repeats using assembled genomes as input
- Use long NGS reads (PacBio, Oxford Nanopore) as input (*work in progress*)
- Identify some tandem repeats with very short monomers or other low-complexity repeats

Think about proper design of your analysis - RE is just a program designed for a specific type of input (e.g. it needs WGS data as input and it will not work with BAC clones)

# Repeat composition of individual species

**Repeat type:** Ty3/gypsy Ty1/copia Satellite DNA rDNA other/unknown



# Comparative analysis

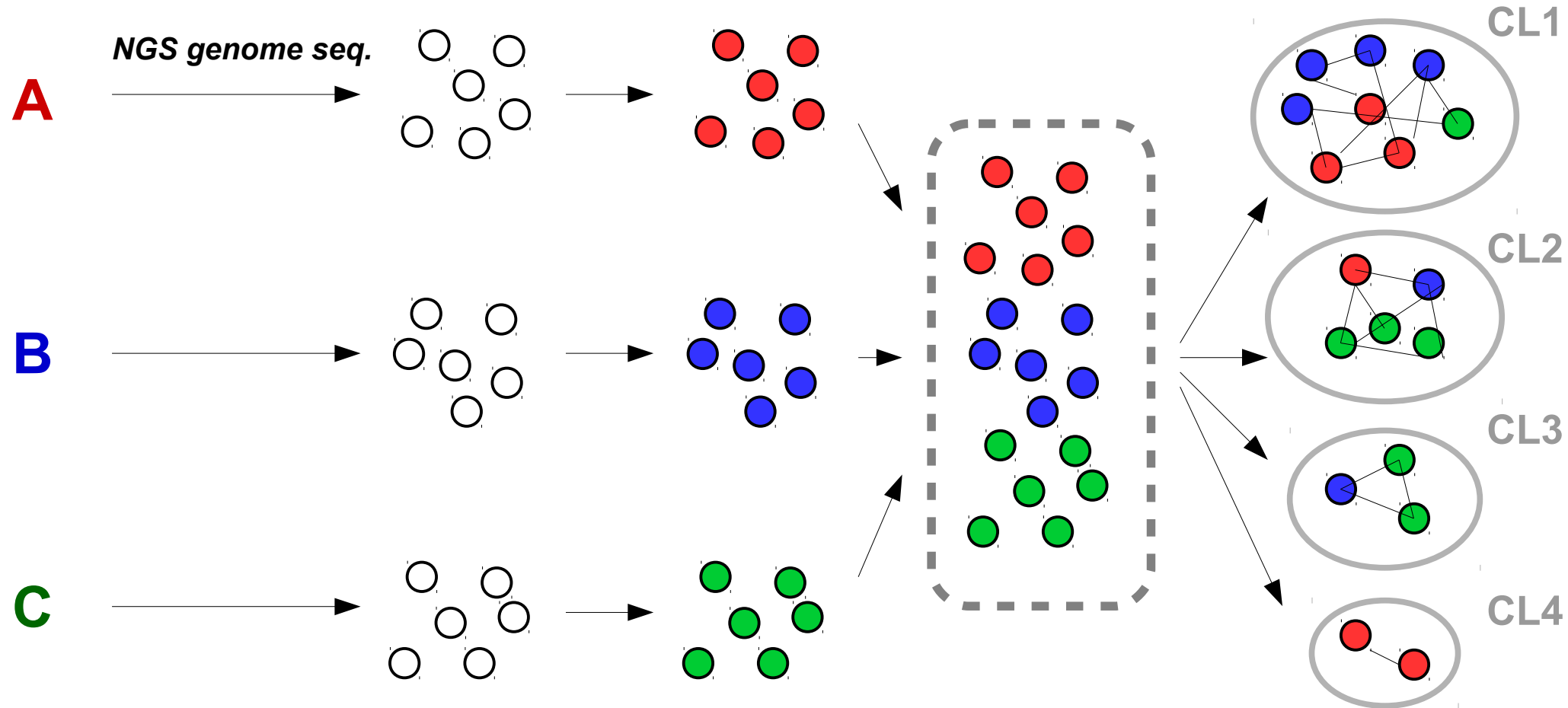
SPECIES  
(samples)

Reads

Modify names

Merge

CLUSTERS



# Comparative analysis

Sequence reads from **multiple species are mixed** and subjected to clustering analysis

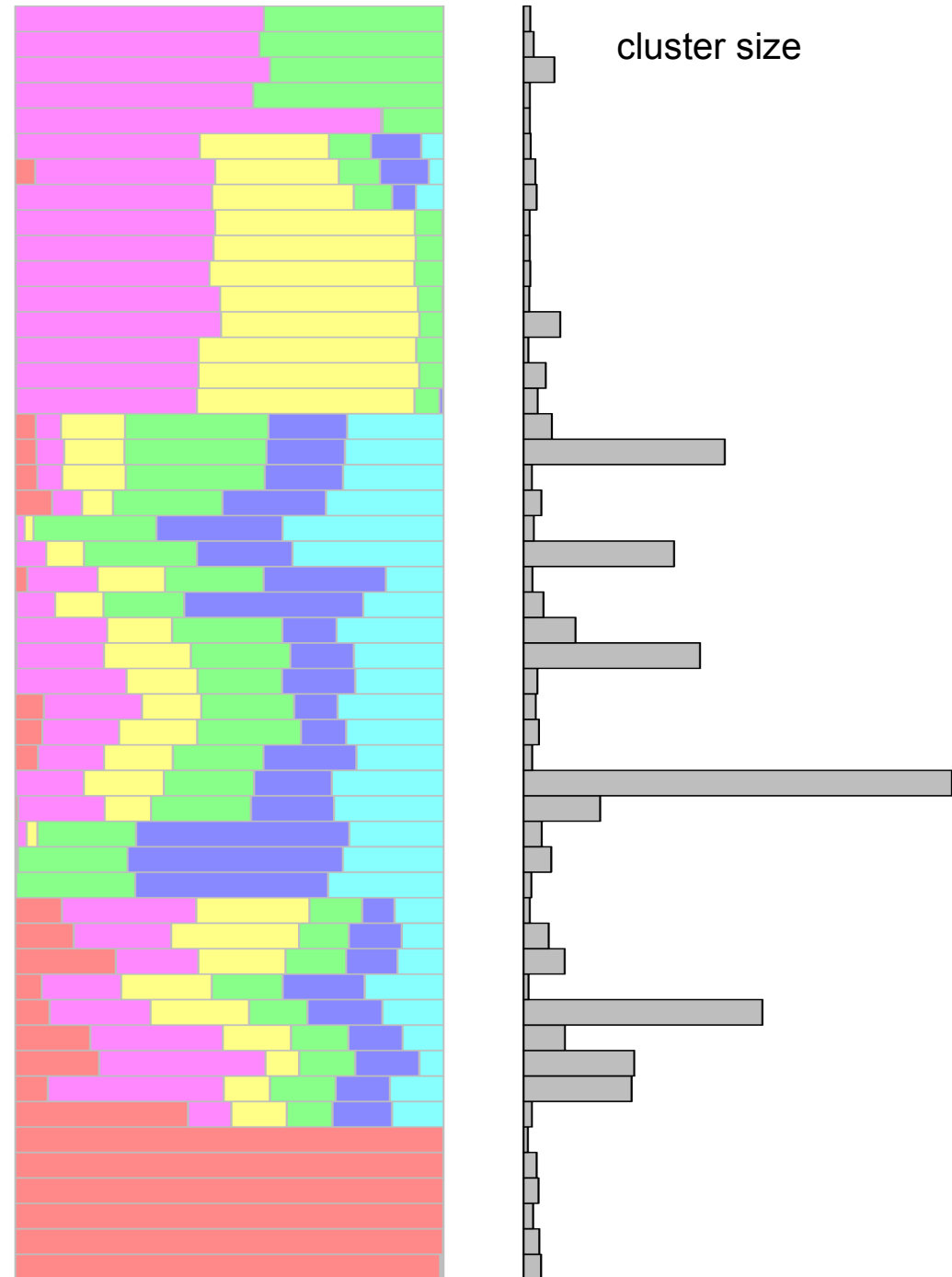
- Efficient identification of homologous repeats from different species
- Easy quantification



## SPECIES:



proportion  
of reads from  
individual  
species >>>



# Comparative analysis

**A**

**SPECIES:**

*Ensete gillettii*

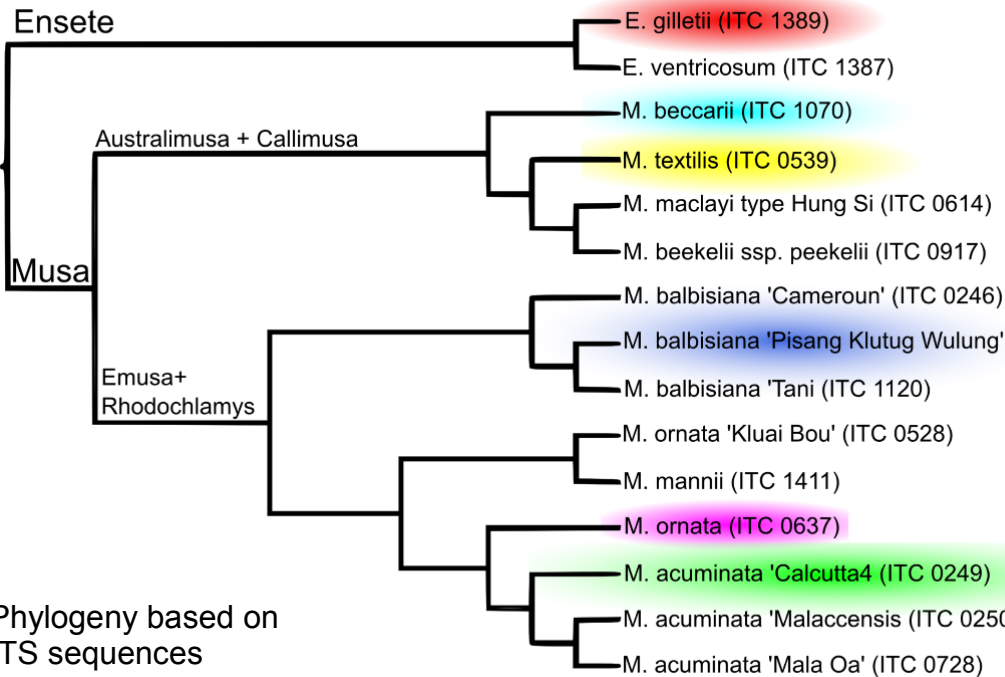
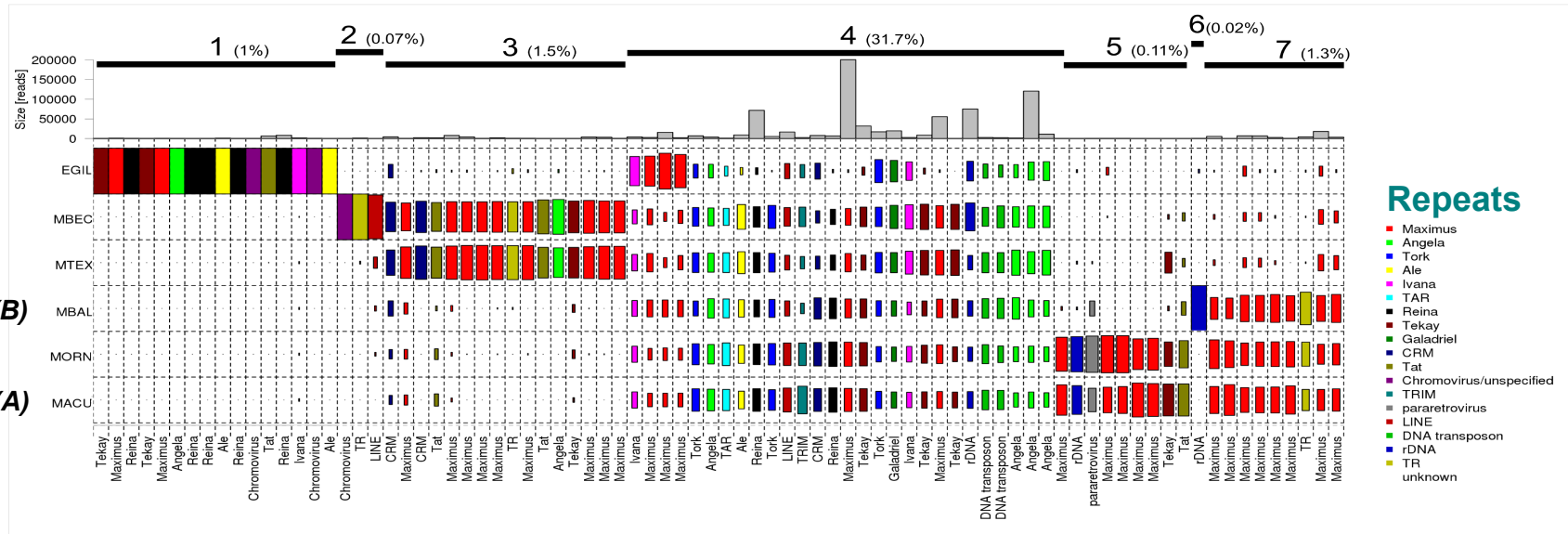
*Musa beccarii*

*M. textillis* (T)

*M. balbisiana* (B)

*M. ornata*

*M. acuminata* (A)

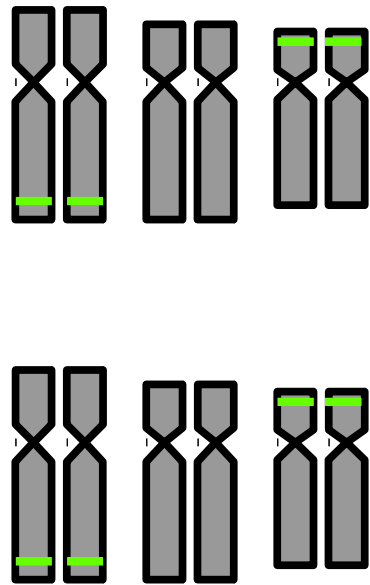


Phylogeny based on  
ITS sequences





# Comparative analysis (searching for B-specific repeats)



Next-gen  
sequencing

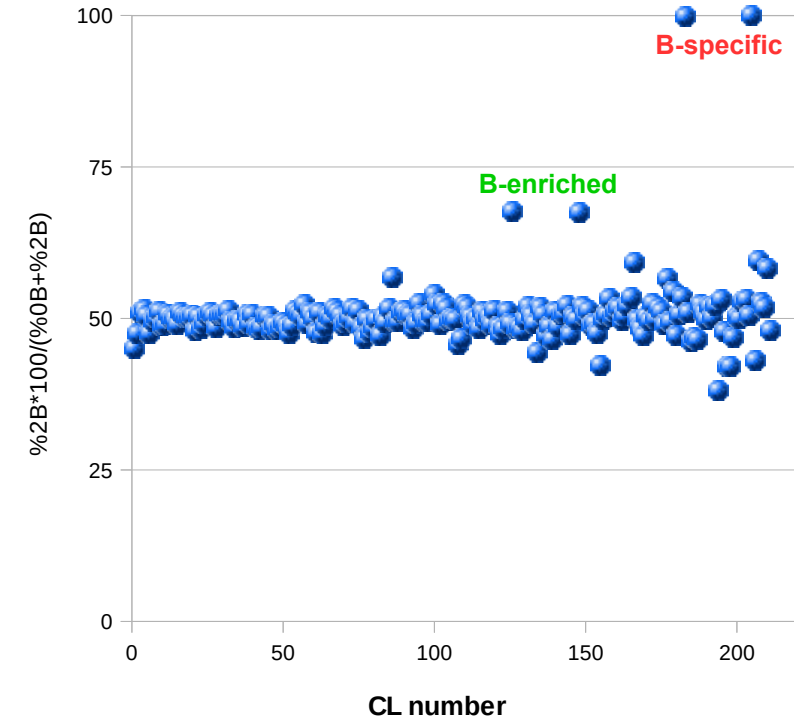
As + Bs

Next-gen  
sequencing

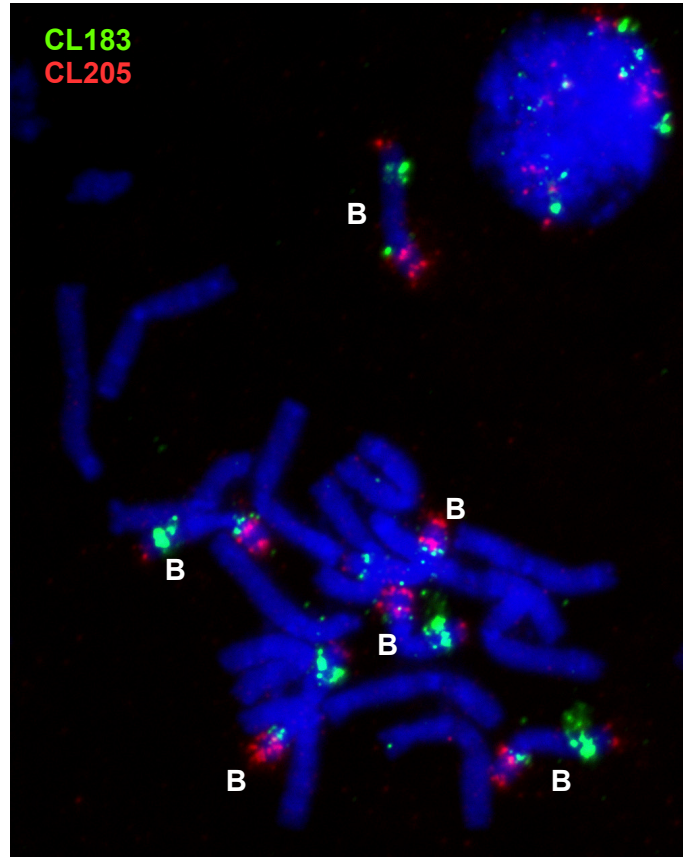
As - 0 -

Mix coded  
reads  
and  
perform  
**CLUSTERING**

## Detection of clusters enriched on Bs

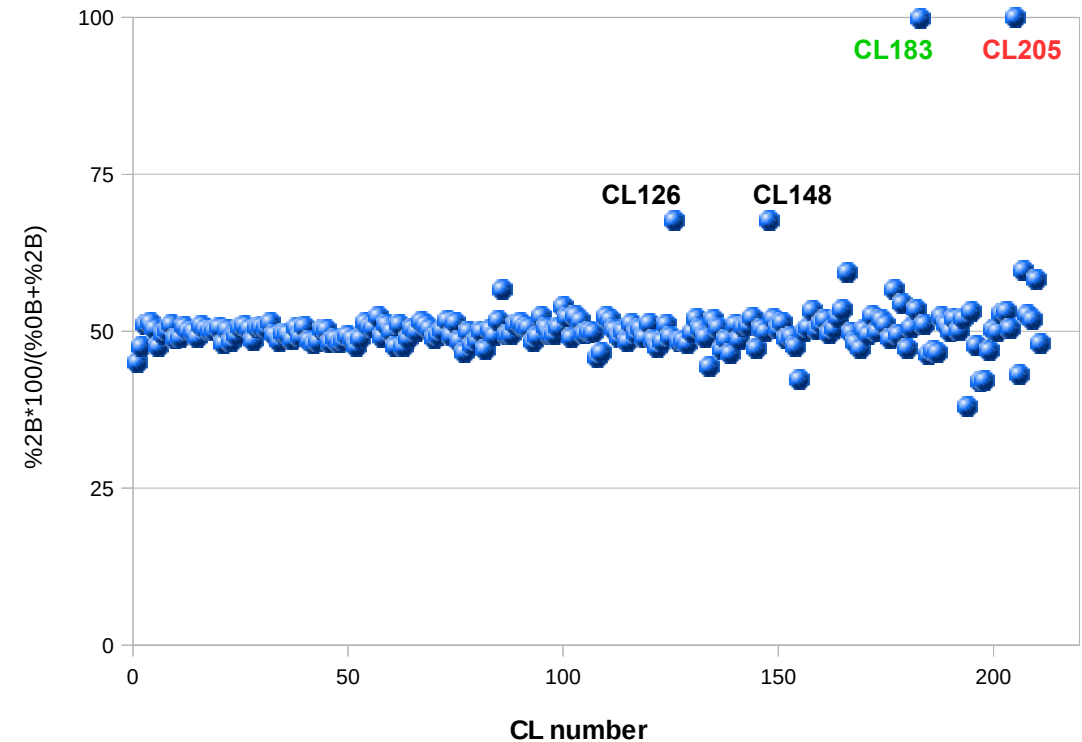


# *Aegilops speltoides*, comparative analysis of 0B / 2B plants



*FISH by Alevtina Ruban*

## Detection of clusters enriched on Bs



*Sequencing by Houben lab*

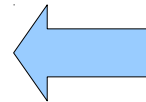
## Repeats significantly enriched on Bs

CL	reads	[%]	0B	%	2B	%	$\frac{\%2B \cdot 100}{\%2B + \%0B}$	Repeat type
126	5216	0.191	1688	0.120	3528	0.252	<b>68</b>	satellite (86 bp)
148	2877	0.105	934	0.067	1943	0.139	<b>68</b>	tandem (?)
183	731	0.027	1	0.000	730	0.052	<b>100</b>	satellite (~1.1 kb)
205	358	0.013	0	0.000	358	0.026	<b>100</b>	satellite (185 bp)

# Cluster-centered downstream applications

## Cluster annotation and quantification

CL	reads	genome %	class	type
1	304159	4.229	gypsy	Tat
2	234749	3.264		
3	216307	3.007	gypsy	chromo
4	202822	2.820	copia	Maximus
5	149693	2.081	gypsy	Athila
6	145911	2.029	gypsy	Tat
7	143766	1.999	gypsy	chromo
8	142608	1.983	copia	Maximus
9	141836	1.972	LINE	
10	123886	1.722	gypsy	chromo
11	79345	1.103		
12	72781	1.012	copia	Angela
13	67096	0.933	gypsy	Tat
14	65455	0.910	gypsy	Athila
15	62334	0.867	gypsy	Tat
16	53845	0.749	copia	Ivana/Oryco
17	49341	0.686		
18	45062	0.626		
19	44762	0.622		
20	43332	0.602	tandem	
21	42344	0.589	gypsy	chromo
22	40125	0.558	gypsy	Tat
23	39923	0.555		
24	36353	0.505	gypsy	chromo
25	35977	0.500		
26	35674	0.496		
27	34829	0.484	rDNA	5S
28	34534	0.480	gypsy	chromo
29	34302	0.477	gypsy	chromo
30	33114	0.460		
31	32930	0.458		



## CLUSTERS

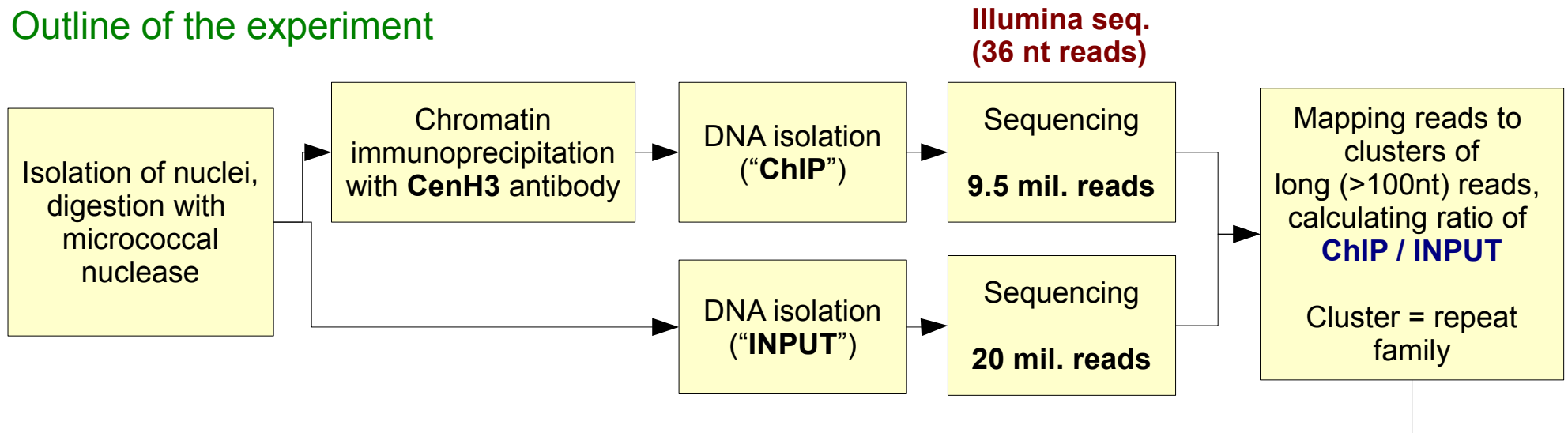
- Represent specific repeat families/variants or their parts
- They are collections of sequence reads capturing full sequence variability of repeat populations

Using clusters as reference for similarity-based classification of:

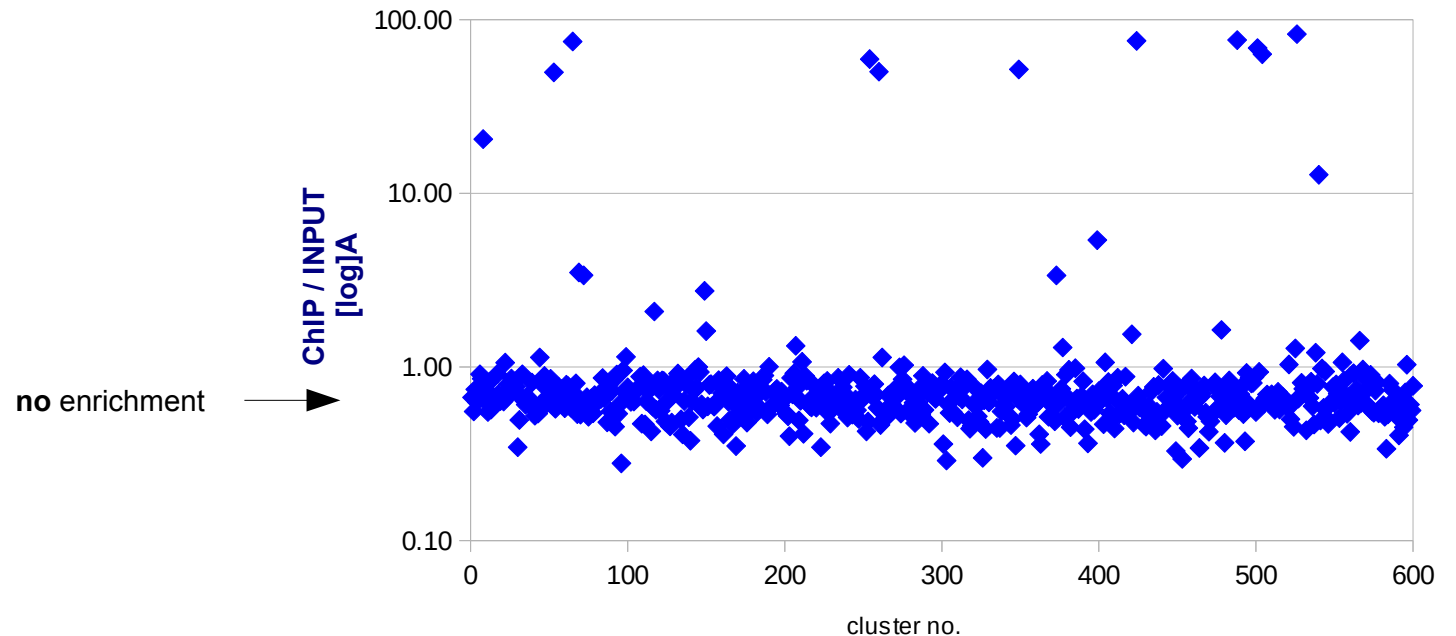
- RNA-seq reads (mRNA, smRNAs,...)
  - Detection of transcribed repeats
  - Comparative analysis (tissues,...)
- ChIP-seq reads
  - Association of repeats with specific types of chromatin

# Identification of centromeric repeats by ChIP-seq

## Outline of the experiment

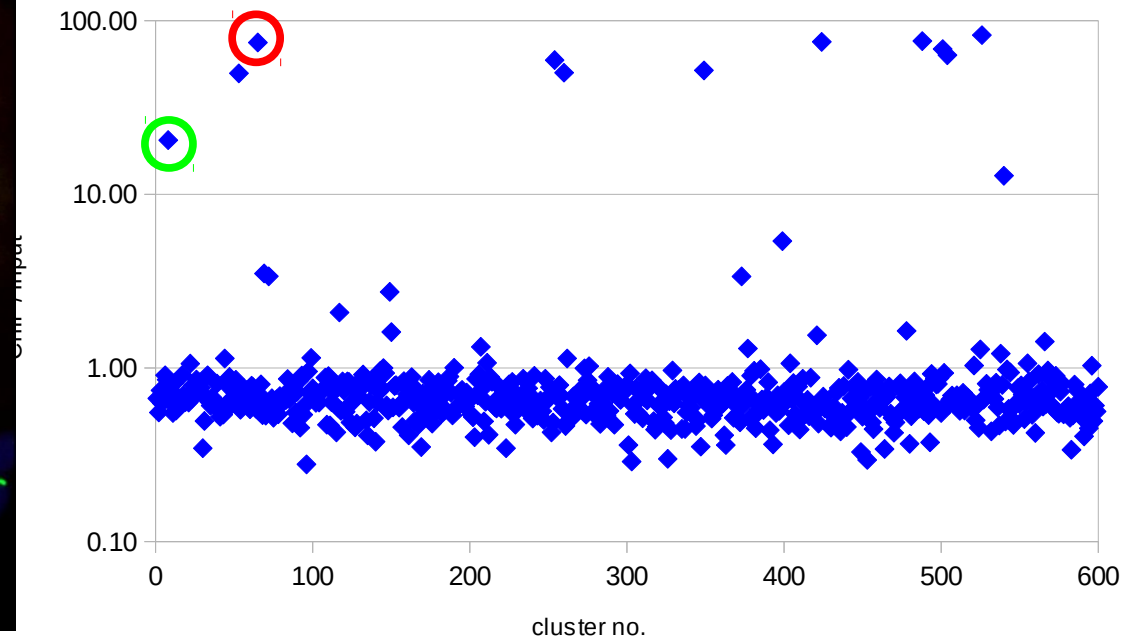
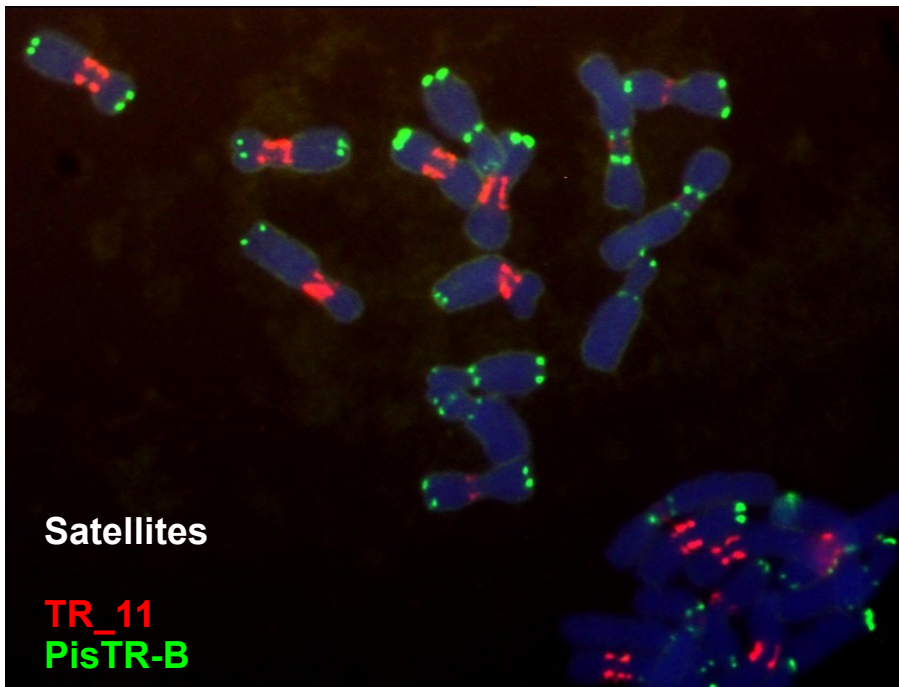
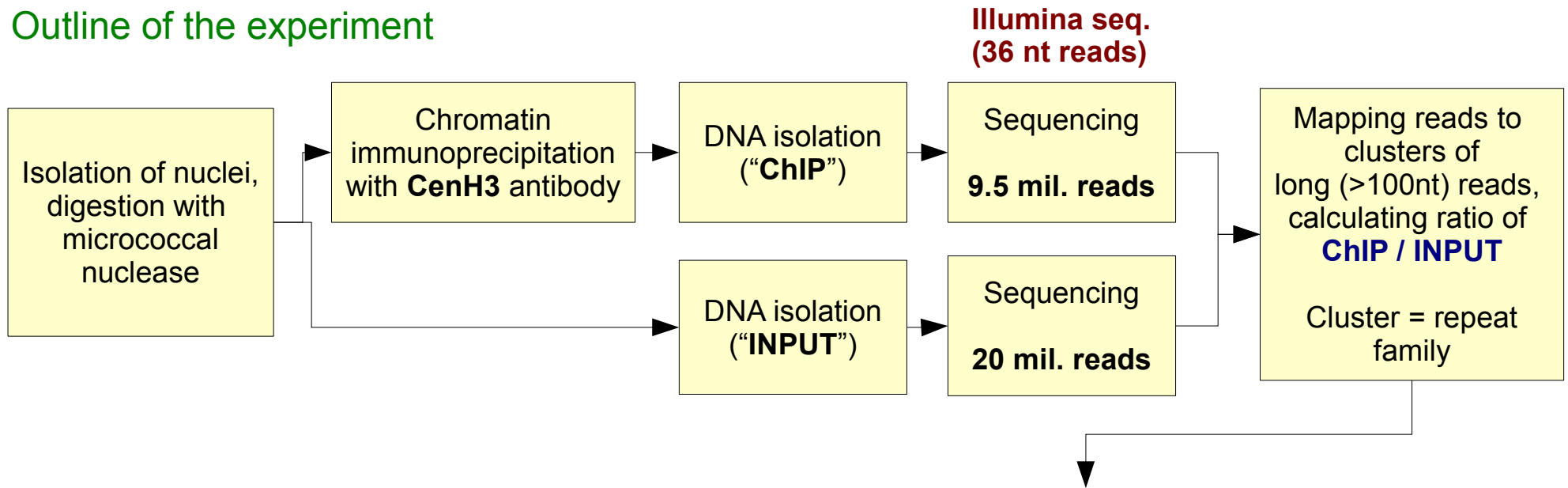


**100-fold enrichment**



# Identification of centromeric repeats by ChIP-seq

## Outline of the experiment



# History



• First paper on repeat clustering from NGS data (Macas et al. 2007)

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

*command-line version*

• **Public RE server** (Novak et al. 2013)

*RE runs under Galaxy (cluster of 12 servers)*



• Automated repeat classification

*New tools*  
• TAREAN  
• ChIP-seq  
• PROFREP  
• RE database

**2017**

**RepeatExplorer 2**

*Featuring:*

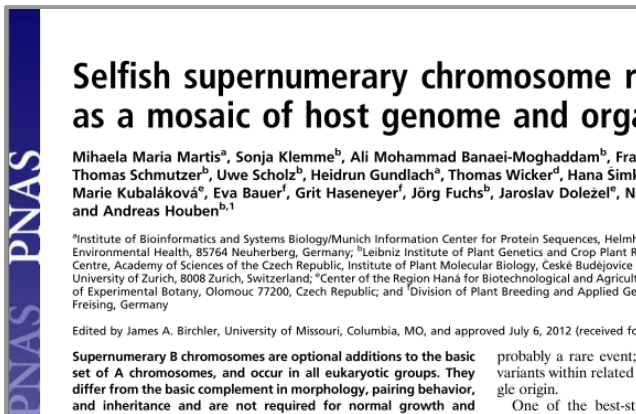
- New protein db
- TAREAN
- Super-clusters



# Repetitive DNA characterization using RepeatExplorer

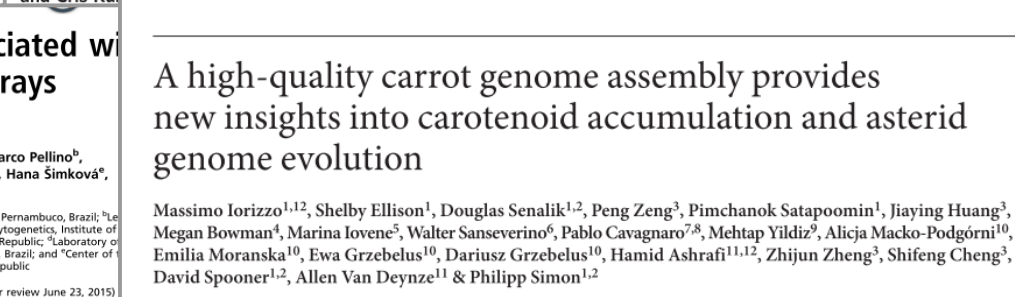
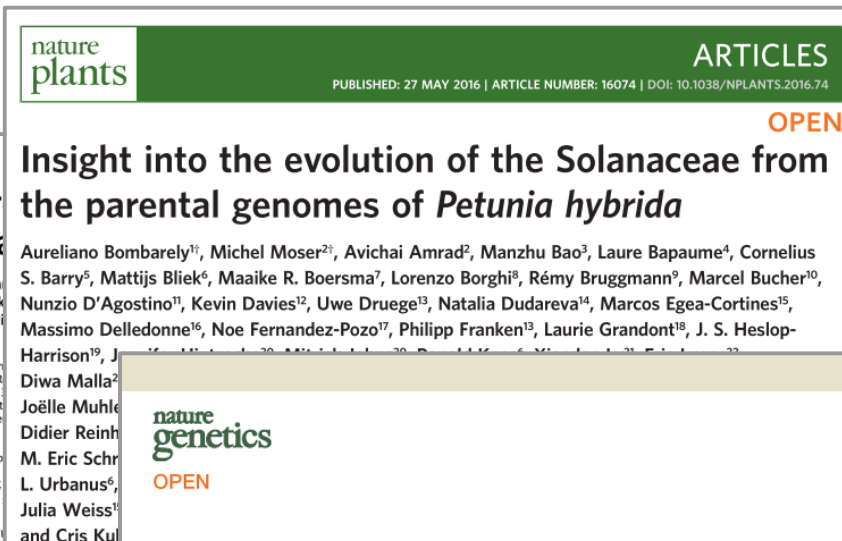
## Plants

- Over 100 species characterized so far
- Comparative studies
- Whole genome assembly projects



Holocentric chromosomes lack a primary constriction, in contrast to monocentrics. They form kinetochores distributed along almost the entire poleward surface of the chromatids, to which spindle fibers attach. No centromere-specific DNA sequence has been found for any holocentric organism studied so far. It was proposed that centromeric repeats, typical for many monocentric species, could not occur in holocentrics, most likely because of differences in the centromere organization. Here we show that the holocentric centromeres of the Cyperaceae *Rhynchospora pubera* are

cases, a longitudinal CENH3-positive centromere structure was served during mitosis. In the rush *Luzula* (Juncaceae), the long centromere forms a groove (here referred as the "centromeric groove") in each sister chromatid along almost the whole meta-chromosome except for the most terminal regions (9–11). Re-a similar centromere organization was found in the sedge *Rhynchospora pubera* (12). The absence of CENH3 and the centromeric protein C (CENP-C) in some lineages of holocentric (13) challenges the general notion of a conserved molecular



We report a high-quality chromosome-scale assembly and analysis of the carrot (*Daucus carota*) genome, the first sequenced genome to include a comparative evolutionary analysis among members of the euasterid II clade. We characterized two new polyploidization events, both occurring after the divergence of carrot from members of the Asterales order, clarifying the evolutionary scenario before and after radiation of the two main asterid clades. Large- and small-scale lineage-specific duplications have contributed to the expansion of gene families, including those with roles in flowering time, defense response, flavor, and pigment accumulation. We identified a candidate gene, DCAR\_032551, that conditions carotenoid accumulation (Y) in carrot taproot and is coexpressed with several isoprenoid biosynthetic



# Repetitive DNA characterization using RepeatExplorer

## Plants

- Over 100 species characterized so far
- Comparative studies
- Whole genome assembly projects

## Mammals

Bats, deer

## Fish

- Austrolebias charrua*, *Cynopoecilus melanotaenia*

## Insects

- Locust, grasshoppers, kissing bugs

## Worms

- Soil helminths

