Using RepeatExplorer output for repeat

annotation and quantification

RepeatExplorer pipeline



The pipeline combines all principles of repeat identification:

De novo ("repetitiveness") (graph-based clustering can handle millions of reads)

Structure-based approaches

- shape of cluster graphs
- investigation of assembled contigs

Similarity searches

- protein domains (blastx)
- CDD (RPS-BLAST)
- RepeatMasker

Repeat annotation is achieved by combining all three approaches

Summary table

Top clusters

	cluster	total length [bp]	number of reads	Genome proportion[%]	cumulative GP [%]	Repeat Masker	Domain hits	Layout	All missing mates [%]	Missing mates with no similarity hit [%]	Portion of similarity hits to other clusters[%]	Outside reads with similarity [%]
1	<u>CL1</u>	1611900	16119	18.00	18	rRNA (8531hits, 50.7%) DNA.MULE.MuDR. (368hits, 2.08%) LTR.Gypsy (7hits, 0.0161%) Satellite (1hits, 0.0023%)		R	4.746	4.746	0.00000	0.000
2	<u>CL2</u>	988500	9885	11.10	29	LTR.Copia (5818hits, 51.2%) Low_complexity (50hits, 0.202%) Simple_repeat (51hits, 0.168%) LTR.Gypsy (6hits, 0.0245%)	Ty1-RT Ty1/copia Ivana/Oryco (1346 hits 13.6%) Ty1-PROT Ty1	and the second sec	9.944	9.924	0.0000	0.000
3	<u>CL3</u>	922400	9224	10.30	39	LTR (2524hits , 19.1%) Satellite (312hits, 2.28%) LTR.Gypsy (54hits, 0.378%) DNA.CMC.EnSpm (63hits, 0.287%) DNA.hAT.Tip100 (3hits, 0.0194%)	Ty3-RT Ty3/gypsy chromovirus (2 hits 0.0217%)		25.520	25.520	0.00012	0.043
4	<u>CL4</u>	900800	9008	10.10	50	Low_complexity (340hits, 1.23%) LTR.Copia (100hits, 0.502%) Simple_repeat (1hits, 0.00422%) LTR.Gypsy (1hits, 0.00389%)	Ty3-RT Ty3/gypsy Ogre/Tat (1 hits 0.0111%)		6.083	6.083	0.00000	0.000

Summary table – new features

Automatic Classification Summary

	repeat_name	Content [%	6] number_of_cluste
1	ltr/copia/unspecified	0.515	3
2	ltr/copia/maximus	2.280	8
3	ltr/copia/tork	0.260	1
4	ltr/copia/angela	0.155	1
5	ltr/copia/ivana	0.268	1
6	ltr/gypsy/unspecified	3.200	7
7	ltr/gypsy/chromo	5.080	8
8	ltr/gypsy/ogre_tat	40.500	66
9	ltr/gypsy/athila	1.010	3

Automatic classification / super-clusters are not yet available in the public RE

Top clusters

	cluster	total length [bp]	number of reads	Genome proportion[%]	cumulativ GP [%]	e Automatic classification	Super cluster	Repeat Masker	Protein domain hits	blastn hits	Layout	Repeat Masker custom library	group - reads	numbe r of	All missing mates [%]	Missir mates with n simila hit [%
1	<u>CL1</u>	3670900	36709	1.900	1.9	nd	1	Satellite (34392hits, 86.3%) Low_complexity (45hits, 0.0331%) LINE.Penelope (17hits, 0.0283%) DNA (5hits, 0.00627%) Simple_repeat (3hits, 0.00204%) .TR.Gypsy (1hits, 0.00185%) CC.Helitron (1hits,		organelle/plastid (36 hits 0.0981%)		Satellite.VicTR.A (31921hits, 82.3%) Telomeric.Arabidopsis (19hits, 0.031%)	group VSA_ VSP2 VPN_ VML_	content 0 (0%) 0 (0%) 15853 (2.77%) 20856 (2.59%)	3.10	0.62
2	<u>CL2</u>	3232900	32329	1.670	3.6	ltr/gypsy /chromo	2	TR.Gypsy (13535hits, 38.9%) TR.Gypsy.peabody (9877hits, 28.3%) Low_complexity (1453hits, 1.7%) TR (88hits, 0.224%) Simple_repeat (144hits, 0.184%) TR.gypsy (29hits, 0.0663%) Organelle	Ty3-INT Ty3/gypsy chromovirus (7349 hits 22.7%) Ty3-gypsy chromovirus (4605 hits 14.2%) Ty3-RH Ty3-RH Ty3-RH Ty3-RH Ty3-RH Ty3-Qypsy chromovirus (2807 hits 8.68%) Ty3-CHDII Ty3-CHDII Ty3-Qypsy chromovirus (725 hits 2.24%) Ty3-PROT Ty3/gyps	organelle/mitochondria (13 hits 0.0402%) organelle/plastid (1 hits 0.00309%)	and the second sec	LTR.Gypsy.peabody (15124hits, 43.8%) Organelle.Mitochondrion (194hits, 0.484%) LTR.Gypsy.centromeric (65hits, 0.152%) LTR.Gypsy.Ogre (2hits, 0.00424%) LTR.Gypsy (1hits, 0.00294%) Organelle.Plastid (group VSA_ VSP2 VPN_ VML_	content 2832 (1.6%) 8141 (2.17%) 7966 (1.39%) 13390 (1.66%)	15.20	5.40
3	<u>CL3</u>	3192600	31926	1.650	5.2	ltr/gypsy /ogre_tat	3	TR.Gypsy (1551hits, 3.57%) TR.Gypsy.Ogre (1614hits, 3.37%) Low_complexity (695hits, 0.803%) Simple_repeat (67hits, 113%)	DTA-CD1 NA NA (1 hits 0.00313%) Ty3-INT Ty3/gypsy chromovirus (1 hits 0.00313%)			LTR.Gypsy.Ogre (2281hits, 4.76%) LTR.Gypsy.Ogre_PA (47hits, 0.0929%) LTR.Gypsy.peabody (3hits, 0.00417%) LTR.Gypsy (2hits, 0.00244%) Satellite VieTR 4 (2hits,	group VSA_ VSP2 VPN_	content 1809 (1.02%) 1822 (0.487%) 3838 (0.67%)	17.87	6.06

Detailed information for the largest clusters

RepeatExplorer output files





Protein domains

Domain	ID	Туре	Lineage	Hits	MeanScore
Ty3-INT	MedT2_ID89	Ty3/gypsy	chromovirus	96	49.7
Ty3-INT	PopT1_ID21	Ty3/gypsy	chromovirus	74	46.7
Ty3-RT	PetI1_ID72	Ty3/gypsy	chromovirus	71	60.4
Ty3-RT	CRR3_ID17	Ty3/gypsy	chromovirus	56	48.4
Ty3-GAG	PetI1_ID72	Ty3/gypsy	chromovirus	49	35.5
Ty3-INT	MedT1_ID93	Ty3/gypsy	chromovirus	42	58.8
Ty3-INT	PopT2_ID20	Ty3/gypsy	chromovirus	41	45.3
Ty3-PROT	lpoB1_ID87	Ty3/gypsy	chromovirus	37	45.8
Ty3-PROT	Petl1_ID72	Ty3/gypsy	chromovirus	36	45.9
Ty3-RH	BraR5_ID117	Ty3/gypsy	chromovirus	31	57.7
Ty3-INT	CRA2a_ID5	Ty3/gypsy	chromovirus	31	53.4
Ty3-GAG	Petl2_ID76	Ty3/gypsy	chromovirus	30	35.0
1+1 - uu+	D 114 1070	- a.	· ·		10.0

	SW	perc	perc	perc	query	positi	ion in	query	matching	repeat	posit:	ion in	repeat	
	score	div.	del.	ins.	sequence	begin	end	(left)	repeat	class/family	begin	end	(left)	ID
	602	17.0	0.0	0.0	SJM 1000031f	1	100	(0)	C Gypsy-33 ST-I-int	LTR/Gypsy	(4879)	1441	1342	1
	256	14.3	0.0	0.0	SJM_1000031r	1	42	(58)	+ Gypsy-33_ST-I-int	LTR/Gypsy	1224	1265	(5055)	2 *
<u><u> </u></u>	581	6.8	0.0	0.0	SJM_1000031r	28	100	(0)	+ Gypsy-33_ST-I-int	LTR/Gypsy	1197	1269	(5051)	2
Similarity	885	3.0	0.0	0.0	SJM_1000270f	1	100	(0)	+ Gypsy-33_ST-I-int	LTR/Gypsy	3002	3101	(3219)	3
	802	5.0	5.0	0.0	SJM_1000270r	1	100	(0)	C Gypsy-33_ST-I-int	LTR/Gypsy	(2983)	3337	3233	4
searches	680	11.0	0.0	0.0	SJM_1002242f	1	100	(0)	+ Gypsy-33_ST-I-int	LTR/Gypsy	2328	2427	(3893)	5

Repeat annotation

For repeat annotation in individual clusters, consider:

- Shape of cluster graph
- Significant similarity to conserved types of repeats (45S and 5S rDNA, *cpDNA*)
- Similarity to conserved protein domains
- Structural features of assembled contigs (dot-plot)
- Variability and masked regions of reads in contig assembly
- (blastn and blastx to GenBank)
- Cluster connections via paired-end reads
- LTR + primer binding site (PBS) detection
- Automated classification

Be prepared to leave some clusters without annotation !

Graph shapes





Linear graphs



DNA transposons, LINEs, ...

Linear graphs



assembled contig

insertion sites



LTR-retrotransposons





Cluster fragmentation leads to linear graphs !



LTR-retrotransposon Ogre in *Vicia pannonica* (~ 40% of the genome)

Cluster fragmentation leads to linear graphs !



LTR-retrotransposon Ogre in *Vicia pannonica* (~ 40% of the genome)

Cluster fragmentation leads to linear graphs !

rDNA (2045hits, 80.8%) rDNA.185 (363hits, 14.4%) X25SrDNA (84hits, 3.31%) rDNA.26S (5hits, 0.198%) CDNA.26S (5hits, 0.198%)

45S rDNA

plastid DNA (contamination)

Both can be easily recognized by similarity searches

It's getting complicated...



It's getting complicated...



...but it has some meaning

(in this case, presence of element variants differing in LTR sequence and in deletion within gag-pol region)



It's getting complicated...



...but it has some meaning

Using different algorithms for calculating graph topology



TANDEM REPEATS

Read length << monomer



Read length >= monomer









MITE (foldback)







ЯIT







contig

A new tool for detection of LTR / PBS sites



A new tool for detection of LTR / PBS sites



A new tool for detection of LTR / PBS sites



The tool searches assembled contigs (ACE files) for differences in proportions of masked reads accompanied by occurrence of TG/CA



A new tool for detection of LTR / PBS sites



Program output:

	CL	contig	pos.	site	site_depth	out_mas	maske∳	maske∳	region_in	region_out	blast to tRNA	%	lengti	1	5	ite		tRNA		E-val
															f	rom	to	from	to	
	19	400	364	TGCGA CA	106.6	30.4	0.0362	0.2975	GCGAGGAA	GATGGCGA	At-chr2.tRNA28-A	rg⊳ 100) 18	0	0	3	20	23	6	7E-007
				(window size	7)	1	1	I	ł	tRNA-A	rg	I	1 1	1				•		I
г	AAT	* TTTT	TCCG	CGACCATCO	GGA * GGAA	TCGTAT	rttt*c(GAGATG	CGACAGATG	GCGACTCTGC	TGGGGAC * * TA * G	CTCCA	AGCAA	AAGAG	GAGI	GAA	GCCI	TAATT	TAG	
	• • • •		1000																	
1	AAT	* 111111	TCCG	CGTCCATCO	GCGA*GG <mark>G</mark> A	TCGTAT	rttt*co	GAGATG	CGACAA TAA	CCATTAGATG	TGAAGACACAAC'I	TGATT	AGAGG	GACT	_					
.1	AAA	* TTTT	TCCG	CAACCATC	TGA * GGAA	TCATAT:	rttt*co	GAGATG	CGACATATG	G <mark>TGACTCT</mark> AC	TGGGGAC * * TA * G	CTC <mark>T</mark> A	AGCAA	AAGGO	GAG					
C	CAT	* TTTT	TCCG	AGACCATCO	GCGA * GG <mark>G</mark> A	TCGTAT	rttt*c(GAGATG	CGACAGATG	ACGACTCTGC	TGGGGAC * * TATC	TTGTC	CAAGC	AAGAG	5A					
Т	ААТ	* TTTT	TCCG	CGACCATC	GCGA * GGAA	TCCTAT	rttt <mark>*</mark> C(GAGATG	CGACAGATG	GCGACTC <mark>C</mark> GC	T <mark>A</mark> GGGA <mark>T</mark> **TA*0	CTCCA	ag <mark>a</mark> aa	AAGAG	GAG					
Г	AAT	* TTTT	TCCG	CGACCATCO	GCGA*GGAA	TCGTAT	rttt <mark>*</mark> Co	GAGATG	CGACAG <mark>T</mark> TC	TGAGTTTGCT	TGGCAAGGCTTGG	GCAAG	GTGTG	TGTC2	ΑT					
A	AAT	ΓΑ <mark>ΤΤΤ</mark>	TCCG	CGACCATCO	GC <mark>A</mark> A*GG <mark>G</mark> A	TCGTAT	rttt*c(GAGATG	CGACAGTAT	CGCGTAACAA	TCGTATCAAATCI	TACAC	TTGAA	CAAAG	ст					
Т	AAT	* TTTT	TTCG	CGACCATCO	GGA*GGAA	TCGTAT	FTTT * T	GAGATG	<mark>cga</mark> tactag	GCCCATGTAA	TCACTTCTTCTG	CCCAT	AATTG	AGTAC	CAA					
Т	ААТ	* TTTT	TCCG	CGACCAACG	GCGA*GGAA	TCGTAT	rttt*c	GAGATG	CGACAGTTA	CTTCTGAACA	TTAATCAATATGA	TTTCC	CAGAA	AGAAG	CTTG	;				
Т	ААТ	* TTTT	TCCG	CGACCATT	TGA * GGAA	TCGTAT	TTTT*C	AGATG	CGACAGATG	GCGACTCTGC	TGGGGATACTTAG	CTCCA	AGCAA	AAGAG	GAAT					
A	CAT	* TTTT	TCCG	CGACCATC	GGA*GGAA	TCGTAT	rttt*co	GAGATG	CGACAGATG	GCGACTCTGC	TGAGGAC * * TA * G	CTCCA	AGCAA	AAGAG	GAGT	GAA	G			
Δ	TTT	* TTTT	TCCA	CGACCATCO	CGA*GGAA	TCATAT	rttt*co	AGATG	CGACAGATG	GTGACTCTAC	TGGGGAC * * TA * G	стесс	AGCAA	AAGAO	AGT	GAA	e i			
2	CAT	* 77777	TCC	CGACCATCO		TCGTAT	rrrr*co	ZAGATO	GACAGATO	GCGACTCTG	TGGGGAC**TA*0	CTCTA	AGAAA	AAGAO	23.67	GAN				
-	3 3 1		Teee	CGACCATCO	CGA *CCAA	TCCTA		TAGATO		ACAATTCACC	ACAACTATCAATZ	a tra tre	2 222 <i>0</i>	TACT	- AO 1					
1	AAI		maga	CGACCATCO	GGA GGAA			JAGAIG	GACATIAG	AGAATICACC		anarc		TACT						
1	AAA		Teeg	CGACCATCA	CGA * GGAA	TCGTAT	rrrrr*co	SAGATG	CGACAGATG	GCGACTCTGC	TGGGGAT · TA · C	CTCCA	AGCAA	AAGA	FAG'I	GAA	GCC			
Т	AAA	*TATT	TTCG	CGGCCATCO	CGA*GGAA	TCGTAT:	rttra co	GAGATG	CGACAGATG	GCGACTCTGC	TGGGGAC * *TA *G	CTCCA	AGCAA	AAGAG	5AG <mark>0</mark>	GAA	GCC1			
Т	AAA	TTTAT	TICG	CGACCTTC	GCGA*GGAA	TCGTAT:	TTTT*C	GAGATG	CGACAGATG	GCGACTCGTC	TGGGGAC * * TA * G	CTCCA	AGCAA	AAGAG	FAGT	GAA	GCC			

"Chimeric" clusters ?

Reconstruction of rDNA gene cluster in Silene latifolia



Repeat quantification

-		7192836	64.1			
		(= 100%)				
	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	copia	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	copia	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	copia	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	copia	lvana/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			?
	~~	00100				DT

Based on numbers of reads in individual clusters

(more info during practical training)

Using paired-end read information for improving cluster annotation

Repeat quantification

Cluster annotation and quantification

1					
		7192836	64.1		
		(= 100%)			
	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	lvana/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		
ł	~~~	00,100			1

Proportions of various repeat types in a genome

