00000000 10000000 010 01000001 10000001 0 00010000 00000000 10022133 010000 20000000 0001 00000000 10 00000000 1 0001000 1000 0 0

# Short reads only denovo assembly of grape genomes

Michele Vidotto

# IGA CLUB

July 6, 2016







### The genome assembly problem

# It's a bit like trying to do the hardest puzzle you can imagine!



michele.vidotto@uniud.it

# Genome assembly pipeline - preprocessing



michele.vidotto@uniud.it

### Genome assembly pipeline - graph construction



### Genome assembly pipeline - postprocessing



Michele Vidotto

# The ALLPATHS-LG assembler recipe

• Best interaction between experimental design and algorithmic approach



- Overlapping paired reads merged into single long fragment
- Longer K-mer can be used for graph construction (k=96)





## The ALLPATHS-LG peculiarity

- Multiple K-spectrum based error correction of fragment libraries (k=24,25)
- Alignment and error correction of the jumping reads on the assembly graph
- Scaffolding gap filling and modules included

Michele Vidotto

• HAPLOIDIFY function identifies bubbles in graph with haploid:diploid copy number and resolves them by choosing one path



# The ALLPATHS-LG pipeline

- Modular system. Each module performing an assembly step
- 52 modules needed for the grape assemblies



#### Michele Vidotto

# Our ALLPATHS-LG recipe

Cov. (X)		Cov. (X)
	Rkatsiteli	
20	overlapping fragment	22
33	fragment	32
53		55
19	jumping	29
16	long jumping	12
36		40
	Kishmish Vatkana	
28	overlapping fragment	29
35	fragment	23
62		52
16	jumping	24
14	long jumping	10
30		35
	Sangiovese	
55	overlapping fragment	33
55	fragment	34
23		68
10	jumping	23
33	long jumping	22
	Cov. (X) 20 33 53 19 16 36 28 35 62 16 16 14 30 4 30 55 55 55 23 10 33	Cov. (X)RkatsiteliRkatsiteli200overlapping fragment331fragment332jumping193jumping194long jumping363364verlapping fragment375fragment384jumping395jumping396397jumping398overlapping fragment399jumping301jumping302fragment303jemping fragment304jumping fragment305fragment306jumping fragment307jumping fragment308jumping fragment309jumping310jumping



	Cov. (X)		Cov. (X)
PN40024		RPV3	
overlapping	11	overlapping fragment	29
fragment		fragment	30
fragment	66		59
	77	jumping	18
jumping	15	long jumping	10
	15		29

- ~30X Overlapping fragment libraries a)
- b) ~20X Non-overlapping fragment libraries
- **c**) ~20X 3kb Jumping libraries + ~10X **10kb** (long) Jumping libraries



Thanks to Gabriele Magris

Michele Vidotto

### Fragment libraries – TruSeq Vs Nextera





### TruSeq PCR-Free



#### Michele Vidotto

#### Thanks to Irena Iurman and Nicoletta Felice

# Jumping libraries – contamination problem



Thanks to Irena Iurman and Nicoletta Felice

### Jumping libraries – duplication problem



	READS	% DUPLICATION	READ PAIR DUPLICATES	OPTICAL DUPLICATES PAIR	READS TO REACH 0.8 UNIQ FRAC (NB)	ESTIMATED LIBRARY SIZE	PRED FRAC OBSERVED (NB)
a)	5,399,709	4.80%	258,945	104,049	27,910,043	88,751,341	1.20%
b)	50,520,592	28.86%	14,581,610	5,564,168	30,297,840	96,516,581	9.60%

•Poisson model: the Poisson model assumes that each molecule in the library has an equal probability of being sequenced

•Negative binomial model: the NB model assumes there is variability in the probability of each molecule being sequenced, and that this variability follows a gamma distribution

•Log-series distribution: the LSD is a limiting case of the NB model where the gamma distribution has infinite variance, so there are an infinite number of molecules that are infinitely unlikely to be sequenced

### Jumping libraries – insert distribution problem

FR

15000

FR

25000

30000

35000



3Kb insert size always good!



Insert Size Histogram for All\_Reads in file exp\_397\_cabernet-franc\_GCCAAT\_L001\_unclas-pair\_U.bam



Michele Vidotto

# The high heterozygous grape genome

- Diploid genome
- PN4004 draft 487 Mb

Highly heterozygous individuals



• 19 chr

















#### michele.vidotto@uniud.it

### Assemblies results

	traminer		gouais blanc		caberne	et franc	pn40024	
Genome size estimated (bp)	394,44	1,994	420,689,048		483,072,260		444,716,040	
Estimated CN=1 (bp)	71.1	L%	69.6%		59.6%		65.0 %	
Estimated CN>1 (bp)	28.9	9%	30.4	4%	40.	40.4%		) %
Coverage estimated (X)	53	3	58	8	5	50		1
	contigs	scaffolds	contigs	scaffolds	contigs	scaffolds	contigs	scaffolds
Total (bp)	412,472,184	678,662,631	415,882,525	743,730,069	487,795,456	711,707,043	377,475,581	521,491,794
Average (bp)	3,728	78,223	3,722	77,764	4,770	57,005	6,814	86,684
Max (bp)	137,674	2,817,510	113,144	2,522,649	180,444	1,496,800	148,414	1,727,814
Min (bp)	32	885	69	905	58	894	118	885
Sequences (#)	110,653	8,676	111,742	9,564	102,258	12,485	55,399	6,016
N50 (#)	12,360	484	12,724	524	9,708	1,178	7,151	553
L50 (bp)	7,841	413,025	7,541	408,071	11,145	178,391	14,075	264,099
Gaps (#)		67,484		66,217		55,190		40,378
Median gap length (bp)		4,012		5,545		4,405		3,825

	rkatsiteli		kishmish vatkana		sangio	ovese	rpv3		
Genome size estimated (bp)	407,30	2,602	454,927,661		420,525,421		435,235,299		
Estimated CN=1 (bp)	69.3	1 %	56.5	5 %	68.4	4 %	58.6 %		
Estimated CN>1 (bp)	30.9	9%	43.5	5 %	31.6	5 %	41.4 %		
Coverage estimated (X)	5	2	46		64		54		
	contigs	scaffolds	contigs	scaffolds	contigs	scaffolds	contigs	scaffolds	
Total (bp)	426,363,529	753,744,202	410,415,417	610,410,426	412,308,934	691,235,812	397,904,367	608,826,013	
Average (bp)	3,739	74,710	3,803	71,426	4,044	44,121	3,560	69,240	
Max (bp)	135,733	2,436,454	119,547	2,355,113	184,133	2,528,353	90,275	2,693,357	
Min (bp)	46	889	27	895	76	881	46	894	
Sequences (#)	114,026	10,089	107,916	8,546	101,962	15,667	111,757	8,793	
N50 (#)	14,103	612	12,979	457	10,243	627	15,087	410	
L50 (bp)	7,250	352,572	7,751	364,852	8,975	301,183	7,032	425,483	
Gaps (#)		69,862		64,039		61,417		76,094	
Median gap length (bp)		5,200		5,728		3,581		1,770	

### Michele Vidotto

### Assemblies evaluation – kmers comparison



Michele Vidotto

### Assemblies evaluation – kmers comparison



#### Michele Vidotto

### Assemblies evaluation – alignment on the reference



						kishmish		
	traminer	gouais blanc	cabernet franc	pn40024	rkatsiteli	vatkana	sangiovese	rpv3
genes:	86.61%	85.96%	88.17%	86.88%	86.07%	84.16%	85.47%	85.45%
exons:	93.52%	92.89%	94.76%	94.43%	92.68%	91.91%	93.14%	92.62%
repeats:	81.74%	80.38%	83.88%	79.66%	82.16%	79.17%	79.89%	79.95%
Inc antisense:	90.10%	89.85%	91.08%	92.79%	89.73%	86.92%	89.61%	88.88%
Inc intergenic:	81.60%	79.69%	83.11%	82.45%	81.35%	76.02%	78.28%	79.14%
Inc intronic:	92.02%	90.30%	93.73%	93.35%	91.06%	90.68%	90.68%	88.78%

#### Michele Vidotto

### Assemblies evaluation – fragment reads relignment



		Gouais	Cabernet		Kishmish			~	$\land$	
	Traminer	Blanc	Franc	Rkatsiteli	Vatkana	Sangiovese	PN40024	– <del>د</del>		
Proper										
pairs (%)	94.05	93.76	96.45	95.47	94.15	89.33	94.03	d 4 -		
	Per base coverage			Per base coverage			Der henn geworen	N -		

٩W







Sangiovese



150

PN40024

#### Michele Vidotto

### High emizygosity emerges

Per base coverage



	Traminer		Gouais Blanc		Cabernet Franc		Rkatsiteli		Sangiovese	
	% of	% total	% of	% total	% of	% total	% of	% total	% of	% total
	contigs	length	contigs	length	contigs	length	contigs	length	contigs	length
hemizygous	63.55%	40.36%	49.21%	27.46%	67.81%	41.46%	60.32%	35.95%	39.75%	20.62%
non-hemizygous	36.24%	59.55%	49.44%	72.14%	31.97%	58.46%	39.37%	63.91%	60.13%	79.33%
low coverage	0.06%	0.03%	1.26%	0.37%	0.08%	0.02%	0.15%	0.07%	0.04%	0.01%
uncovered	0.14%	0.06%	0.09%	0.03%	0.15%	0.05%	0.16%	0.07%	0.09%	0.03%

### Michele Vidotto

### Emizygosity contigs are heterozygous and repetitive



michele.vidotto@uniud.it

### Aknowledgment

Prof. Michele Morgante Fabio Marroni Gabriele Magris Sara Pinosio

### IGATS

Davide Scaglione Simone Scalabrin Irena Iurman Vittorio Zamboni

### **UniUDLab and Administration**

GiusiZaina Nicoletta Felice Cristian DelFabbro











**European Research Council** 

Supported by the ERC project NOVABREED -Novel variation in plant breeding and the plant pan-genomes (Grant agreement no.: 294780)

### K-mer frequency based error correction algorithm



### Assembly metrics

1. Order contigs/scaffolds in descreasing length

2.Calculate contigs total length (TOT\_LEN)

**N50** = first *n* given  $\sum_{i=1}^{n} length(i) \ge TOT_{LEN}/2$ 

**L50** = length of the last contig/scaffold added to cover 50% or more of the assembly

**NG50** = first *n* then cui  $\sum_{i=1}^{n} length(i) \ge GENOME_{LEN}/2$ 

LG50 = ...



# The ALLPATHS-LG receipt

- a) ~45X Overlapping fragment libraries
- b) ~45X Jumping libraries
- c) ~1X Long Jumping libraries



- b) ~20X Non-overlapping fragment libraries
- c) ~20X 3kb Jumping libraries + ~10X 10kb (Long?) Jumping libraries

