1 De novo phased assembly of the Vitis riparia grape genome

- 2
- 3 Nabil Girollet¹, Bernadette Rubio^{1,2} and Pierre-François Bert¹
- 4

5 Affiliations

- ¹EGFV, Bordeaux Sciences Agro INRA Université de Bordeaux, ISVV, 210 chemin de
- 7 Leysotte, 33882 Villenave d'Ornon, France

⁸ ²IFV, Institut Français de la Vigne et du Vin, Domaine de l'Espîguette, 30240 Le Grau du Roi,

9 France

10 Corresponding Author: Pierre-François Bert (pierre-francois.bert@u-bordeaux.fr)

11

12 Abstract

13 Grapevine is one of the most important fruit species in the world. In order to better 14 understand genetic basis of traits variation and facilitate the breeding of new genotypes, we 15 sequenced, assembled, and annotated the genome of the American native Vitis riparia, one 16 of the main species used worldwide for rootstock and scion breeding. A total of 164 Gb raw 17 DNA reads were obtained from Vitis riparia resulting in a 225X depth of coverage. We 18 generated a genome assembly of the V. riparia grape de novo using the PacBio long-reads 19 that was phased with the 10x Genomics Chromium linked-reads. At the chromosome level, a 20 500 Mb genome was generated with a scaffold N50 size of 1 Mb. More than 34% of the 21 whole genome were identified as repeat sequences, and 37,207 protein-coding genes were 22 predicted. This genome assembly sets the stage for comparative genomic analysis of the 23 diversification and adaptation of grapevine and will provide a solid resource for further 24 genetic analysis and breeding of this economically important species.

25

26 Background & summary

Since few decades and the development of sequencing technologies, the number of species
whose genome has been totally sequenced has increased exponentially. There is a large

variability for the quality of all the sequences assemblies. In 2017, 72 plant reference quality genome assemblies were reported in NCBI¹. For plant breeding, the availability of a contiguous genome sequence provides a tool to better identify genes underlying traits and how they may be regulated by various environmental parameters in different genetic backgrounds. At the simplest, it allows for association of genetic markers for selection and introgression of traits across germplasm to enable the development of novel products for consumers^{2 3}.

As an important crop, Vitis vinifera was one of the first higher plant species whose genome 36 was sequenced by a French-Italian consortium⁴. The consortium decided to sequence a near 37 38 homozygous V. vinifera cultivar related to Pinot Noir (PN40024) in order to facilitate the 39 sequence assembly by limiting sequence variability. To date, this genome still stands as the 40 reference for the grapevine community, but grapevine intra species and interspecies diversity 41 makes using a single reference genome inadequate for studying the function of other 42 genotypes. In order to address the variations in a cultivated V. vinifera variety, the Pinot Noir 43 genome was sequenced using Sanger sequencing providing a high quality draft of the 44 genome with about 10X coverage⁵. Next Generation Sequencing reads are too short to 45 resolve abundant repeats in particular in plants genome, leading to incomplete or ambiguous 46 assemblies⁶. Few attempts to produce high quality grapevine genomes were undertaken in 47 grapevine and produced valuable data to study the genetic variations of V. vinifera cv. Tannat⁷ and cv. Thompson seedless⁸ through comparison with the reference genomes. 48

The last few years have seen rapid innovations in sequencing technologies and improvement in assembly algorithms that enabled the creation of highly contiguous genomes. The development of third generation sequencing technologies that deliver long reads from single molecules and carry the necessary information to phase haplotypes over several kilobases have greatly improved the feasibility of *de novo* assemblies^{9 10 11}. Sequences of *V. vinifera* cv. Cabernet Sauvignon were first released¹² using PacBio sequencing and FALCON, and FALCON-Unzip pipeline¹². This generated a 591 Mbp haplotype genome from a set of 718

56 primary contigs, and a set of correlated 2,037 haplotigs spanning 367 Mbp. The total p-contig 57 size was larger than the estimated genome size of V. vinifera (~500 Mbp) suggesting that in 58 some cases FALCON-Unzip underestimated the alternative haplotype sequences because of high heterozygosity between homologous regions, which is common in grapevine^{13 14}. Later, 59 60 the PacBio assembly and annotation of V. vinifera cv Chardonnay variety provided after 61 curation of artefactual contig assignment, 854 p-tigs and 1883 h-tigs, totaling 490 Mb and 378 Mb¹⁵. More recently, another version of the Chardonnay genome was proposed with a 62 different level of curation at 605 Mb¹⁶. 63

64 An evaluation of genetic diversity based on a panel of 783 V. vinifera varieties using 10K 65 SNPs revealed a high level of diversity (He = 0.32) and confirmed the close pedigree 66 relationship within the cultivated grapevine due to the wide use of the most interesting parents during domestication and early selection by humans¹⁷. Considering that grape 67 68 cultivation currently faces severe pathogen pressures and climate change, we assume that 69 the exploitation of the natural genetic diversity may ensure the long-term sustainability of the grape and wine industries¹⁸. Grapes belong to the genus *Vitis*, which includes over 60 inter-70 71 fertile species. The most common grape cultivars derive their entire ancestry from the 72 species V. vinifera, but wild relatives have also been exploited to create hybrid cultivars, often with increased disease resistances¹⁹. 73

To date, no wild *Vitis* genomes have been released so far and the only whole genome sequences for grape are from *V. vinifera* varieties and yet there is a clear need for genetic resources ²⁰. Here, we report the first *de novo* assembly and genome annotation of the North American native grape *V. riparia*. Using the latest sequencing technologies, we show that 10x Genomics Chromium data can be combined with long read PacBio sequencing to effectively determine genome phasing. The phased haplotypes of *V. riparia* genome will greatly contribute to give more insight into the functional consequences of genetic variants.

81 Methods

82 Sample collection, library construction and whole genome sequencing

83 The Vitis riparia Gloire de Montpellier (RGM) selection was obtained in 1880 by L. Vialla and 84 R. Michel from North American collections and is the only commercially available pure V. 85 riparia stock. RGM clone #1030 and the European native Vitis vinifera Cabernet sauvignon 86 (CS) clone #15 were grown at INRA, Bordeaux (France). A F1 segregating population of 114 individuals named CSxRGM1995-1 was derived from the cross between CS and RGM²¹. 87 This population was genotyped using the GBS approach²² to create a high resolution genetic 88 89 map to assist in anchoring and orienting the assembled V. riparia genome scaffolds. 90 Total DNA was isolated and extracted using QIAGEN Genomic-tips 100/G kit (Cat No./ID:

91 10243) following the tissue protocol extraction. Briefly, 1g of young leaf material was ground 92 in liquid nitrogen with mortar and pestle. After 3h of lysis and one centrifugation step, the 93 DNA was immobilized on the column. After several washing steps, DNA is eluted from the 94 column, then desalted and concentrated by alcohol precipitation. The pellet is resuspended 95 in TE buffer.

96 Three PacBio libraries with a 20-kb insert size were also constructed and sequenced on RSII 97 platforms (97.71 Gb data; ~118-fold covering), following the standard PacBio protocol of 98 Sequencing Kit 1.2.1 (Pacific Biosciences, USA). Four 10x Chromium Genomics libraries were constructed using the Chromium[™] Genome Solution (10X Genomics, USA), and 2x150 99 100 bp sequenced on Illumina HiSeq3000, producing ~350 million paired-end linked-reads (~ 101 107-fold covering). Finally, 2 libraries for 2x100 bp sequencing were built with different insert 102 sizes: 500 bp for paired-end (PE) and 6 kb for mate-pair (MP), based on the standard 103 Illumina protocol and sequenced on the Illumina HiSeq 2500. The raw reads were trimmed 104 before being used for subsequent genome assembly. For Illumina HiSeq sequencing, the 105 adaptor sequences, the reads containing more than 10% ambiguous nucleotides, as well as 106 the reads containing more than 20% low-quality nucleotides (quality score less than 5), were 107 all removed. After data cleaning and data preprocessing, we obtained a total of 164 Gb of 108 clean data (52 Gb PacBio data, 59 Gb 10X Genomics, 33 Gb PE reads and 20 Gb MP 109 reads,), representing 331X coverage of the V. riparia genome (Table 1).

110

111 Genome size and heterozygosity estimation

112 Lodhi and Reisch³⁵ estimated the genome size in grape to be approximately 475 Mb based 113 on measurements using flow cytometry for 19 species including wild Vitis species, V.vinifera 114 and V. labrusca cultivars. The measurements showed intraspecific variation in genome size 115 between different varieties of Vitis vinifera ranging from 1C = 415 to 511 Mb, and between 116 different North America Vitis species ranging from 1C = 411 to 541 Mb, with V. riparia around 117 470 Mb. Genome sequencing of different V. vinifera varieties gave values in the same range 118 or greater depending on the methods of sequencing and assembly. In order to verify these values, we estimated genome size of V. riparia by the k-mer method^{36 37} using data from pair-119 120 end and mate-pair Illumina sequencing. By analyzing the 21-mers depth distribution, a total 121 of ~50 billion k-mers were estimated with a peak frequency of 100, corresponding to a 122 genome size of 494 Mb and the estimated repeat sequencing ratio was 33.74%. In this 123 study, V. riparia heterozygosity was estimated to be 0.46% (mean distance 1 SNP each 217 124 bp between heterozygous SNPs) from 10x Chromium Genomics data processing.

125

126 De novo Genome assembly and scaffolding of the Vitis riparia genome

127 We employed a hybrid *de novo* whole-genome assembly strategy, combining both short 128 linked-reads and PacBio long reads data. Genome assembly was first performed on full PacBio cleaned reads using FALCON v0.3.0³⁸. Error correction and pre-assembly were 129 130 carried out with the FALCON/FALCON Unzip pipeline after evaluating the outcomes of using 131 different parameters in FALCON during the pre-assembly process. Based on the contig N50 132 results, a length cutoff of 5kb and a length cutoff pr of 8kb for the assembly step were ultimately chosen. The draft assembly was polished using Quiver³⁹, which mapped the 133 PacBio reads to the assembled genome with the BLASER pipeline⁴⁰. Haplotypes were 134 135 separated during assembly using FALCON-Unzip and the preliminary genome assembly was 136 approximately 530 Mb (1,964 primary-contigs) and 317 Mb (3,344 haplotigs). A summary of 137 the assembly statistics can be found in Table 1. Assembly was then processed with Purge

Haplotigs¹³ to investigate the proper assignment of contigs, followed by 2 rounds of polishing
to correct residual SNP and INDELs errors with Pilon v1.22 software⁴¹ using high-coverage
(~106X) Illumina paired-end and mate pair data.

The 10x Chromium Genomics linked-reads were used to produce a separate *V. riparia* assembly using the Supernova assembler option *--style=pseudohap2* and created two parallel pseudohaplotypes⁴². The mean input DNA molecule length reported by the Supernova assembler was 45kb and the assembled genome size was 424Mb with a N50 scaffold of 711kb.

Subsequently, the PacBio assembly was scaffolded with the 10x Chromium Genomics one using the hybrid assembler LINKS⁴³ with 7 iterations, producing 870 scaffolds spanning 500 Mb with N50 = 964 kb and L50 = 255 (Table 2). Finally, genome phasing was reconstituted using Long Ranger analysis pipeline that processes Chromium sequencing output to align reads and call and phase SNPs, indels, and structural variants on the basis of molecular barcodes information.

152

153 Genotyping by Sequencing and genetic mapping

154 Two 96-plex GBS libraries (Keygene N.V. owns patents and patent applications protecting its 155 Sequence Based Genotyping technologies) were constructed for the two parents (two 156 replicates for each) and the 114 F1 plants of the cross CS x RGM. Raw reads were 157 checked with FastQC²³, demultiplexed with a custom script and cleaned with CutAdapt²⁴. 158 Cleaned reads were then mapped to the V. riparia RGM scaffolds previously obtained, the V. vinifera Cabernet Sauvignon contigs¹² and V. vinifera PN40024 genome assemblies⁴ for SNP 159 calling. Aligned on these genomes were performed using BWA²⁵. SAMtools²⁶ and Picard 160 tools²⁷ and SNP genotypes were detected with GATK²⁸ using the *hardfilter* parameters²⁹. In 161 162 the variant call format (VCF) output file only sites with less than 20 % missing data and a minimum allele frequency (MAF) □≥ □ 0.2 were retained. The SNP set was parsed into two 163 data sets based on a pseudo-test cross mapping strategy³⁰ using major_minor and 164 get_pseudo_test_cross scripts from Hetmapps³¹. The segregation ratios of markers in the 165

population were examined by Chi-square analysis. Markers with segregation ratios that differed from expected 1:1 at P<0.05 were classified as segregation distortion markers and discarded. The RGM and CS sets contain 1591 and 2359 SNPs respectively. Linkage groups (LGs) were determined using software JoinMap® 4.1^{32 33} and Rqtl³⁴. LG were formed with a logarithm of odds (LOD) threshold of 6 and a maximum recombination frequency of 0.45. The 19 LGs that corresponded to the 19 chromosomes of grapevine were reconstructed and leaded to a total genetic map length of 2,268 cM and 2,514 cM for RGM and CS respectively.

174 **Pseudo-molecule construction**

175 The PacBio / 10x Chromium Genomics hybrid scaffolding was organized into pseudo-176 molecules using GBS markers information from the CS x RGM genetic map. Scaffolds were anchored and oriented SNP using AllMaps⁴⁴ with the unequal weights2 parameters for a 177 178 single run for the entire genome. Final pseudo-molecules were named according to Vitis 179 vinifera PN40024 reference genome using SNP identification through SNP calling on this 180 reference. Since PN40024 genome is the only one available who has been scaffolded into pseudo-molecules, collinearity with *V. riparia* was evaluated using D-GENIES⁴⁵ and showed 181 182 extremely high conservation along the 19 chromosomes of the species (Figure 1) even if the North American and Eurasian Vitis species diverged approximately 46.9 million years ago⁴⁶. 183

184

185 **Genome annotation and gene prediction**

186 Consistent with observations that long reads sequencing technologies are a better solution 187 for resolving repeat sequences, we found that known repetitive elements accounted for 170 188 Mb (33.94%) of the genome in V. riparia. This is a lower proportion among grape genomes 189 when comparing published values to date. However, when comparisons are performed with the same analysis workflow and tools^{47 48}, the percentages obtained between the two 190 191 genotypes were in the same range (Online-only Table 1). Similar to other grape genomes, 192 long terminal repeat (LTR) elements constituted the highest proportion of all repeated 193 elements in V. riparia, (21.44%) with Copia and Gypsy families accounting for 8.33% and 194 12.66% respectively. The Long Interspersed Nuclear Elements (LINEs) and Miniature 195 Inverted-repeat Transposable Elements (MITEs) represented 3.61% and 6.02% of the whole 196 genome respectively.

After repeat masking, the genome was ab initio annotated using MAKER-P pipeline^{49 50}, 197 SNAP⁵¹ and Augustus⁵² gene finder with 3 rounds of Maker and an Augustus prediction. 198 199 Structural annotation was then followed with an Interproscan functional annotation and 200 putative gene function assignation using BLAST on UniProtKB. MAKER-P quality metrics with a threshold of AED<0.5 were chosen to retain the set of predicted genes. We finally 201 202 of 37,207 protein-coding generated а gene set genes (11,434_AED<0.1; 203 8,638_0.1≤AED<0.2; 5,748_0.2≤AED<0.3; 5,418_0.3≤AED<0.4; 5,969_0.4≤AED<0.5) with

204 31,240 of them coupled with an evidence of protein function.

205

To facilitate genomic investigations for the community, a JBrowse Genome Browser⁵³ was set up for *V. riparia* pseudo-molecules and is available from https://www6.bordeauxaquitaine.inra.fr/egfv/.

209

210 Code Availability

- 211 1. GBS demultiplexing
- 212 https://github.com/timflutre
- 213 2. Filters FASTQ files with CASAVA 2.20
- 214 fastq_illumina_filter --keep N -v -v -o good_reads.fq raw_reads.fastq
- 215 3. Cutadapt (regular 3' adapter)
- 216 https://cutadapt.readthedocs.io/en/stable/guide.html
- 217 cutadapt -a AGATCGGAAGAGCGGTTCAGCAGGAATGCCGAG
- 218 -A AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT
- 219 -G CTCGGCATTCCTGCTGAACCGCTCTTCCGATCT
- 220 -g ACACTCTTTCCCTACACGACGCTCTTCCGATCT -u7 -U7 -m10

- 4. Burrows-Wheeler Alignment BWA-MEM
- 222 http://bio-bwa.sourceforge.net/bwa.shtml
- bwa mem ref.fa read1.fastq.gz read2.fastq.gz > aligned.reads.sam with these options :
- -M Mark shorter split hits as secondary (for Picard compatibility)
- -R Complete read group header line with '\t' used in STR to be converted to a TEB in the
- 226 output SAM. An example is '@RG\tID:\tSM:\tPL:\tLB:'
- 5. Picard tools
- 228 https://broadinstitute.github.io/picard/
- 229 SortSam : java jar picard.jar SortSam with these options: INPUT (BAM file), OUTPUT (BAM
- 230 file), SORT_ORDER
- 231 MarkDuplicates : java jar picard.jar MarkDuplicates with these options: INPUT (BAM file),
- 232 OUTPUT (BAM file), METRIC_FILE (file)
- 233 BuildBamIndex : java jar picard.jar BuildBamIndex with these options: INPUT (BAM file)
- 234 6. GATK tools
- 235 HaplotypeCaller : java -jar GenomeAnalysisTK.jar -T HaplotypeCaller -R ref.fasta -I
- 236 file.bam -genotyping_mode DISCOVERY -drf DuplicateRead -emitRefConfidence GVCF -o
- 237 file.g.vcf
- 238 https://software.broadinstitute.org/gatk/documentation/tooldocs/3.8-
- 239 0/org_broadinstitute_gatk_tools_walkers_haplotypecaller_HaplotypeCaller.php
- 240 CombineGVCFs : java -jar GenomeAnalysisTK.jar -T CombineGVCFs -R ref.fasta -drf
- 241 DuplicateRead -G Standard -G AS_Standard --variant sample1 to sample'n'.g.vcf -o
- 242 cohort_file.g.vcf
- 243 https://software.broadinstitute.org/gatk/documentation/tooldocs/3.8-
- 244 0/org_broadinstitute_gatk_tools_walkers_variantutils_CombineGVCFs.php
- 245 GenotypeGVCFs : java -jar GenomeAnalysisTK.jar -T GenotypeGVCFs -R ref.fasta -drf
- 246 DuplicateRead –G Standard –G AS_Standard --variant cohort_file.g.vcf –o final_file.vcf
- 247 https://software.broadinstitute.org/gatk/documentation/tooldocs/3.8-
- 248 0/org_broadinstitute_gatk_tools_walkers_variantutils_GenotypeGVCFs.php

- 249 SelectVariants : java -jar GenomeAnalysisTK.jar -T SelectVariants -R ref.fasta -V
- 250 final_file.vcf -selectType SNP -o file_snps.vcf
- 251 https://software.broadinstitute.org/gatk/documentation/tooldocs/3.8-
- 252 0/org_broadinstitute_gatk_tools_walkers_variantutils_SelectVariants.php
- 253 VariantFiltration : java -jar GenomeAnalysisTK -T VariantFiltration -R ref.fasta -V
- 254 file_snps.vcf --filterExpression « QD < 2.0 || FS > 60.0 || MQ < 40.0 || MQRankSum < -12.5 ||
- 255 ReadPosRankSum < -8.0 » --filteredName «FILTER » -o filtered_snps.vcf
- 256 https://software.broadinstitute.org/gatk/documentation/tooldocs/3.8-
- 257 0/org_broadinstitute_gatk_tools_walkers_filters_VariantFiltration.php
- 258 7. VCF filtering
- 259 vcftools --vcf filtered_snps.vcf --remove-filtered-all --recode --out filteredFinal_snps.vcf
- 260 8. Falcon and Falcon_Unzip Assembly for SMRT sequencing
- 261 https://github.com/PacificBiosciences/FALCON/wiki
- 262 https://github.com/PacificBiosciences/FALCON_unzip/wiki
- 263 Main parameters: length_cutoff = 5000, length_cutoff_pr = 5000
- 264 pa_HPCdaligner_option = -v -dal128 -e0.70 -M40 -l2500 -k17 -h500 -w7 -s100
- 265 ovlp_HPCdaligner_option = -v -dal128 -M40 -k19 -h500 -e.96 -l1500 -s100
- 266 pa_DBsplit_option = -a -x500 -s200
- 267 ovlp_DBsplit_option = -s200
- 268 falcon_sense_option = --output_multi --output_dformat --min_idt 0.80 --min_cov 4
- 269 max_n_read 400 --n_core 16
- 270 falcon_sense_skip_contained = False
- 271 overlap_filtering_setting = --max_diff 120 --max_cov 120 --min_cov 4 --n_core 24
- 272 9. Purge Haplotigs
- 273 https://bitbucket.org/mroachawri/purge_haplotigs/src/master/
- 274 purge_haplotigs readhist -b aligned.bam -g genome.fasta
- 275 10. Supernova Assembly for 10x Chromium sequencing
- 276 https://support.10xgenomics.com/de-novo-assembly/software/overview/latest/welcome

277 Option pseudohap2 style output

- 11. Scaffolding Falcon assembly with LINKS using Supernova outputs Assembly
- 279 https://github.com/bcgsc/LINKS
- 280 LINKS -f.fa -s fileofname.fofn -b cns1-linked_draft -d 5000 -t 100 -k 19 -l 5 -a 0.3
- 281 LINKS -f.fa -s fileofname.fofn -b cns2-linked_draft -d 6000 -t 80 -k 19 -l 15 -a 0.3
- 282 LINKS -f.fa -s fileofname.fofn -b cns3-linked_draft -d 7000 -t 60 -k 19 -l 20 -a 0.3
- 283 LINKS -f.fa -s fileofname.fofn -b cns4-linked_draft -d 10000 -t 30 -k 19 -l 20 -a 0.3
- LINKS -f.fa -s fileofname.fofn -b cns5-linked_draft -d 15000 -t 30 -k 19 -l 20 -a 0.3
- 285 LINKS -f.fa -s fileofname.fofn -b cns6-linked_draft -d 50000 -t 30 -k 19 -l 30 -a 0.3
- 286 LINKS -f.fa -s fileofname,fofn -b cns7-linked_draft -d 75000 -t 30 -k 19 -l 40 -a 0.3
- 287 12. Improving quality with PILON and Illumina sequencing
- 288 https://github.com/broadinstitute/pilon/wiki/Requirements-&-Usage
- 289 13. Allmaps pseudomolecules scaffolding
- 290 https://github.com/tanghaibao/jcvi/wiki/ALLMAPS
- 291 14. Assembly evaluation with BUSCO v3
- 292 https://busco.ezlab.org/
- 293 15. Vitis TE(s) Identification using RepeatMasker
- 294 http://www.repeatmasker.org/
- 295 16. Annotation with MAKER_P pipeline, SNAP and Augustus gene finder
- 296 http://www.yandell-lab.org/publications/pdf/maker_current_protocols.pdf
- 297 https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-5-59
- 298 https://github.com/Gaius-Augustus/Augustus
- 299 DB (vitis) AND "Vitis"[porgn] from https://www.ncbi.nlm.nih.gov
- 300 EST DB (vitis) AND "Vitis"[porgn] from https://www.ncbi.nlm.nih.gov
- 301 First run: rm_pass=0, est2genome=1 and protein2genome=1
- 302 gff3_merge -d master_datastore_index.log
- 303 maker2zff -c 0 -e 0 -o 0 -x 0.05 maker1.gff
- 304 fathom -categorize 1000 genome.ann genome.dna

- 305 fathom -export 1000 -plus uni.ann uni.dna
- 306 forge export.ann export.dna
- 307 hmm-assembler.pl RGM . > snap1.hmm
- 308 Second run: rm_pass=1, est2genome=0, protein2genome=0, maker_gff=maker1.gff ,
- snaphmm=snap1.hmm leading to a maker2.gff3 and a snap2.hmm files.
- 310 gff3_merge -d master_datastore_index.log
- 311 maker2zff -c 0 -e 0 -o 0 -x 0.05 maker2.gff
- 312 fathom -categorize 1000 genome.ann genome.dna
- 313 fathom -export 1000 -plus uni.ann uni.dna
- 314 forge export.ann export.dna
- 315 hmm-assembler.pl RGM . > snap2.hmm
- 316 Run Augustus:
- 317 zff2gff3.pl genome.ann | perl -plne 's/t(\S+)\$/t\.\t\$1/' >genome.gff3
- 318 autoAug.pl --genome=../pilon2.fasta --species=RGM18 --cdna=sequence_est_ncbi.fasta -
- 319 trainingset=genome.gff3 --singleCPU -v --useexisting
- 320 Third run: rm_pass=1, est2genome=0, protein2genome=0, maker_gff=maker2.gff ,
- 321 snaphmm=snap2.hmm, augustus_species=RGM18 leading to a maker3.gff3,
- 322 maker3.transcripts.fasta and maker3.proteins.fasta structural prediction.
- 323 gff3_merge -d master_datastore_index.log
- 324 fasta_merge -d master_datastore_index.log
- 17. Interproscan functional annotation and putative gene function assignation
- 326 Download protein DB from http://www.uniprot.org
- 327 makeblastdb -in protein_db.fasta -input_type fasta -dbtype prot
- 328 blastp -db protein_db.fasta -query maker3.proteins.fasta -out maker3.proteins.blastp -evalue
- 329 0.000001 -outfmt 6 -max_hsps 1
- 330 maker_functional_gff protein_db.fasta maker3.proteins.blastp maker3.gff3 >>
- 331 maker3.putative.gff3

- 332 maker functional fasta protein db.fasta maker3.proteins.blastp maker3.proteins.fasta >>
- 333 maker3.putative.proteins.fa
- 334 maker_functional_fasta protein_db.fasta maker3.proteins.blastp maker3.transcripts.fasta >>
- 335 maker3.putative.transcripts.fa
- 336 Run Interproscan
- 337 interproscan.sh --iprlookup --goterms -f tsv -i maker3.putative.proteins.fa -pa -b
- 338 RGM.annotated.proteins
- 18. Assembly validation using WBA mem
- 340 bwa mem -M -t 20 VitRiparia.fasta reads_pe.R1.fastq reads_pe.R2.fastq >
- 341 aln_pe_reads.sam
- 342 samtools view -bS aln_pe_reads.sam -o aln_pe_reads.bam #| samtools sort -
- 343 aln_reads.sorted.bam
- 344 samtools sort -o aln_pe_reads.sorted.bam aln_pe_reads.bam
- 345 bamtools stats -in aln_pe_reads.sorted.bam > bamstat_pe.reads
- 346

347 Data records

348 The V. riparia genome project was deposited at NCBI under BioProject number 349 PRJNA512170 and BioSample SAMN10662253. The DNA sequencing data from Illumina, 350 PacBio and 10x Genomics have been deposited in the Sequence Read Archive (SRA) 351 database under accession numbers SRX5189632 to SRX5189680. This Whole Genome 352 Shotgun project has been deposited at DDBJ/ENA/GenBank under the accession 353 SJAQ00000000. The version described in this paper is version SJAQ01000000. Genetic 354 mapping data and structural and functional annotation file of the Vitis riparia assembly are 355 available on figshare (https://figshare.com/s/0a52d4408214e9f1e280).

356

357 Technical validation

To evaluate the accuracy and completeness of the *V. riparia* assembly, genome features were compared to those of *V. vinifera* (Table 2). We found that both contig and scaffold N50

360 lengths of Vitis riparia reached considerable continuity. The Guanine-Cytosine content (GC =

361 34.32 %) was similar to those of *V. vinifera* Chardonnay (34.43%).

To further assess the accuracy of the *V. riparia* genome assembly, the NGS-based short reads from whole-genome sequencing data were also aligned against the genome assembly using BWA mem⁵⁴. We found that 98.4% of the reads were reliably aligned to the genome assembly, and 95.8% of the reads were properly aligned to the genome with their mates. Paired-end reads data were not used during the contig assembly, thus the high alignment ratio demonstrated the high quality of contig assembly.

368 The assembled genome was also subjected to Benchmarking Universal Single-Copy Orthologs⁵⁵, which quantitatively assesses genome completeness using evolutionarily 369 370 informed expectations of gene content from near-universal single-copy orthologs, using the 371 genes in the embryophyta release 9 dataset (embryophyta.odb9). The BUSCO results 372 showed that 96.5% of conserved BUSCO proteins were detected in the V. riparia assembly, 373 including 1.1% of fragment BUSCO proteins (Table 2). Overall, these metrics compare well with other recently published grape genomes, providing a high quality genome sequences for 374 375 the following functional investigations.

376

377 References

- Peterson, D.G. & Arick, M. Sequencing Plant Genomes. (Progress in Botany. Springer,
 Berlin, Heidelberg 2018).
- Nguyen, K.L., Grondin, A., Courtois, B., & Gantet, P. Next-generation sequencing
 accelerates crop gene discovery. *Trends Plant Sci* 24, 263-274 (2018).
- 382 3. Scheben, A., Yuan, Y. & Edwards, D. Advances in genomics for adapting crops to climate
 383 change *Curr. Plant Biol.* 6, 2-10 (2016).
- 4. Jaillon, O. *et al.* The grapevine genome sequence suggests ancestral hexaploidization in
 major angiosperm phyla. *Nature* 449, 463-467 (2007).
- 5. Velasco, R. *et al.* A High Quality Draft Consensus Sequence of the Genome of a
 Heterozygous Grapevine Variety. *PLoS ONE* 2:e1326 (2007).

- 388 6. Alkan, C., Sajjadian, S. & Eichler, E.E. Limitations of next-generation genome sequence
- 389 assembly. *Nat. Methods* **8**,61-65 (2010).
- 390 7. Da Silva, C. et al. The high polyphenol content of grapevine cultivar Tannat berries is
- 391 conferred primarily by genes that are not shared with the reference genome. Plant Cell 25,
- 392 4777-4788 (2013).
- 393 8. Di Genova, A. et al. Whole genome comparison between table and wine grapes reveals a
- 394 comprehensive catalog of structural variants. *BMC Plant Biol* **14**, 7 (2014).
- 395 9. Jiao, W.B. & Schneeberger, K. The impact of third generation genomic technologies on
- 396 plant genome assembly. Curr. Opin. Plant Biol. 36, 64-70 (2017).
- 397 10. Li, C., Lin, F., An, D., Wang, W. & Huang, R. Genome Sequencing and Assembly by
- Long Reads in Plants. Genes 9, 6 (2018).
- 11. Treangen, T.J. & Salzberg, S.L. Repetitive DNA and next-generation sequencing:
 computational challenges and solutions. *Nat. Rev. Genet.* 13, 36-46 (2012).
- 401 12. Chin, C.S. *et al.* Phased diploid genome assembly with single-molecule real-time
 402 sequencing. *Nat. Methods* **13**, 1050-1054 (2016).
- 13. Roach, M. J., Schmidt, S. A., & Borneman, A. R. Purge Haplotigs: allelic contig
 reassignment for third-gen diploid genome assemblies. *BMC bioinformatics* **19**, 460 (2018).
- 405 14. Vinson, J.P. *et al.* Assembly of polymorphic genomes: algorithms and application to
 406 *Ciona savignyi. Genome Research* 15, 1127–35 (2005).
- 407 15. Roach, M.J. et al. Population sequencing reveals clonal diversity and ancestral
- 408 inbreeding in the grapevine cultivar Chardonnay. *PLoS Genet* **14**:e1007807 (2018)
- 409 16. Zhou, Y.S. et al. Structural variants, clonal propagation, and genome evolution in
- 410 grapevine (*Vitis vinifera*) bioRxiv 508119; doi: https://doi.org/10.1101/508119 (2018).
- 411 17. Laucou V. et al. Extended diversity analysis of cultivated grapevine Vitis vinifera with 10K
- 412 genome-wide SNPs. PLoS One 13:e0192540 (2018).
- 18. Myles, S. et al. Genetic structure and domestication history of the grape. Proc Natl Acad
- 414 *Sci U S A* **108**, 3530-3535 (2011).

- 415 19. Migicovsky, Z. et al. Genomic ancestry estimation quantifies use of wild species in grape
- 416 breeding. *BMC Genomics* **17**, 478 (2016).
- 417 20. FAO Commission on genetic resources for food and agriculture assessment. The state of
- 418 the world's biodiversity for food and agriculture (2019).
- 419 21. Marguerit, E. et al. Genetic dissection of sex determinism, inflorescence morphology and
- 420 downy mildew resistance in grapevine. *Theor. Appl. Genet.* **118**, 1261-1278 (2009).
- 421 22. Elshire, R. J. et al. A robust, simple genotyping-by-sequencing (GBS) approach for high
- 422 diversity species. *PLoS ONE* **6**:e19379 (2011).
- 423 23. Andrews, S. Fastqc: a quality control tool for high throughput sequence data,
- 424 http://www.bioinformatics.babraham.ac.uk/projects/fastqc (2010).
- 425 24. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing
- 426 reads. *EMBnet.journal* **17**, 10-12 (2011).
- 427 25. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25,
 428 2078–2079 (2009).
- 429 26. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows–Wheeler
 430 transform. *Bioinformatics* 26, 589–595 (2010).
- 431 27. Picard Tools By Broad Institute. Available from: http://broadinstitute.github.io/picard.
- 432 28. DePristo M.A. et al. A framework for variation discovery and genotyping using next-
- 433 generation DNA sequencing data. *Nat Genet* **43**, 491-498 (2011).
- 434 29. Van der Auwera, G.A., et al. From FastQ Data to High-Confidence Variant Calls: The
- 435 Genome Analysis Toolkit Best Practices Pipeline. Current Protocols in Bioinformatics 43, 1-
- 436 11 (2013).
- 437 30. Grattapaglia, D., Bertolucci, F.L.G & Sederoff, R. Genetic mapping of QTLs controlling
- 438 vegetative propagation in *Eucalyptus grandis* and *E. urophylla* using a pseudo-testcross
- 439 mapping strategy and RAPD markers. *Theor. Appl. Genet.* **90**, 933–947 (1995).
- 440 31. Hyma, K.E. et al. Heterozygous mapping strategy (HetMappS) for high resolution
- 441 Genotyping-By-Sequencing markers: a case study in grapevine. *PLoS One* **10**: e0134880
- 442 (2015).

- 443 32. Stam, P. & Van Ooijen J.W. JOINMAP version 2.0: software for the calculation of genetic
- 444 linkage maps (1995).
- 33. Van Ooijen, J.W. JoinMap® 4.0, Software for the calculation of genetic linkage maps in
- 446 experimental populations. Kyazma B.V. Wageningen, Netherlands (2006)
- 447 34. Broman, K.W., Wu, H., Sen, S., & Churchill G.A. R/qtl: QTL mapping in experimental
- 448 crosses. *Bioinformatics* **19**, 889-890 (2003).
- 449 35. Lodhi, M.A. & Reisch, B.I. Nuclear DNA content of Vitis species, cultivars, and other
- 450 genera of the Vitaceae. Theor. Appl. Genet. 90, 11-16 (1995).
- 451 36. Liu, B., Shi, Y., Yuan, Y., Hu, X., Zhang, H., Li, N., Li, Z., Chen, Y., Mu, D., and Fan, W.
- 452 Estimation of genomic characteristics by analyzing k-mer frequency in de novo genome
- 453 projects. arXiv preprint arXiv:1308 (2012).
- 454 37. Marcais, G. & Kingsford, K. A fast, lock-free approach for efficient parallel counting of
- 455 occurrences of k-mers. *Bioinformatics* **27**, 764-770 (2011).
- 456 38. PacBio FALCON GitHub page https://github.com/PacificBiosciences/FALCON, Accessed
- 457 18 Mar 2018.
- 458 39. Pacific Biosciences, SMRT tools. https://www.pacb.com/wp-content/uploads/SMRT-
- 459 Tools-Reference-Guide-v4.0.0.pdf.
- 460 40. Chaisson, M.J. & Tesler, G. Mapping single molecule sequencing reads using basic local
- 461 alignment with successive refinement (BLASR): application and theory. BMC Bioinformatics
- 462 **13**, 238 (2012).
- 463 41. Walker, B.J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection
 464 and genome assembly improvement. *PLoS One* **9**: e112963 (2014).
- 465 42. Weisenfeld, N.I., Kumar, V., Shah, P., Church, D.M. & Jaffe D.B. Direct determination of 466 diploid genome sequences. *Genome research* **27**, 757-767 (2017).
- 467 43. Warren, R.L. *et al.* LINKS: Scalable, alignment-free scaffolding of draft genomes with
 468 long reads. *Gigascience*. 4, 35 (2015).
- 469 44. Tang, H. et al. ALLMAPS: robust scaffold ordering based on multiple maps. Genome
- 470 *Biol.* **16**. 3-10 (2015).

- 471 45. Cabanettes, F. & Klopp, C. D-GENIES: dot plot large genomes in an interactive, efficient
- 472 and simple way. *PeerJ* **6**:e4958 (2018).
- 473 46. Ma, Z.Y. et al. Phylogenomics, biogeography, and adaptive radiation of grapes.
- 474 Molecular phylogenetics and evolution **129**, 258-267 (2018).
- 475 47. Tarailo-Graovac, M. & Chen, N. Using RepeatMasker to identify repetitive elements in
- 476 genomic sequences. *Curr Protoc Bioinformatics* **25**, 4-10 (2009).
- 477 48. Bao, W., Kojima, K.K. & Kohany, O. Repbase Update, a database of repetitive elements
- in eukaryotic genomes. *Mobile DNA* **6**:11 (2015).
- 479 49. Campbell, M.S., Holt, C., Moore, B. & Yandell, M. Genome Annotation and Curation
- 480 Using MAKER and MAKER-P. *Curr Protoc Bioinformatics* **48**, 1-39 (2014).
- 481 50. Campbell, M.S. et al. MAKER-P: a tool kit for the rapid creation, management, and
- 482 quality control of plant genome annotations. *Plant Physiol* **164**, 513–24. (2014).
- 483 51. Korf, I. Gene finding in novel genomes. *BMC Bioinformatics* **5**:59. (2004).
- 484 52. Stanke, M. et al. AUGUSTUS: a web server for gene finding in eukaryotes. Nucleic Acids
- 485 *Research* **32**, 309-312 (2004).
- 486 53. Buels, R. et al. JBrowse: a dynamic web platform for genome visualization and analysis.
- 487 Genome Biology **12**, 17-66 (2016).

orthologs. Bioinformatics 31, 3210-3212 (2015).

- 488 54. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
 489 arXiv:1303.3997v1. (2013)
- 490 55. Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V. & Zdobnov, E.M.
 491 BUSCO: assessing genome assembly and annotation completeness with single-copy
- 493

492

494 Data Citations

- 495 1. NCBI Sequence Read Archive https://www.ncbi.nlm.nih.gov/bioproject/PRJNA512170
 496 (2018)
- 497 2. NCBI Assembly https://www.ncbi.nlm.nih.gov/assembly/GCA_004353265.1 (2019)

- 498 3. Girollet, N., Rubio, B. and Bert P.-F. De novo phased assembly of the *Vitis riparia* grape
- genome. figshare https://figshare.com/s/0a52d4408214e9f1e280 (2019)
- 500

501 Conflicts of interests

- 502 The authors declare that they have no competing interests.
- 503

504 Acknowledgements

- 505 We thank Dr. Nathalie Ollat and Dr. Gregory Gambetta for critical review of the manuscript.
- 506 We thank the INRA GeT-PlaGE platform (http://get.genotoul.fr) for sequencing the genome,
- 507 the INRA Genotoul bioinformatics platform (http://genotoul.toulouse.inra.fr) for providing
- 508 computational resources and support.
- 509 We thank Dr. Timothée Flutre and Amandine Launay for their help in coordinating the 510 FruitSelGen project and in acquiring GBS data.
- 511 This work was supported by a grant overseen by the French National Research Agency 512 (ANR) as part of the ANR-09-GENM-024 "Vitsec" program and by the University of
- 513 Bordeaux.

514

515 Author contributions

Author P.-F. B. conceived the project. N. G. and P.-F. B. assembled the genomes, performed the genome annotation and downstream analyses. B. R. performed GBS analysis and genetic mapping. P.-F. B. wrote the paper. All authors read, edited and approved the final manuscript.

520

522 Table 1: Data count and library informations for Vitis riparia genome sequencing

23								
		Sequencing	Insert	Read	Number of	Number	Sequence	Application
		platform	size	length	sequences	of bases	depth	
			(bp)	(bp)	(million)	(billion)		
		PacBio	NA	7054	8.3	59	118X	Genome
								assembly
		10X Chromium	400	2 x 150	350	52	107X	Genome
								scaffolding
								and phasing
		Illumina	400	2 x 100	331	33	66X	Genome
			(pair					survey and
			end)					genomic
			6,000	2 x 100	200	20	40X	base
			(mate					correction
			pair)					
	Total					164	331X	
24								
25								
26								
27								
28								
29								
30								
31								
32								
33								

Table 2: Summary of the V. riparia genome assembly and comparison with with V. vinifera

535 varieties.

536

	Vitis vinifera			Vitis riparia
	PN40024	Cabernet	Chardonnay	Riparia Gloire
		sauvignon		de Montpellier
Technology	Sanger	PacBio	PacBio	PacBio / 10X
				Chromium
Genome	12X	140X	115X	225X
coverage				
Contig length	NA	591 p-tigs	490 p-tigs	530 p-tigs
(Mb)		368 h-tgs	378 h-tgs	317 h-tgs
Number of	14,665	718 p-tigs	854 p-tigs	1,964 p-tigs
contigs		2,037 h-tgs	1,883 h-tgs	3,344 h-tgs
Number of	2,065	NA	NA	870
scaffolds				
N50 (kb)	103	2,170	935	964
Total length (Mb)	486	591	490	500
Number of coding	42,414	36,687	29,675	37,207
genes	(Cost.v3)			
BUSCO	C:95.8%	C:94.0%	C:95.0%	C:95.4%
	F:1.5%	F:2.0%	F:1.6%	F:1.1%
	M:2.7%	M:4.0%	M:3.4%	M:3.5%

537

539 Table 3: Repeated elements present in the Vitis riparia genome.

		number of	length occupied	percentage of
		elements	(base pairs)	sequence
Retroelements		95,441	125,292,108	25.05%
SINEs		0	0	0.00%
	Penelope	0	0	0.00%
LINEs:		17,557	18,048,495	3.61%
	CRE/SLACS	0	0	0.00%
	L2/CR1/Rex	0	0	0.00%
	R1/LOA/Jockey	0	0	0.00%
	R2/R4/NeSL	0	0	0.00%
	RTE/Bov-B	0	0	0.00%
	L1/CIN4	17,557	180,48,495	3.61%
LTR elemer	its:	77,884	107,243,613	21.44%
	BEL/Pao	0	0	0.00%
	Ty1/Copia	27,636	41,679,808	8.33%
	Gypsy/DIRS1	48,176	63,337,662	12.66%
	Retroviral	0	0	0.00%
DNA transpos	ons	80,801	30,107,834	6.02%
	hobo-Activator	14,757	6,431,569	1.29%
	Tc1-IS630-Pogo	0	0	0.00%
	En-Spm	436	1,086,471	0.22%
	MuDR-IS905	0	0	0.00%
	PiggyBac	0	0	0.00%
	Tourist/Harbinger	34,070	82,90 ,818	1.66%
	Other (Mirage, P-element,	0	0	0.00%
Rolling-circles	Rolling-circles		0	0.00%
Unclassified	Unclassified		67,393	0.01%
Total interspe	ersed repeats		155,467,335	31.09%
Small RNA		202	43,383	0.01%
Satellites			1,975,618	0.40%
Simple repeats		222,322	9,483,835	1.90%
Low complexi	ty	54,205	2,848,156	0.57%

Fig. 1: Comparison of *Vitis riparia* hybrid scaffolds with the reference PN40024 assembly. Hybrid scaffolds (Y-axis) were aligned to all 19 PN40024 chromosomes (X-axis) using D-GENIES and alignments were subsequently filtered for 1-on-1 alignments and rearrangements with a 20 Kbps length cutoff.





