

1 **An improved reference of the grapevine genome reasserts the origin of the**  
2 **PN40024 highly-homozygous genotype**

3

4 Authors:

5 Amandine Velt<sup>1</sup>°, Bianca Frommer<sup>2</sup>°, Sophie Blanc<sup>1</sup>, Daniela Holtgräwe<sup>2</sup>, Éric Duchêne<sup>1</sup>,  
6 Vincent Dumas<sup>1</sup>, Jérôme Grimplet<sup>3</sup>, Philippe Hugueney<sup>1</sup>, Catherine Kim<sup>4,5</sup>, Marie Lahaye<sup>1</sup>,  
7 José Tomás Matus<sup>6</sup>, David Navarro-Payá<sup>6</sup>, Luis Orduña<sup>6</sup>, Marcela K. Tello-Ruiz<sup>4</sup>, Nicola  
8 Vitulo<sup>7</sup>, Doreen Ware<sup>4,8</sup> and Camille Rustenholz<sup>1</sup>#  
9 ° equally contributed

10 #Corresponding author: Camille Rustenholz, [camille.rustenholz@inrae.fr](mailto:camille.rustenholz@inrae.fr) (C.R.)

11 Phone: 0033 3 89 22 49 81

12 Fax: 0033 3 89 22 49 33

13 Affiliations:

14 <sup>1</sup>SVQV, INRAE - University of Strasbourg, 68000 Colmar, France; [amandine.velt@inrae.fr](mailto:amandine.velt@inrae.fr)  
15 (A.V.); [sophie.blanc@inrae.fr](mailto:sophie.blanc@inrae.fr) (S.B.); [eric.duchene@inrae.fr](mailto:eric.duchene@inrae.fr) (É.D.); [vincent.dumas@inrae.fr](mailto:vincent.dumas@inrae.fr)  
16 (V.D.); [philippe.hugueney@inrae.fr](mailto:philippe.hugueney@inrae.fr) (P.H.); [marie.lahaye@inrae.fr](mailto:marie.lahaye@inrae.fr) (M.L.);  
17 [camille.rustenholz@inrae.fr](mailto:camille.rustenholz@inrae.fr) (C.R.)

18 <sup>2</sup>Genetics and Genomics of Plants, CeBiTec & Faculty of Biology, Bielefeld University,  
19 33615 Bielefeld, Germany; [daniela.holtgraewe@cebitec.uni-bielefeld.de](mailto:daniela.holtgraewe@cebitec.uni-bielefeld.de) (D.H.);  
20 [frommer@cebitec.uni-bielefeld.de](mailto:frommer@cebitec.uni-bielefeld.de) (B.F.)

21 <sup>3</sup>Unidad de Hortofruticultura, Centro de Investigación y Tecnología Agroalimentaria de  
22 Aragón (CITA), 50059 Zaragoza, Spain; [jgrimplet@cita-aragon.es](mailto:jgrimplet@cita-aragon.es) (J.G.)

23 <sup>4</sup>Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, 11724, USA; [ckim@cschl.edu](mailto:ckim@cschl.edu)  
24 (C.K.); [telloruiz@cschl.edu](mailto:telloruiz@cschl.edu) (M.K.T.-R.); [ware@cschl.edu](mailto:ware@cschl.edu) (D.W.)

25 <sup>5</sup>Haverford College, 370 Lancaster Avenue, Haverford, PA 19041, USA; [ckim@cschl.edu](mailto:ckim@cschl.edu)  
26 (C.K.)

27

1 <sup>6</sup> Institute for Integrative Systems Biology (I2SysBio), Universitat de València-CSIC,  
2 Paterna, 46908, Valencia, Spain; [tomas.matus@uv.es](mailto:tomas.matus@uv.es) (J.T.M.);  
3 [david.navarro.paya@gmail.com](mailto:david.navarro.paya@gmail.com) (D.N-P.); [luis.orduna.rubio@gmail.com](mailto:luis.orduna.rubio@gmail.com) (L.O.)

4 <sup>7</sup> Dipartimento di Biotecnologie, Università degli Studi di Verona, 37134, Verona, Italy;  
5 [nicola.vitulo@univr.it](mailto:nicola.vitulo@univr.it) (N.V.)

6 <sup>8</sup> USDA ARS NEA Robert W. Holley Center for Agriculture and Health, Agricultural  
7 Research Service, Ithaca, NY 14853, USA; [ware@csihl.edu](mailto:ware@csihl.edu) (D.W.)

8  
9 ORCID:

10 Amandine Velt: 0000-0003-2368-839X

11 Bianca Frommer: 0000-0002-5792-0102

12 Sophie Blanc: 0000-0003-2501-548X

13 Daniela Holtgräwe: 0000-0002-1062-4576

14 Éric Duchêne: 0000-0003-2712-1892

15 Vincent Dumas: 0000-0001-6089-3048

16 Jérôme Grimplet: 0000-0002-3265-4012

17 Philippe Hugueney: 0000-0002-1641-9274

18 José Tomás Matus: 0000-0002-9196-1813

19 David Navarro-Payá: 0000-0003-1623-3039

20 Luis Orduña: 0000-0003-2756-4028

21 Marcela K. Tello-Ruiz: 0000-0002-7499-5368

22 Nicola Vitulo: 0000-0002-9571-0747

23 Doreen Ware: 0000-0002-8125-3821

24 Camille Rustenholz: 0000-0001-5355-3408

25

26 Running Head: PN40024.v4, an improved grapevine reference genome

27 Keywords: *Vitis vinifera*, genotype PN40024, reference genome, long reads, improved  
28 annotation

## 1 Abstract

2 The genome sequence of the diploid and highly homozygous *V. vinifera* genotype PN40024  
3 serves as the reference for many grapevine studies. Despite several improvements to the  
4 PN40024 genome assembly, its current version PN12X.v2 is quite fragmented and only  
5 represents the haploid state of the genome with mixed haplotypes. In fact, being nearly  
6 homozygous, this genome contains several heterozygous regions that are yet to be resolved.  
7 Taking the opportunity of improvements that long-read sequencing technologies offer to fully  
8 discriminate haplotype sequences, an improved version of the reference, called PN40024.v4,  
9 was generated. Through incorporating long genomic sequencing reads to the assembly, the  
10 continuity of the 12X.v2 scaffolds was highly increased with a total number decreasing from  
11 2,059 to 640 and a reduction in N bases of 88%. Additionally, the full alternative haplotype  
12 sequence was built for the first time, the chromosome anchoring was improved and the number  
13 of unplaced scaffolds was reduced by half. To obtain a high-quality gene annotation that  
14 outperforms previous versions, a liftover approach was complemented with an optimized  
15 annotation workflow for *Vitis*. Integration of the gene reference catalogue and its manual  
16 curation have also assisted in improving the annotation, while defining the most reliable  
17 estimation of 35,230 genes to date. Finally, we demonstrated that PN40024 resulted from nine  
18 selfings of cv. ‘Helfensteiner’ (cross of cv. ‘Pinot noir’ and ‘Schiava grossa’) instead of a single  
19 ‘Pinot noir’. These advances will help maintain the PN40024 genome as a gold-standard  
20 reference, also contributing towards the eventual elaboration of the grapevine pangenome.

## 21 Introduction

22 Cultivated grapevine (*Vitis vinifera* ssp. *vinifera*) was the fourth plant whose genome was  
23 sequenced and assembled (Jaillon *et al.* 2007). Because of the grapevine’s high level of  
24 heterozygosity (one Single Nucleotide Polymorphism (SNP) per 100 bp and one Indel per  
25 450 bp, (Velasco *et al.* 2007)), the genotype selected for sequencing was PN40024, whose  
26 ~475 Mb genome (Lodhi and Reisch 1995) is nearly homozygous (estimated at ~93%).  
27 PN40024 was indeed generated through nine rounds of selfing and supposedly originated from  
28 ‘Pinot noir’, hence its identification as ‘PN’. This unique genome characteristic allowed a high-  
29 quality whole-genome shotgun assembly based on 8X coverage Sanger reads (Jaillon *et al.*  
30 2007). In 2009, a 4X coverage was added, which improved the overall coverage of the genome

1 (from 68.9% for the 8X version to 91.2% for the 12X.v0)  
2 (<http://urgi.versailles.inra.fr/Species/Vitis/Data-Sequences/Genome-sequences>; FN597015-  
3 FN597047 at EMBL, release 102; File S1 Figure S1). In 2017, a third assembly version, named  
4 12X.v2, was published as the result of a large anchoring effort using six dense parental genetic  
5 maps (Canaguier *et al.* 2017). Despite these advances, no additional sequencing efforts have  
6 been made and although it's of very high quality, the 12X.v0 Sanger contigs are numerous  
7 (14,642), the 12X.v2 scaffolds are composed of large N gaps (3.1% of the cumulative scaffold  
8 size) and the 19 pseudomolecules are quite fragmented (19.3 scaffolds on average per  
9 pseudomolecule).

10 In recent years, the advent of third generation sequencing technologies, especially those from  
11 the Pacific Biosciences (PacBio) platform, have allowed the assembly of grapevine diploid  
12 genomes with a higher level of contiguity compared to the 12X.v2 version of the PN40024  
13 genome (for example, cv. 'Cabernet Sauvignon' genome assembly (Massonnet *et al.* 2020)).

14 Along with the versions of each genome assembly, several versions of gene annotations were  
15 made available (File S1 Figure S1). The first version of the grapevine genome assembly, 8X,  
16 was published along with the prediction of 30,434 gene models based on the GAZE software  
17 (Howe *et al.* 2002; Jaillon *et al.* 2007). For the 12X.v0, three different versions of gene  
18 predictions were made available: the v0 version (26,346 gene models), based on the GAZE  
19 software (Howe *et al.* 2002), the CRIBIv1 version (29,971 gene models), based on the JIGSAW  
20 software (Allen and Salzberg 2005), and the CRIBIv2 version (31,845 gene models), with an  
21 effort made on the discovery of splicing variants (Vitulo *et al.* 2014). For the 12X.v2, the  
22 International Grapevine Genome Program (IGGP) led the initiative of merging annotations  
23 from NCBI Refseq, CRIBIv1 and VCost, which was based on the Eugene software (Sallet *et*  
24 *al.* 2019) and was generated in the frame of the COST Action FA1106. This version, called  
25 VCost.v3, resulted in an exhaustive view of the PN40024 grapevine gene content with its  
26 42,413 gene models (Canaguier *et al.* 2017). However, after several years as the reference  
27 annotation by the grapevine scientific community, it appeared that the great increase in number  
28 of gene models for VCost.v3 compared to all the previous annotation versions was caused by  
29 many small and fragmented predictions that were probably erroneous.

30 By combining the top-quality Sanger contigs from the 12X version and long reads generated  
31 here by Single-Molecule Real-Time (SMRT) sequencing (PacBio), we provide an improved  
32 version of the PN40024 genome sequence assembly, referred to as PN40024.v4. Along with

1 this new assembly, we also provide a new version of the gene annotation, PN40024.v4.2, based  
2 on a newly-developed annotation workflow, RNA-Seq datasets and an exhaustive manual  
3 curation of a set of catalogued genes of functional interest to the community. Finally, we  
4 demonstrate that PN40024 originates from selfings of the ‘Helfensteiner’ cultivar instead of  
5 ‘Pinot noir’.

6

## 7 **Methods & Materials**

### 8 **Plant material, DNA extractions and sequencing**

9 DNA extractions of young leaves of cv. ‘Pinot noir’ clone 162 (ID code FRA038-193.Col.162),  
10 cv. ‘Schiava grossa’ (synonymous ‘Trollinger’, ID code FRA038-2525.Col.1) and cv.  
11 ‘Helfensteiner’ (ID code FRA038-2744.Col.1) were performed as described by Merdinoglu *et al.*  
12 *(Merdinoglu et al. 2005)*. Illumina DNA PCR-Free Prep kit was used to prepare the  
13 resequencing libraries according to provider procedures. Paired-end Illumina HiSeq 4000  
14 sequencing at about 15x coverage was performed for ‘Pinot noir’ and ‘Schiava grossa’,  
15 respectively. Paired-end Illumina NovaSeq 6000 sequencing at about 15x coverage was  
16 performed for ‘Helfensteiner’.

17 One gram of young leaves (1cm<sup>2</sup>) of PN40024 (ID code FRA038-40024.Col.1) was collected  
18 and DNA was extracted using QIAGEN Genomic-tips 100/G kit. SMRT sequencing on a  
19 Sequel I machine (3 SMRTCells; PacBio) and dedicated library preparation were performed  
20 according to provider procedures.

21 Genotyping-by-sequencing (GBS) was performed on the population ‘Riesling’ x  
22 ‘Gewurztraminer’ (exhaustively described by Duchêne *et al.* (Duchêne *et al.* 2020) using the  
23 procedure described by Girollet *et al.* (Girollet *et al.* 2019).

24 All data generated in the frame of this study were submitted under the ENA Study Accession  
25 PRJEB45423.

### 26 **Genome assembly**

27 Raw SMRT reads ([ERR7997743](https://www.ncbi.nlm.nih.gov/sra/ERR7997743)) were self-corrected using CANU (v.1.6) (Koren *et al.* 2017),  
28 followed by a correction with PN40024 Illumina reads ([SRR8835144](https://www.ncbi.nlm.nih.gov/sra/SRR8835144)) using LORDEC

1 (v.0.5.3) (Salmela and Rivals 2014). The corrected reads were mapped on PN12X scaffolds  
2 ([https://urgi.versailles.inra.fr/download/vitis/VV\\_12X\\_embl\\_102\\_Scaffolds.fsa.zip](https://urgi.versailles.inra.fr/download/vitis/VV_12X_embl_102_Scaffolds.fsa.zip)) using  
3 minimap2 (v2.17-r954-dirty) (Li 2018). A total of 163,446 reads (15%) were aligned on less  
4 than 80% of their length and/or with less than 80% identity and were thus considered as missing  
5 from PN12X scaffolds. These unmapped reads were assembled using Flye (v2.4-gc9db046)  
6 (Kolmogorov *et al.* 2019). We aligned these new contigs on the Uniprot *Arabidopsis* database  
7 (release 2019\_01) using blastx (Altschul *et al.* 1990). Contigs longer than 5 kb and having  
8 hit(s) with *Arabidopsis* proteins with >60% identity and >60% length coverage, were selected  
9 for the next step. The fasta files of these new contigs and the PN12X scaffolds were  
10 concatenated to generate the new assembly. Firstly, the repeats were masked using Red  
11 (v05/22/2015) (Girgis 2015). Then, Haplomerger2 (v20180603) (Huang *et al.* 2017) was used  
12 following three steps according to developer procedures: i) break the misjoins and output the  
13 new diploid assembly; ii) separate/merge two haplotypes and output haploid assemblies (REF  
14 and ALT); and iii) remove tandem errors from haploid assemblies. Some scaffolds / contigs  
15 were deleted by Haplomerger2 during the assembly process but sequences longer than 10 kb  
16 were retrieved and added to the REF scaffolds. The two haploid assemblies (REF/ALT) were  
17 then scaffolded with the OPERA-LG tool (v2.0.6) (Gao *et al.* 2016), which uses both, corrected  
18 SMRT reads and Illumina reads. A first gap-filling step (two rounds) was carried out with  
19 Illumina reads using GapCloser (v1.12) (Luo *et al.* 2012) and a second gap-filling step (three  
20 rounds) was carried out with corrected SMRT reads using LR\_Gapcloser (v1.0) (Xu *et al.*  
21 2019). A final polishing step was performed with the Illumina reads using PILON (v1.23-1-  
22 g41e0b8e) (Walker *et al.* 2014) (Figure 1A and 1B).

23 The anchoring of the new haploid scaffolds was performed using the six genetic maps used for  
24 the same purpose by Canaguier *et al.* (Canaguier *et al.* 2017) and two new genetic maps from  
25 cv. ‘Riesling’ and cv. ‘Gewurztraminer’ derived from GBS. To transfer the markers from  
26 Canaguier *et al.* (Canaguier *et al.* 2017) from PN12X.v2 to the scaffolds of PN40024.v4,  
27 BLAST (v2.2.28) (Altschul *et al.* 1990) or ipress (ipress from exonerate v2.2.0) (Slater and  
28 Birney 2005) was used to align the markers and find the position of each on the scaffolds of  
29 PN40024.v4 REF and ALT. A total of 2,333 markers for REF and 2,326 markers for ALT were  
30 used from these six maps to anchor the scaffolds. For the two new genetic maps from ‘Riesling’  
31 and ‘Gewurztraminer’, 5,884 (‘Riesling’) and 5,840 (‘Gewurztraminer’) SNP markers were  
32 available for REF and 5,866 (‘Riesling’) and 5,832 (‘Gewurztraminer’) for ALT. The SNP  
33 markers were derived from GBS data (ERR8657388 to ERR8657647) and were analyzed with

1 Fast-GBS (Torkamaneh *et al.* 2017) with modifications to allow paired-end read analysis  
2 (<https://forgemia.inra.fr/sophie.blanc/gbs>). The two genetic maps were built using R ASMap  
3 package with the “kosambi” parameter (Taylor and Butler 2017). A first run of Allmaps  
4 (v0.9.13) (Tang *et al.* 2015) was performed with the “merge” command to merge all genetic  
5 maps and then “split”, “gaps”, “refine” and “build” commands to create breakpoints (58 for  
6 REF scaffolds and 47 for ALT scaffolds), with default parameters. Subsequently, all maps were  
7 recreated for new scaffolds and then orientation and anchoring of new haploid scaffolds on the  
8 19 pseudochromosomes were performed using Allmaps with the “merge” command to merge  
9 all maps and “path” command to anchor, with default parameters (Figure 1C).

10

### 11 **Quality assessment of the PN40024.v4 genome sequence assembly**

12 A quality analysis of the genome assembly was done with Merqury v1.3 (Rhie *et al.* 2020).  
13 Since PN40024 is a ‘Helfensteiner’ selfing (demonstrate below) and since ‘Helfensteiner’  
14 originated from a cross between ‘Pinot noir precoce’ and ‘Schiava grossa’, ‘Schiava grossa’  
15 was used as the maternal parent. The run was carried out on the scaffolds using genomic paired-  
16 end short reads of PN40024 as the child data (SRR8835144), short reads of ‘Pinot noir’ as the  
17 paternal data (ERR8014965) and short reads of cv. ‘Schiava grossa’ as the maternal data  
18 (ERR8014964). A k-mer database was built for the three read datasets with k=19, the Merqury  
19 hap-mer databases were computed and the PN40024.v4 genome assembly was evaluated using  
20 ‘num\_switch 100’ and ‘short\_range 20,000’. For comparison reasons, the Merqury quality  
21 analysis was carried out on PN12X.v2 using the same k-mer databases.

22 The “*Flowering locus T*” (*FT*) and the “*Adenine phosphoribosyltransferase 3*” (*APRT3*) genes  
23 are absent and truncated in PN12X.v2, respectively. To check whether these genes could be  
24 retrieved in the new genome assembly, cDNA sequences of *FT* (NM\_001280978.1) and  
25 *APRT3* (GSVIVT00007310001, PN8X version) were used to perform blastn (Altschul *et al.*  
26 1990) against PN8X, PN12X, PN12X.v2 and PN40024.v4 genome assemblies. High Scoring  
27 Pairs (HSPs) were then accumulated for each analysis and the mean percentage identity, query  
28 overlap, hit query start and end were calculated.

29

30 PN40024 (SRR8835144), and the cultivars ‘Silvaner Gruen’ (SRR5891620), ‘Cabernet Franc’  
31 (SRR5891774), ‘Cabernet Sauvignon’ (SRR5891776), ‘Chardonnay’ (SRR5891778), ‘Muscat  
32 Hamburg’ (SRR5891787), ‘Semillon’ (SRR5891866), ‘Pinot Noir’ (SRR5891886), ‘Merlot’



1 (SRR5891890), ‘Sauvignon Blanc’ (SRR5891893), ‘Muscat of Alexandria’ (SRR5891985)  
2 and ‘Riesling’ (SRR5891989) genomic paired-end resequencing datasets were aligned against  
3 PN40024.v4 REF, PN12X.v2 and ‘Cabernet Sauvignon’ haplotype 1 (Massonnet *et al.* 2020)  
4 pseudomolecule assemblies (without chrUn or unplaced contigs / scaffolds) using bwa-mem2  
5 (v2.0) (Vasimuddin *et al.* 2019) “mem” command with default parameters. Samtools’ (v1.9)  
6 (Li *et al.* 2009) “flagstat” command was used with default parameters to compute alignment  
7 statistics.

8 PN12X scaffolds were mapped against PN40024.v4 REF pseudomolecules using NUCmer  
9 (MUMmer v3.1) (Kurtz *et al.* 2004) with “-maxmatch -l 100 -c 500” parameters. The output  
10 file was filtered using MUMmer show-coords command with “-l -g -I 99.5” parameters. The  
11 resulting file was formatted into BED format and merged with the bed file corresponding to N  
12 gap regions in the PN40024.v4 assembly. Pseudomolecule regions over 100 bp that did not  
13 correspond to either PN12X scaffolds or N gap regions were identified as ‘newly assembled’  
14 PacBio long read-based regions.

15 The identification of variants between PN40024 paired-end Illumina resequencing  
16 (SRR8835144) and PN40024.v4 REF and ALT pseudomolecules was performed as described  
17 in section “Origin of PN40024”. The homozygous calls “1/1” were considered as assembly  
18 errors. The densities of the heterozygous calls “0/1” along the REF and ALT pseudomolecules  
19 were used to define seven heterozygous regions of the PN40024 genome.

## 20 **Origin of PN40024**

21 PN40024 ([SRR8835144](https://sra.ebi.ac.uk/submit/SRR8835144)), ‘Pinot noir’ ([ERR8014965](https://sra.ebi.ac.uk/submit/ERR8014965)), ‘Schiava grossa’ ([ERR8014964](https://sra.ebi.ac.uk/submit/ERR8014964)),  
22 ‘Helfensteiner’ ([ERR8014963](https://sra.ebi.ac.uk/submit/ERR8014963)) and ‘Araklinos’ ([SRR8835172](https://sra.ebi.ac.uk/submit/SRR8835172)) paired-end resequencing  
23 datasets were all analyzed using the same pipeline. Datasets were aligned against PN40024.v4  
24 REF assembly using bwa-mem2 (v2.0) (Vasimuddin *et al.* 2019) “mem” command with default  
25 parameters. Samtools (v1.9) (Li *et al.* 2009) “view” and “sort” commands with default  
26 parameters were used to convert and sort the output BAM files. GATK (v4.1.4.0) (McKenna  
27 *et al.* 2010) “MarkDuplicatesSpark”, “HaplotypeCaller” and “GenotypeGVCFs” commands  
28 with default parameters were used to generate variant files in VCF format. The GATK  
29 “VariantFiltration” command was used to filter out variants meeting at least one of the  
30 following criteria: QD < 8.0, QUAL < 100.0, FS > 60.0, SOR > 3.0, DP < 3, DP > 30, AD <  
31 2. The final variant files were obtained using GATK “SelectVariants” command with “--  
32 exclude-filtered --exclude-non-variants” parameters. The homozygous SNP calls “1/1” were



1 selected for each analyzed genotype. All SNPs corresponding to a homozygous call in  
2 PN40024 genotypes were excluded from the analysis as they represent assembly errors. The  
3 remaining homozygous SNPs were used to draw density plots on the PN40024.v4  
4 pseudomolecules. The regions that are rich in homozygous SNPs for a given genotype  
5 correspond to regions for which this genotype does not share a haplotype with PN40024.

6 The haplotypic blocks were defined after segmentation of homozygous SNP densities along  
7 the chromosomes using the R package changepoint (v2.2.2) (Killick and Eckley 2014) with  
8 command “cpt.mean” and the parameters method="PELT" and penalty="AIC". Some manual  
9 curation of the segments was performed to join directly adjacent segments of the same origin  
10 (‘Pinot noir’ or ‘Schiava grossa’). The size of the segments was used to calculate the proportion  
11 of ‘Pinot noir’, ‘Schiava grossa’ and common haplotypes.

## 12 **Gene prediction**

13 Before performing gene prediction, the PN40024.v4 genome assembly was repeat masked with  
14 RepeatMasker v4.1.2 (Smit *et al.* 2013) using crossmatch as search engine. Predictions with a  
15 SW-Score <1,000 were filtered out and predictions with a Smith-Waterman (SW)-Score  
16 between 1,000 and 2,000 were only kept if the reported percentage of substitutions were <20%.  
17 The PN40024.v4 genome assembly was softmasked with BEDTool (v2.26.0) (Quinlan and  
18 Hall 2010).

19  
20 To annotate the PN40024.v4 genome assembly, publicly available *V. vinifera* stranded (File  
21 S2 Table S1) and unstranded (File S2 Table S2) paired-end RNA-Seq datasets of different  
22 tissues and treatments were collected. RNA-Seq data were trimmed with Trimmomatic (v0.39)  
23 (Bolger *et al.* 2014). The annotation pipeline was first tested on the PN40024 12X.v0 genome  
24 assembly using VCost.v3 gene annotation as quality reference. The gene predictors SNAP  
25 (Korf 2004) and BRAKER2 (Hoff *et al.* 2016, 2019; Brůna *et al.* 2021) were trained and tested  
26 on the softmasked 12X.v0 genome assembly. The RNA-Seq data was mapped on 12X.v0 and  
27 on PN40024.v4 REF and ALT sequences with GMAP/GSNAP v2020-09-12 setting “-B 5 --  
28 novelsplicing 1” (Wu and Watanabe 2005). Primary mappings were extracted with SAMTools  
29 v1.9 (Li *et al.* 2009). Based on the primary mappings, stranded and unstranded reference-  
30 guided transcriptome assemblies were computed with PsiCLASS v1.0.1 using default  
31 parameters (Song *et al.* 2019).

1 Additionally, *A. thaliana* protein sequences (UniProt/SwissProt release 2020\_02),  
2 eudicotyledone protein sequences (UniProt/SwissProt release 2020\_02, OrthoDB10 v1),  
3 Viridiplantae and Vitales sequences (UniProt/SwissProt release 2020\_02) were aligned on  
4 12X.v0 and on PN40024.v4 REF and ALT with pBLAT v1.9 (Wang and Kong 2019), a parallel  
5 implementation of the original blat algorithm (Kent 2002). The genome regions on which the  
6 protein data mapped were extracted and the protein sequences were aligned to these regions  
7 with exonerate v2.4.0 (Slater and Birney 2005). Only the proteins that aligned on the reference  
8 genome with an identity of 25%, a similarity of 50% and with a sequence alignment coverage  
9 of at least 80%, were retained and included in the gene prediction.

10 The gene predictor GlimmerHMM v3.0.4 (Majoros *et al.* 2004) was trained on 12X.v0 and on  
11 PN40024.v4 REF and ALT using 7,500 (12X.v0) and 15,000 (PN40024.v4) random  
12 PsiCLASS transcripts of the 12X.v0 or PN40024.v4 REF or ALT stranded transcriptome  
13 assembly, respectively. The training was followed by gene prediction with GlimmerHMM with  
14 default settings.

15 Moreover, the gene predictor SNAP v2006-07-28 was trained on the 12X.v0 genome assembly.  
16 For this, the 12X.v0 genome assembly, the stranded transcriptome assembly, the Viridiplantae  
17 protein sequences and the eudicotyledone protein sequences were given to MAKER2 v3.01.03  
18 (Holt and Yandell 2011; Campbell *et al.* 2014) and initial data alignment with BLAST (ncbi-  
19 blast-2.10.1+) (Altschul *et al.* 1990; Camacho *et al.* 2009) and exonerate was performed  
20 followed by MAKER2 *ab initio* gene prediction. MAKER2 was run with  
21 “max\_dna\_len=300000” and “split\_hit=20000”. A SNAP hmm file was generated with the  
22 MAKER2 gff file and a second MAKER2 run was performed with enabled SNAP gene  
23 prediction and the SNAP hmm file as input. Hmm file generation and SNAP gene prediction  
24 with MAKER2 and the new hmm file was repeated. The hmm file generated with the 12X.v0  
25 assembly was used to run SNAP gene prediction on the PN40024.v4 REF and ALT genome  
26 sequences.

27 An AUGUSTUS species model was computed with BRAKER2 v2.1.5-master\_20200915 and  
28 the 12X.v0 genome assembly. BRAKER2 was run with enabled softmasking and in *etpmode*  
29 calling GeneMark-ETP+ v4.61 (Lomsadze *et al.* 2005, 2014; Brúna *et al.* 2020) for initial gene  
30 prediction followed by AUGUSTUS training and gene prediction (AUGUSTUS version  
31 master\_v3.3.3\_20200914) (Stanke *et al.* 2006, 2008). With BRAKER2, the programs  
32 DIAMOND v0.9.24.125 (Buchfink *et al.* 2015), SAMtools v1.9-180-gf9e1caf (Li *et al.* 2009),

1 SPALN version 2.3.3f (Gotoh 2008; Iwata and Gotoh 2012), ProtHint version 2.5.0 and  
2 BamTools v2.5.1 (Barnett *et al.* 2011) were called. The stranded RNA-Seq primary mappings,  
3 the eudicotyledon protein sequences (OrthoDB10 v1) and the Viridiplantae protein sequences  
4 were used as input. The gene prediction on PN40024.v4 REF and ALT was performed with  
5 BRAKER2 v2.1.5-master\_20210218, the generated AUGUSTUS species model and  
6 AUGUSTUS version master\_v3.4.0\_20210218. Again, the stranded RNA-Seq mappings and  
7 the same protein sequences were used as input. The BRAKER2 parameter settings were left  
8 the same as above.

9 The last *ab initio* gene prediction was done on the PN40024.v4 genome assembly with GeneID  
10 v1.4.5-master-20200916 and the publicly available *V. vinifera* parameter set using default  
11 settings. To add the VCost.v3 gene annotation to the set of predictions, an annotation liftover  
12 was performed with liftoff v1.5.1 (Shumate and Salzberg 2021) with default parameters onto  
13 the PN40024.v4 genome assembly.

14 To combine the predictions and evidence data into an overall gene model set, the  
15 GlimmerHMM, SNAP, BRAKER2 and GeneID *ab initio* gene prediction as well as the lifted  
16 VCost.v3 annotation, the stranded and unstranded transcriptome assemblies, the GFF file with  
17 the aligned protein data, the repeat annotation GFF file and the PN40024.v4 genome assembly  
18 was given to EvidenceModeler v1.1.1 (Haas *et al.* 2008). The used weights are listed in File  
19 S2 Table S3.

20 Subsequently, the raw gene models were quality filtered. Gene models only supported by *ab*  
21 *initio* predictors were kept if at least two gene prediction programs predicted them, if the start  
22 and stop codon was present and if the gene length was equal or larger than 300 bp. However,  
23 *ab initio* supported gene models not matching these constraints were kept if they had a database  
24 hit with the UniProt/SwissProt or NCBI non-redundant database. To obtain that, a blastp search  
25 of the protein sequences against the two databases was run, allowing hits with an e-value  $<1e^{-6}$ .  
26 <sup>6</sup>. Of the gene models only supported by evidence data or by VCost.v3 lifted annotation, those  
27 gene models with missing start and stop and a gene length  $<300$  bp were discarded.

28 The gene models generated by EvidenceModeler were finally processed by PASA (v2.4.1,  
29 default parameters) using the stranded transcriptome assembly as a reference to add UTR  
30 regions and to calculate alternatively spliced models. Genes with overlapping UTRs were

1 shortened. tRNAs were predicted with tRNAscan-SE-2.0 (Chan *et al.* 2021) on the  
2 PN40024.v4 genome assembly.

3 To retain gene naming of VCost.v3 gene models, a reciprocal best blast hit (RBH) search  
4 between protein sequences of PN40024.v4.1 gene models and protein sequences of VCost.v3  
5 gene models was carried out. For the RBH search, only the longest protein sequence per gene  
6 was used, the e-value was set to  $1e^{-4}$  and the query coverage and identity was set to 70%.  
7 Moreover, only RBHs with genes on the same pseudochromosome and showing collinearity  
8 with other genes were considered valid. Thus, genes with a valid RBH were named according  
9 to the VCost.v3 gene, novel genes received the prefix '04' at the start of the gene number and  
10 genes predicted for alternative heterozygous sequence regions received the suffix '\_alt' (File  
11 S2 Table S4).

12 The PN40024.v4.1 gene models were functionally annotated with Blast2GO (v1.5.1) (Conesa  
13 *et al.* 2005; Götz *et al.* 2008). For this, protein domains of the PN40024.v4.1 proteins were  
14 identified with InterProScan (v5.52-86.0) (Jones *et al.* 2014) with options "--goterms --  
15 pathways -dp" using the databases/tools CDD-3.18 (Lu *et al.* 2020), Coils-2.2.1 (Lupas *et al.*  
16 1991), Gene3D-4.3.0 (Sillitoe *et al.* 2019), Hamap-2020\_05 (Pedruzzi *et al.* 2013),  
17 MobiDBLite-2.0 (Necci *et al.* 2017), PANTHER-15.0 (Mi *et al.* 2021), Pfam-33.1 (Mistry *et al.*  
18 *et al.* 2021), PIRSF-3.10 (Wu *et al.* 2004), PIRSR-2021\_02, PRINTS-42.0 (Attwood *et al.* 2012),  
19 ProSitePatterns-2021\_01, ProSiteProfiles-2021\_01 (Sigrist *et al.* 2013), SFLD-4 (Akiva *et al.*  
20 2014), SMART-7.1 (Letunic and Bork 2018), SUPERFAMILY-1.75 (Gough *et al.* 2001;  
21 Wilson *et al.* 2009), TIGRFAM-15.0 (Haft *et al.* 2013). PN40024.v4.1 protein sequences were  
22 aligned with diamond "blastp" (v2.0.11) (Buchfink *et al.* 2015) to the NCBI nr database  
23 (nr.07\_07\_2021.fasta) with options "--sensitive --top 5 -e 1e-5 -f 5". InterProScan and diamond  
24 results were used as input for Blast2GO.

### 25 **Quality assessment of the PN40024.v4.1 gene annotation**

26 To estimate completeness of the PN40024.v4.1 gene model set, plant core genes were predicted  
27 with BUSCO v5.1.2 (Simão *et al.* 2015; Waterhouse *et al.* 2018) using database  
28 eudicots\_odb10.

29 Samples previously analyzed by Palumbo *et al.* (Palumbo *et al.* 2014) were used to perform  
30 differential gene expression analysis by using either PN12X.v2 assembly with VCost.v3  
31 annotations or PN40024.v4 assembly with PN40024.v4.1 annotations. We analyzed cv.

1 ‘Sangiovese’ (SRR1631822; SRR1631823; SRR1631824), cv. ‘Barbera’ (SRR1631834;  
2 SRR1631835; SRR1631836) and cv. ‘Refosco’ samples (SRR1631858; SRR1631859;  
3 SRR1631860) for the stage “Berries beginning to touch” (~EL35 according to Eichhorn and  
4 Lorenz phenological scale (Eichhorn and Lorenz 1977)). The RNA-Seq data were downloaded  
5 from the SRA with SRA Toolkit (v2.10.8) (SRA Toolkit Development Team;  
6 <https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=software>) and analyzed with an *in-*  
7 *house* pipeline using FASTQC (v0.11.5) (Andrews 2010), STAR (v2.5.3a) (Dobin *et al.* 2013),  
8 Samtools (v1.4.1) (Li *et al.* 2009), Bamtools (v2.4.0) (Barnett *et al.* 2011), featureCounts  
9 (v1.5.3) (Liao *et al.* 2014) and SARTools (v1.7.3) (Varet *et al.* 2016).

## 10 **Manual gene model curation**

11 For manual gene model curation, an Apollo Webserver v2.6.4  
12 (<https://github.com/GMOD/Apollo/blob/master/README.md>) (Dunn *et al.* 2019) was set up  
13 for the PN40024.v4 genome assembly and provided with different data tracks such as  
14 PN40024.v4.1 and previous gene annotations, RNA-Seq mappings and exonerate protein  
15 mappings (see section Material and Methods - Gene prediction). By these means, gene models  
16 were manually inspected and curated if needed or also new genes were added following  
17 dedicated guidelines offered to the community ([https://integrape.eu/resources/data-](https://integrape.eu/resources/data-management/)  
18 [management/](https://integrape.eu/resources/data-management/)). Using Apollo, the plant core genes classified as fragmented or missing by  
19 BUSCO were manually curated and adapted if necessary. In the frame of this study, we also  
20 began to manually curate genes present in the grape reference catalogue ((Navarro-Payá *et al.*  
21 2022); <https://grapedia.org/genes/>). A home-made python script was used to generate the  
22 PN40024.v4.2 version of gene annotations including those manually curated  
23 ([https://gitlab.com/MSVteam/pn40024-visualization-tools/-/tree/master/update\\_gff3\\_script](https://gitlab.com/MSVteam/pn40024-visualization-tools/-/tree/master/update_gff3_script)).

24

## 25 **Results and Discussion**

### 26 **Improved metrics for the genome assembly of PN40024**

27 A hybrid strategy was developed to assemble the genome of PN40024 genotype using 27X of  
28 SMRT long reads along with the PN12X scaffolds and 15X PN40024 Illumina paired-end  
29 resequencing data (Figure 1). This new assembly was named PN40024.v4. Six hundred and

1 forty scaffolds were produced with a N50 size of 6.5 Mb for a cumulative size of 474.5 Mb for  
2 the PN40024.v4 REF haplotype (Table 1). Compared to the former PN12X.v2, the number of  
3 scaffolds was reduced by a factor of three and the N50 was doubled. Moreover, the number of  
4 unknown bases, marked as N in the new scaffold sequences, represents 1.8 Mb and 0.4% of  
5 the assembly size versus 15.0 Mb and 3.1% for PN12X.v2 scaffolds. Thus, PN40024.v4 REF  
6 is more contiguous and has more informative sequences than PN12X.v2. Also, the PN40024.v4  
7 assembly size is closer to the grapevine genome size of 475 Mb, estimated using flow  
8 cytometry by Lodhi and colleagues (Lodhi and Reisch 1995). Phasing efforts on the partially  
9 heterozygous genotype resulted in the reconstruction of the second PN40024 haplotype  
10 (PN40024.v4 ALT) with 485 scaffolds and a total genome assembly size of ~463 Mb (File S2  
11 Table S5). Thus, the PN40024.v4 genome assembly now represents both haplotypes of the  
12 diploid PN40024 genome.

13 There are 7,640 newly assembled PacBio long read-based regions that were identified as  
14 missing from PN12X.v2 scaffolds. Their cumulative size is 24.1 Mb, *i.e.* 5.1% of the total  
15 PN40024.v4 genome assembly size (average = 3,152 bp; median = 558 bp; max = 32,650 bp).

16 A total of 2,333 markers were used from the six Canaguier's maps (Canaguier *et al.* 2017), in  
17 addition to 5,884 and 5,840 SNP markers from cv. 'Riesling' and cv. 'Gewurztraminer' GBS  
18 maps, respectively, to anchor these scaffolds. We were able to anchor 165 PN40024.v4 REF  
19 scaffolds to the 19 pseudochromosomes, for a cumulative size of ~462 Mb (97.4%) (Table 1).  
20 The 19 PN40024.v4 REF pseudomolecules are composed of 8.7 scaffolds on average (min =  
21 3; max = 26; median = 6) whereas 19.3 scaffolds on average composed the PN12X.v2  
22 pseudomolecules (min = 5; max = 82; median = 13). The remaining unplaced scaffolds were  
23 ordered according to their size to generate "chrUn" sequence representing 12.5 Mb (-47%  
24 compared to PN12X.v2 unplaced scaffolds). Thus, PN40024.v4 anchoring was improved as  
25 the pseudomolecules are less fragmented and as the size of chrUn has almost been halved.

26 At the chromosome scale, 10 pseudomolecules became shorter in PN40024.v4 compared to  
27 PN12X.v2 (average loss = ~448 kb; median = ~255 kb; min = 2,961 bp; max = 1,133,439 bp).  
28 Chromosome 6 showed the biggest reduction as the location of a large fragment has been  
29 correctly assigned to chromosome 9 (File S1 Figure S2). Nine pseudomolecules became larger  
30 (average gain = ~869 kb; median = ~582 kb; min = 15,118 bp; max = 2,045,414 bp), notably  
31 chromosome 9, 7 and 15, which gained 1.5 Mb, 1.9 Mb and 2.0 Mb, respectively.

1 By aligning PN40024 Illumina paired-end reads against PN40024.v4 genome assembly, we  
2 identified 101,778 heterozygous variations. Using their density along the chromosomes, we  
3 were able to identify seven well-defined heterozygous regions in PN40024.v4 genome  
4 assembly as it was the first time that a software dedicated to diploid assembly (Haplomerger2)  
5 was used to assemble the PN40024 genome. These regions were located on chromosomes 2,  
6 3, 4, 7, 10, 11 and 15 with the two largest regions being on chromosome 7 and 10 (11.4 Mb  
7 and 5.5 Mb, respectively) (Figure 3). Their overall cumulative size of 20.6 Mb represents 4.3%  
8 of PN40024.v4, which is less than the residual heterozygosity size of 7%, estimated by Jaillon  
9 and colleagues based on genetic markers (Jaillon *et al.* 2007). Using the same procedure, we  
10 identified six heterozygous regions in PN12X.v2 assembly on the same chromosomes as  
11 PN40024.v4 except for the one on chromosome 15. Their overall cumulative size of 16.6 Mb  
12 represents 3.4% of PN12X.v2 and 4 Mb less than the heterozygous regions anchored on the  
13 PN40024.v4 chromosomes. These sequences were badly resolved and mostly located in the  
14 unanchored fraction of PN12X.v2 assembly (File S1 Figure S2). Thus, we conclude that  
15 PN40024.v4 is a better diploid assembly compared to PN12X.v2.

#### 16 **Quality of the PN40024.v4 genome assembly**

17 The BUSCO analysis performed on the PN40024.v4 genome assembly confirmed that the gene  
18 space was more complete with 98.1% of the 2,326 total searched Eudicots BUSCO genes being  
19 complete, compared to PN12X.v2 with 97.6% (Figure 6). The *FT* (*Flowering locus T*) gene is  
20 conserved among all flowering plants as it promotes transition from vegetative growth to  
21 flowering. However, its sequence could only be found on an unanchored scaffold in the PN8X  
22 version and was totally missing in PN12X.v0 and PN12X.v2. It is now present on chromosome  
23 7 of the PN40024.v4 assembly and also on its allelic region, chromosome 7\_ALT sequence.  
24 Similarly, the *APRT3* gene, located in the sex determination locus of grapevine, was present  
25 on chromosome 2 in the PN8X version and was truncated in PN12X.v0 and PN12X.v2. It is  
26 now fully retrieved on chromosome 2 of PN40024.v4 assembly and on its allelic region,  
27 chromosome 2\_ALT sequence. These two examples, along with the BUSCO analysis, show  
28 that the PN40024.v4 assembly is more complete, especially in the residual heterozygous  
29 regions that are now more accurately exposed.

30 The alignment metrics of PN40024 genomic Illumina paired-end reads have always been better  
31 against PN40024.v4 compared to PN12X.v2, either for overall percentage of mapped reads  
32 (97.58% versus 96.58%) or for properly mapped pairs of reads (85.81% versus 82.82%)



1 (Figure 2). This confirms that the PN40024.v4 assembly is more complete and with a more  
2 accurate structure than PN12X.v2. Moreover, we compared alignments of 11 genomic Illumina  
3 paired-end read datasets from various cultivars against PN40024.v4 and PN12X.v2 assemblies,  
4 but also against ‘Cabernet Sauvignon’ (Massonnet *et al.* 2020) haplotype 1, whose assembly  
5 metrics and technology were similar to PN40024.v4. Again, PN40024.v4 performs best for  
6 each dataset, even when ‘Cabernet Sauvignon’ was aligned against its own assembly  
7 (Figure 2). These results confirm that PN40024.v4 shows a quality suitable to become the new  
8 grapevine reference genome assembly, as it performs well with aligning genomic reads of  
9 various *V. vinifera* cultivars.

10  
11 The error rate at nucleotide level was assessed by calling homozygous variations between  
12 PN40024 genomic Illumina paired-end reads aligned against the PN40024.v4 genome  
13 assembly. We identified 28.7 errors / Mb compared to 8.4 errors / Mb in the PN12X.v2 genome  
14 assembly. However, they are unevenly distributed along the chromosomes and they mostly co-  
15 localize with the newly assembled long read-based regions and the seven heterozygous regions  
16 (Figure 3). A higher density of errors was also detected in the heterozygous regions of the  
17 PN12X.v2 genome assembly (File S1 Figure S3). We detected 284.4 errors / Mb in  
18 PN40024.v4 heterozygous regions of and 83.1 errors / Mb in PN12X.v2 heterozygous regions,  
19 which is, respectively, about 10 times denser than their average error rate. Thus, the overall  
20 increase of error rate in the PN40024.v4 assembly is mostly due to the use of SMRT long reads  
21 to improve the completeness of the reference genome assembly.

22  
23 Using Merqury, the base level quality value (QV) of the PN40024.v4 genome assembly was  
24 estimated to be 36.02, which is slightly worse than QV of 37.43 of the PN12X.v2 genome  
25 assembly (Table 2). This result confirms that additional SMRT sequences are not as accurate  
26 as Sanger-based sequences and they slightly decrease overall accuracy of the assembly. Also,  
27 the error rate of the PN40024.v4 genome assembly was increased by 0,00006964% compared  
28 to PN12X.v2, but still represents an accuracy of 99.999749801%, a metric associated with  
29 high-quality genome assemblies.

30 Nevertheless, the k-mer completeness was raised from 96.79% to 96.96% for the PN40024.v4  
31 assembly. Based on k-mer profiles of PN40024 and its parents (see “The origin of the PN40024  
32 genotype” section for details), Merqury computed the inheritance spectrum (File S1 Figure S4)  
33 showing a low portion of read-only missing k-mers that are unique for the child read set (paired-

1 end short reads of PN40024). The few missing sequences are probably due to sequencing  
2 errors, k-mers of novel variations or contamination from microbiome in PN40024 short reads,  
3 indicating an almost fully-complete PN40024.v4 genome sequence assembly. Also, as the  
4 spectrum shows a single 2-copy peak around 12x and that no 1-copy peak was observed at half  
5 the size, the k-mer analysis supports the assumptions of an almost homozygous grapevine  
6 genotype.

## 7 **The origin of the PN40024 genotype**

8 So far, the PN40024 genotype was supposed to be originally derived from cv. ‘Pinot noir’  
9 (Jaillon *et al.* 2007). However, we found 1,415,200 homozygous variants between ‘Pinot noir’  
10 and PN40024.v4 (versus 17,696 homozygous variants of PN40024 against its own assembly),  
11 meaning that ‘Pinot noir’ haplotypes were completely missing at these locations. These  
12 homozygous ‘Pinot noir’ variants were unevenly distributed along the chromosomes and  
13 formed blocks (Figure 4). We identified that the haplotypes of unknown origins could be  
14 assigned to ‘Schiava grossa’ (synonyms: ‘Trollinger’ and ‘Frankenthal’) as already suspected  
15 by Jaillon and colleagues (Jaillon *et al.* 2007). There were 953,735 homozygous variants found  
16 between cv. ‘Schiava grossa’ and PN40024.v4 and the formed haplotype blocks were highly  
17 complementary to ‘Pinot noir’ haplotype blocks (Figure 4). As a negative control, the same  
18 analysis was performed with cv. ‘Araklinos’ and 2,273,888 homozygous variants were  
19 identified, evenly distributed along the chromosomes (File S1 Figure S5).

20

21 Using Merqury, only a small portion of hap-mer specific k-mers (parental specific k-mers of  
22 the assembled F1) were found in the PN40024.v4 genome assembly (File S1 Figure S4). With  
23 the use of read data from both parents and child, Merqury was able to compute haplotype blocks  
24 by using the parental specific k-mers as anchors. A total of 1,454 haplotype blocks were  
25 computed for PN40024.v4 sequences with additional 289 haplotype blocks for alternative  
26 heterozygous sequence regions and 2,575 haplotype blocks for the 12X.v2 genome assembly  
27 (Table 3). The N50 was measured to 2.05 Mb (REF), 0.25 Mb (ALT) and 1.76 Mb  
28 (PN12X.v2). Compared to the PN12X.v2 genome assembly, PN40024.v4 presented less  
29 haplotype blocks, but comprised almost all bases showing a higher N50 value, *i.e.* its haplotype  
30 blocks are more contiguous.

1 A greater amount of paternal ('Schiava grossa') than maternal ('Pinot Noir') specific k-mers  
2 were identified. After identifying the origin of each haplotype block using segmentation, it is  
3 estimated that 41% of the genome harbours a 'Schiava grossa'-specific haplotype and 27% a  
4 'Pinot noir'-specific haplotype. It is estimated that 32% of the genome shares a common  
5 haplotype between the two parents, *i.e.*, that these regions could originate either from 'Pinot  
6 noir' or 'Schiava grossa' indicating that ~57% could originate from 'Schiava grossa' and ~43%  
7 from 'Pinot noir'.

8 The switch error rate was determined to 0.96% (REF), to 4.76% (ALT) and to 0.77%  
9 (PN12X.v2). Some of the switches are probably due to sequencing errors in the additional long  
10 read-based sequences. Moreover, as the error rate of ALT sequences was measured to ~4.76%,  
11 portions of the alternative sequences are a mixture of the maternal and paternal haplotype,  
12 confirming that despite the improved separation of the two haplotypes in PN40024.v4, phasing  
13 is still not perfect.

14 By exploring the VIVC database ([www.vivc.de](http://www.vivc.de)), the 'Helfensteiner' cultivar was found to  
15 originate from a cross between 'Pinot noir precoce' (a clone of 'Pinot noir') and 'Schiava  
16 grossa'. By performing the same variant calling analysis, 53,671 homozygous variants were  
17 found between cv. 'Helfensteiner' and PN40024.v4, with 543 homozygous variants / Mb in the  
18 heterozygous regions and 93 homozygous variants / Mb in the homozygous regions (Figure 5).  
19 As a negative control, 'Araklinos' showed 3,967 homozygous variants / Mb in the  
20 heterozygous regions and 4,818 homozygous variants / Mb in the homozygous regions). Thus,  
21 the 'Helfensteiner' homozygous variants are almost six times denser in error-prone regions of  
22 the PN40024.v4 assembly, which makes them probable "false positive" homozygous variants.  
23 Apart from heterozygous regions, no blocks of homozygous variants could be identified,  
24 meaning that one of the two 'Helfensteiner' haplotypes is always present in the PN40024  
25 genome. This confirms that the 'Helfensteiner' variety is the true parent of the first selfing,  
26 from which the PN40024 genotype was created after eight more selfings.

27

### 28 **PN40024.v4.1 gene prediction, functional annotation and manual curation**

29 The PN40024.v4.1 gene annotation of REF haplotype comprises 35,922 gene models of which  
30 35,197 are protein-coding and 725 encode for tRNAs (Table 4). In particular, 1,572 novel  
31 protein-coding genes were annotated in the newly assembled long read-based regions. For  
32 heterozygous regions, 1,855 and 1,809 protein-coding genes were predicted for REF and ALT

1 haplotypes, respectively (Table 5). Most genes were predicted on the ~11 Mb heterozygous  
2 region on chromosome 7 with 830 on the reference sequence and 792 on the alternative  
3 sequence followed by the ~5 Mb region on chromosome 10 with 650 and 623 protein-coding  
4 genes.

5 To check for completeness of the gene models, the plant core genes of the database  
6 eudicots\_odb10 were predicted with BUSCO (Figure 6). Of the 2,326 searched plant core  
7 genes, 2,296 or 98.7% were classified as complete in the PN40024.v4.1 gene annotation. Only  
8 16 were predicted as fragmented and only 14 were not found.

9 Compared to PN12X.v2 VCost.v3 gene annotation, PN40024.v4.1 counts less predictions  
10 (41,182 versus 35,197) but their size is longer on average (4,485 bp versus 4,742 bp) (Table 4).  
11 Also, the BUSCO analysis performed on VCost.v3 showed that 2,257 or 97.0% were classified  
12 as complete (Figure 6). Thus, PN40024.v4.1 gene annotation represents PN40024 gene space  
13 in a more exhaustive and less fragmented manner compared to VCost.v3.

14

15 To help the community in the transfer of information across versions (*i.e.*, correspondences),  
16 we retained as many gene names from VCost.v3 in PN40024.v4.1 as possible. We adopted a  
17 strategy based on RBHs followed by some filtering steps which allowed us to transfer names  
18 for 66% (23,206) of PN40024.v4.1 gene models with the nomenclature VitviXXg0YYYY (XX  
19 being the chromosome number and YYYY a sequential number below 4,000). One third  
20 (11,991) of PN40024.v4.1 gene models could not be named with a VCost.v3 identifier and  
21 were named with the nomenclature VitviXXg04ZZZ (XX being the chromosome number and  
22 ZZZ a sequential number below 1,000). The detailed nomenclature for PN40024.v4.1 gene  
23 annotations is given in File S2 Table S4.

24 The functional annotation of PN40024.v4.1 was performed using Blast2GO and resulted in at  
25 least one Gene Ontology term for 87% (30,689) of the genes and one Enzyme Code for 41%  
26 (14,512) of them. The main classes and ontologies are detailed in File S1 Figure S6 and Figure  
27 S7.

28 A subset of the RNA-Seq data published by Palumbo *et al.* (Palumbo *et al.* 2014) was used to  
29 compare the results of a differential gene expression analysis performed with  
30 PN12X.v2/VCost.v3 and PN40024.v4/PN40024.v4.1. In terms of mapping, the percentage of  
31 aligned reads was equivalent or slightly better when using PN40024.v4 genome assembly  
32 compared to PN12X.v2 (File S2 Table S6). Additionally, the percentage of assigned reads, *i.e.*,

1 the percentage of reads aligned under an annotated gene, was 2.4 to 3% better with  
2 PN40024.v4/PN40024.v4.1 compared to PN12X.v2/VCost.v3, which confirms the improved  
3 quality of PN40024.v4.1 gene annotation. Moreover, after differential gene expression  
4 analysis, the use of PN40024.v4/PN40024.v4.1 allowed the identification of more  
5 differentially expressed genes than PN12X.v2/VCost.v3 (File S1 Figure S8). This result along  
6 with the exhaustive functional annotation of PN40024.v4.1 shows that this new version of the  
7 PN40024 reference genome and annotation is a very efficient resource to perform  
8 transcriptomics and functional enrichment analyses.

9 Despite marked improvement of the PN40024.v4.1 automated annotation with respect to the  
10 previous VCost.v3 annotation, some recently expanded gene families have not been  
11 comprehensively annotated, such as the stilbene synthase (STS) gene family. Therefore, 1,641  
12 genes (1,579 edited and 62 deleted) were manually curated using a purpose-built Apollo server  
13 (<http://138.102.159.70:8080/apollo>) providing a wide range of transcriptomic and genomic  
14 data for PN40024.v4. In an effort to preserve previous VCost.v3 manual curation and  
15 functional annotation efforts, a particular focus was given to genes present in the reference  
16 catalogue (Navarro-Payá et al, 2022). The PN40024.v4.1 automated annotation including the  
17 manually curated features was called PN40024.v4.2, which metrics are presented in (Table 4).  
18 An automated annotation from PN40024.v4.1 that was manually curated was deleted and  
19 replaced by its curated version in PN40024.v4.2. Also, the same rules were applied for gene  
20 name transfer and nomenclature for PN40024.v4.1 and PN40024.v4.2. The BUSCO analysis  
21 performed on PN40024.v4.2 shows that the fragmented plant core genes were reduced to six  
22 and the missing genes to eight (Figure 6). Thus, PN40024.v4.2 gene models comprise 2,308 or  
23 99.2% complete plant core genes.

24  
25 In conclusion, the here provided PN40024.v4 assembly is the most suitable grapevine reference  
26 genome sequence assembly as it notably outperforms PN12X.v2. In terms of genomic and  
27 transcriptomic read mapping, the assembly also outperforms other high-quality *V. vinifera*  
28 genome assemblies, something that occurs even when reads from these recently sequenced  
29 cultivars are used. Having a fully resolved alternative haplotype sequence, more continuous  
30 sequences and resolving many up-to-now unknown bases, PN40024.v4 represents the near  
31 complete diploid genome of the PN40024 genotype. Despite many improvements and advances  
32 in PN40024.v4, the genome sequence is still not perfect in regard to haplotype switching and  
33 newly introduced errors by implementation of long genomic reads. Further improvements

1 should focus on these regions. Nevertheless, the gene annotation of PN40024.v4 should be  
2 used as the most updated resource for transcriptomics and functional enrichment analyses,  
3 while the genes of heterozygous regions that are likely represented on both haplotypes will  
4 allow exploration of heterozygous genetic traits.

## 6 Data availability statement

7 Supplemental files are provided with the manuscript. File S1 contains additional figures and  
8 File S2 additional tables. Raw sequencing data and the PN40024.v4 genome assembly are  
9 available at ENA under BioProject PRJEB45423. Also, the PN40024.v4 genome assembly  
10 with structural and functional gene annotation is available on the INTEGRAPE website  
11 (<https://integrape.eu/resources/genes-genomes/genome-accessions>), on the Grape Genomics  
12 Encyclopedia portal (<http://grapedia.org/>) and under the DOI number doi:10.57745/F9N2FZ  
13 (<https://entrepot.recherche.data.gouv.fr/dataset.xhtml?persistentId=doi:10.57745/F9N2FZ>). A  
14 Sequence Server v2.0.0 interface (<http://138.102.159.70:4567/>) was set up to perform BLAST  
15 analyses. A JBrowse interface (<http://138.102.159.70/jbrowse/>) was set up to visualize  
16 PN40024.v4 assembly and PN40024.v4.1 and v4.2 annotations, but also some previous  
17 annotation versions that were transferred, some RNA-Seq alignments and miscellaneous  
18 tracks. An Apollo interface (<http://138.102.159.70:8080/apollo>; training and account  
19 mandatory) was set up to manually curate gene annotations according to the dedicated  
20 guidelines (<https://integrape.eu/resources/data-management/>). Code used to analyze GBS data  
21 can be found at <https://forgemia.inra.fr/sophie.blanc/gbs> and code used to generate the  
22 PN40024.v4.2 version can be found at [https://gitlab.com/MSVteam/pn40024-visualization-](https://gitlab.com/MSVteam/pn40024-visualization-tools/-/tree/master/update_gff3_script)  
23 [tools/-/tree/master/update\\_gff3\\_script](https://gitlab.com/MSVteam/pn40024-visualization-tools/-/tree/master/update_gff3_script).

## 25 Acknowledgements

26 We thank INRAE department “Biologie et Amélioration des Plantes” for funding; experimental  
27 unit UEAV (INRAE, Colmar) for plant maintenance; Anne-Marie Digby (University of  
28 Verona) for English correction; Emilce Prado for plant DNA extractions; EPGV (INRAE,  
29 Evry) for library prep and DNA sequencing; CNRGV (INRAE, Toulouse) for long read DNA  
30 extractions; Gentyane platform (INRAE, Clermont-Ferrand) for SMRT sequencing; Dr.

1 Timothée Flutre and Amandine Launay for their help in coordinating the FruitSelGen project  
2 and in acquiring GBS data; Mario Pezzotti, Anne-Françoise Adam-Blondon, Michele  
3 Morgante, Gabriele Di Gaspero and Gabriele Magris for helpful discussions throughout the  
4 project; Pablo Carbonell-Bejerano (Instituto de Ciencias de la Vid y el Vino - ICVV) for  
5 critically reviewing the manuscript. This work was supported by the BMBF-funded de.NBI  
6 Cloud within the German Network for Bioinformatics Infrastructure (de.NBI). This article is  
7 based upon work from COST Action CA17111 INTEGRAPE, supported by COST (European  
8 Cooperation in Science and Technology).

## 9 Contributions

10 A.V. performed genome assembly and functional annotation. B.F. performed gene prediction.  
11 B.F., N.V. and D.H. built the annotation pipeline. S.B. performed GBS analysis and built the  
12 genetic maps. É.D. built the genetic map pipeline. A.V., B.F. and C.R. performed the analysis  
13 for quality assessments. C.R. performed the analysis on PN40024's origin. V.D. performed  
14 plant material management and sampling. J.G. and B.F. worked on the gene name transfer.  
15 A.V., M.L. and D.N.-P. built the online tools. A.V., B.F., D.H., J.G., C.K., J.T.M, D.N.-P.,  
16 L.O., M.K.T.-R., D.W. and C.R. worked on gene manual curation and writing of the dedicated  
17 guidelines. É.D., P.H. and C.R. looked for funding. C.R. supervised the project. A.V., B.F.,  
18 D.H. and C.R. drafted and formatted the manuscript. All the authors read and helped improve  
19 the manuscript.

20

## 21 Conflict of interests

22 We declare no conflict of interests.

23

## 24 Funder Information

25 This work was supported by INRAE department "Biologie et Amélioration des Plantes", by  
26 the BMBF-funded de.NBI Cloud within the German Network for Bioinformatics Infrastructure  
27 (de.NBI) and by COST (European Cooperation in Science and Technology).



## Literature Cited

- 1 Akiva, E., S. Brown, D. E. Almonacid, A. E. Barber 2nd, A. F. Custer *et al.*, 2014 The  
2 Structure–Function Linkage Database. *Nucleic Acids Res.* 42: D521–D530.
- 3 Allen, J. E., and S. L. Salzberg, 2005 JIGSAW: integration of multiple sources of evidence  
4 for gene prediction. *Bioinformatics* 21: 3596–3603.
- 5 Altschul, S. F., W. Gish, W. Miller, E. W. Myers, and D. J. Lipman, 1990 Basic local  
6 alignment search tool. *J. Mol. Biol.* 215: 403–410.
- 7 Andrews, S., 2010 FASTQC. A quality control tool for high throughput sequence data.
- 8 Attwood, T. K., A. Coletta, G. Muirhead, A. Pavlopoulou, P. B. Philippou *et al.*, 2012 The  
9 PRINTS database: A fine-grained protein sequence annotation and analysis  
10 resource-its status in 2012. *Database* 2012: bas019.
- 11 Barnett, D. W., E. K. Garrison, A. R. Quinlan, M. P. Strömberg, and G. T. Marth, 2011  
12 BamTools: a C++ API and toolkit for analyzing and managing BAM files. *Bioinforma.*  
13 *Oxf. Engl.* 27: 1691–1692.
- 14 Bolger, A. M., M. Lohse, and B. Usadel, 2014 Trimmomatic: a flexible trimmer for Illumina  
15 sequence data. *Bioinformatics* 30: 2114–2120.
- 16 Brůna, T., K. J. Hoff, A. Lomsadze, M. Stanke, and M. Borodovsky, 2021 BRAKER2:  
17 automatic eukaryotic genome annotation with GeneMark-EP+ and AUGUSTUS  
18 supported by a protein database. *NAR Genomics Bioinforma.* 3: lqaa108.
- 19 Brůna, T., A. Lomsadze, and M. Borodovsky, 2020 GeneMark-EP+: eukaryotic gene  
20 prediction with self-training in the space of genes and proteins. *NAR Genomics*  
21 *Bioinforma.* 2:.
- 22 Buchfink, B., C. Xie, and D. H. Huson, 2015 Fast and sensitive protein alignment using  
23 DIAMOND. *Nat. Methods* 12: 59–60.
- 24 Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos *et al.*, 2009 BLAST+:  
25 architecture and applications. *BMC Bioinformatics* 10: 421.
- 26

- 1 Campbell, M. S., C. Holt, B. Moore, and M. Yandell, 2014 Genome Annotation and Curation  
2 Using MAKER and MAKER-P. *Curr. Protoc. Bioinforma.* 48: 4.11.1-39.
- 3 Canaguier, A., J. Grimplet, G. Di Gaspero, S. Scalabrin, E. Duchêne *et al.*, 2017 A new  
4 version of the grapevine reference genome assembly (12X.v2) and of its annotation  
5 (VCost.v3). *Genomics Data* 14: 56–62.
- 6 Chan, P. P., B. Y. Lin, A. J. Mak, and T. M. Lowe, 2021 tRNAscan-SE 2.0: improved  
7 detection and functional classification of transfer RNA genes. *Nucleic Acids Res.* 49:  
8 9077–9096.
- 9 Conesa, A., S. Götz, J. M. García-Gómez, J. Terol, M. Talón *et al.*, 2005 Blast2GO: a  
10 universal tool for annotation, visualization and analysis in functional genomics  
11 research. *Bioinformatics* 21: 3674–3676.
- 12 Dobin, A., C. A. Davis, F. Schlesinger, J. Drenkow, C. Zaleski *et al.*, 2013 STAR: ultrafast  
13 universal RNA-seq aligner. *Bioinformatics* 29: 15–21.
- 14 Duchêne, É., V. Dumas, G. Butterlin, N. Jaegli, C. Rustenholz *et al.*, 2020 Genetic variations  
15 of acidity in grape berries are controlled by the interplay between organic acids and  
16 potassium. *Theor. Appl. Genet.* 133: 993–1008.
- 17 Dunn, N. A., D. R. Unni, C. Diesh, M. Munoz-Torres, N. L. Harris *et al.*, 2019 Apollo:  
18 Democratizing genome annotation. *PLOS Comput. Biol.* 15: e1006790.
- 19 Eichhorn, K. W., and D. H. Lorenz, 1977 Phanologische Entwicklungsstadien der Rebe.  
20 *Nachrichtenblatt Dtsch. Pflanzenschutzdienstes.*
- 21 Gao, S., D. Bertrand, B. K. H. Chia, and N. Nagarajan, 2016 OPERA-LG: efficient and exact  
22 scaffolding of large, repeat-rich eukaryotic genomes with performance guarantees.  
23 *Genome Biol.* 17: 102.
- 24 Girgis, H. Z., 2015 Red: an intelligent, rapid, accurate tool for detecting repeats de-novo on  
25 the genomic scale. *BMC Bioinformatics* 16: 227.
- 26 Girollet, N., B. Rubio, C. Lopez-Roques, S. Valière, N. Ollat *et al.*, 2019 De novo phased  
27 assembly of the *Vitis riparia* grape genome. *Sci. Data* 6: 127.

- 1 Gotoh, O., 2008 A space-efficient and accurate method for mapping and aligning cDNA  
2 sequences onto genomic sequence. *Nucleic Acids Res.* 36: 2630–2638.
- 3 Götz, S., J. M. García-Gómez, J. Terol, T. D. Williams, S. H. Nagaraj *et al.*, 2008 High-  
4 throughput functional annotation and data mining with the Blast2GO suite. *Nucleic*  
5 *Acids Res.* 36: 3420–3435.
- 6 Gough, J., K. Karplus, R. Hughey, and C. Chothia, 2001 Assignment of homology to genome  
7 sequences using a library of hidden Markov models that represent all proteins of  
8 known structure. *J. Mol. Biol.* 313: 903–919.
- 9 Haas, B. J., S. L. Salzberg, W. Zhu, M. Pertea, J. E. Allen *et al.*, 2008 Automated eukaryotic  
10 gene structure annotation using EVIDENCEModeler and the Program to Assemble  
11 Spliced Alignments. *Genome Biol.* 9: R7.
- 12 Haft, D. H., J. D. Selengut, R. A. Richter, D. Harkins, M. K. Basu *et al.*, 2013 TIGRFAMs and  
13 Genome Properties in 2013. *Nucleic Acids Res.* 41: D387-395.
- 14 Hoff, K. J., S. Lange, A. Lomsadze, M. Borodovsky, and M. Stanke, 2016 BRAKER1:  
15 Unsupervised RNA-Seq-Based Genome Annotation with GeneMark-ET and  
16 AUGUSTUS. *Bioinforma. Oxf. Engl.* 32: 767–769.
- 17 Hoff, K. J., A. Lomsadze, M. Borodovsky, and M. Stanke, 2019 Whole-Genome Annotation  
18 with BRAKER. *Methods Mol. Biol. Clifton NJ* 1962: 65–95.
- 19 Holt, C., and M. Yandell, 2011 MAKER2: an annotation pipeline and genome-database  
20 management tool for second-generation genome projects. *BMC Bioinformatics* 12:  
21 491.
- 22 Howe, K. L., T. Chothia, and R. Durbin, 2002 GAZE: A Generic Framework for the  
23 Integration of Gene-Prediction Data by Dynamic Programming. *Genome Res.* 12:  
24 1418–1427.
- 25 Huang, S., M. Kang, and A. Xu, 2017 HaploMerger2: rebuilding both haploid sub-assemblies  
26 from high-heterozygosity diploid genome assembly. *Bioinformatics* 33: 2577–2579.

- 1 Iwata, H., and O. Gotoh, 2012 Benchmarking spliced alignment programs including Spaln2,  
2 an extended version of Spaln that incorporates additional species-specific features.  
3 Nucleic Acids Res. 40: e161.
- 4 Jaillon, O., J.-M. Aury, B. Noel, A. Policriti, C. Clepet *et al.*, 2007 The grapevine genome  
5 sequence suggests ancestral hexaploidization in major angiosperm phyla. Nature  
6 449: 463–467.
- 7 Jones, P., D. Binns, H.-Y. Chang, M. Fraser, W. Li *et al.*, 2014 InterProScan 5: genome-  
8 scale protein function classification. Bioinforma. Oxf. Engl. 30: 1236–1240.
- 9 Kent, W. J., 2002 BLAT--the BLAST-like alignment tool. Genome Res. 12: 656–664.
- 10 Killick, R., and I. A. Eckley, 2014 changepoint: An R Package for Changepoint Analysis. J.  
11 Stat. Softw. 58: 1–19.
- 12 Kolmogorov, M., J. Yuan, Y. Lin, and P. A. Pevzner, 2019 Assembly of long, error-prone  
13 reads using repeat graphs. Nat. Biotechnol. 37: 540–546.
- 14 Koren, S., B. P. Walenz, K. Berlin, J. R. Miller, N. H. Bergman *et al.*, 2017 Canu: scalable  
15 and accurate long-read assembly via adaptive k-mer weighting and repeat  
16 separation. Genome Res. 27: 722–736.
- 17 Korf, I., 2004 Gene finding in novel genomes. BMC Bioinformatics 5:.
- 18 Kurtz, S., A. Phillippy, A. L. Delcher, M. Smoot, M. Shumway *et al.*, 2004 Versatile and open  
19 software for comparing large genomes. Genome Biol. 5: R12.
- 20 Letunic, I., and P. Bork, 2018 20 years of the SMART protein domain annotation resource.  
21 Nucleic Acids Res. 46: D493–D496.
- 22 Li, H., 2018 Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics 34:  
23 3094–3100.
- 24 Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The Sequence  
25 Alignment/Map format and SAMtools. Bioinforma. Oxf. Engl. 25: 2078–2079.
- 26 Liao, Y., G. K. Smyth, and W. Shi, 2014 featureCounts: an efficient general purpose program  
27 for assigning sequence reads to genomic features. Bioinforma. Oxf. Engl. 30: 923–  
28 930.

- 1 Lodhi, M. A., and B. I. Reisch, 1995 Nuclear DNA content of *Vitis* species, cultivars, and  
2 other genera of the Vitaceae. *Theor. Appl. Genet.* 90: 11–16.
- 3 Lomsadze, A., P. D. Burns, and M. Borodovsky, 2014 Integration of mapped RNA-Seq reads  
4 into automatic training of eukaryotic gene finding algorithm. *Nucleic Acids Res.* 42:  
5 e119.
- 6 Lomsadze, A., V. Ter-Hovhannisyanyan, Y. O. Chernoff, and M. Borodovsky, 2005 Gene  
7 identification in novel eukaryotic genomes by self-training algorithm. *Nucleic Acids*  
8 *Res.* 33: 6494–6506.
- 9 Lu, S., J. Wang, F. Chitsaz, M. K. Derbyshire, R. C. Geer *et al.*, 2020 CDD/SPARCLE: the  
10 conserved domain database in 2020. *Nucleic Acids Res.* 48: D265–D268.
- 11 Luo, R., B. Liu, Y. Xie, Z. Li, W. Huang *et al.*, 2012 SOAPdenovo2: an empirically improved  
12 memory-efficient short-read de novo assembler. *GigaScience* 1: 2047-217X-1–18.
- 13 Lupas, A., M. Van Dyke, and J. Stock, 1991 Predicting coiled coils from protein sequences.  
14 *Science* 252: 1162–1164.
- 15 Majoros, W. H., M. Pertea, and S. L. Salzberg, 2004 TigrScan and GlimmerHMM: two open  
16 source ab initio eukaryotic gene-finders. *Bioinforma. Oxf. Engl.* 20: 2878–2879.
- 17 Massonnet, M., N. Cochetel, A. Minio, A. M. Vondras, J. Lin *et al.*, 2020 The genetic basis of  
18 sex determination in grapes. *Nat. Commun.* 11: 2902.
- 19 McKenna, A., M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis *et al.*, 2010 The Genome  
20 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA  
21 sequencing data. *Genome Res.* 20: 1297–1303.
- 22 Merdinoglu, D., G. Butterlin, L. Bevilacqua, V. Chiquet, A.-F. Adam-Blondon *et al.*, 2005  
23 Development and characterization of a large set of microsatellite markers in  
24 grapevine (*Vitis vinifera* L.) suitable for multiplex PCR. *Mol. Breed.* 15: 349–366.
- 25 Mi, H., D. Ebert, A. Muruganujan, C. Mills, L.-P. Albou *et al.*, 2021 PANTHER version 16: a  
26 revised family classification, tree-based classification tool, enhancer regions and  
27 extensive API. *Nucleic Acids Res.* 49: D394–D403.

- 1 Mistry, J., S. Chuguransky, L. Williams, M. Qureshi, G. A. Salazar *et al.*, 2021 Pfam: The  
2 protein families database in 2021. *Nucleic Acids Res.* 49: D412–D419.
- 3 Navarro-Payá, D., A. Santiago, L. Orduña, C. Zhang, A. Amato *et al.*, 2022 The Grape Gene  
4 Reference Catalogue as a Standard Resource for Gene Selection and Genetic  
5 Improvement. *Front. Plant Sci.* 12:.
- 6 Necci, M., D. Piovesan, Z. Dosztányi, and S. C. E. Tosatto, 2017 MobiDB-lite: fast and highly  
7 specific consensus prediction of intrinsic disorder in proteins. *Bioinforma. Oxf. Engl.*  
8 33: 1402–1404.
- 9 Palumbo, M. C., S. Zenoni, M. Fasoli, M. Massonnet, L. Farina *et al.*, 2014 Integrated  
10 Network Analysis Identifies Fight-Club Nodes as a Class of Hubs Encompassing Key  
11 Putative Switch Genes That Induce Major Transcriptome Reprogramming during  
12 Grapevine Development. *Plant Cell* 26: 4617–4635.
- 13 Pedruzzi, I., C. Rivoire, A. H. Auchincloss, E. Coudert, G. Keller *et al.*, 2013 HAMAP in 2013,  
14 new developments in the protein family classification and annotation system. *Nucleic  
15 Acids Res.* 41: D584-589.
- 16 Quinlan, A. R., and I. M. Hall, 2010 BEDTools: a flexible suite of utilities for comparing  
17 genomic features. *Bioinforma. Oxf. Engl.* 26: 841–842.
- 18 Rhie, A., B. P. Walenz, S. Koren, and A. M. Phillippy, 2020 Merqury: reference-free quality,  
19 completeness, and phasing assessment for genome assemblies. *Genome Biol.* 21:  
20 245.
- 21 Sallet, E., J. Gouzy, and T. Schiex, 2019 EuGene: An Automated Integrative Gene Finder for  
22 Eukaryotes and Prokaryotes. *Methods Mol. Biol. Clifton NJ* 1962: 97–120.
- 23 Salmela, L., and E. Rivals, 2014 LoRDEC: accurate and efficient long read error correction.  
24 *Bioinformatics* 30: 3506–3514.
- 25 Shumate, A., and S. L. Salzberg, 2021 Liftoff: accurate mapping of gene annotations.  
26 *Bioinformatics* 37: 1639–1643.
- 27 Sigrist, C. J. A., E. de Castro, L. Cerutti, B. A. Cuče, N. Hulo *et al.*, 2013 New and  
28 continuing developments at PROSITE. *Nucleic Acids Res.* 41: D344-347.

- 1 Sillitoe, I., N. Dawson, T. E. Lewis, S. Das, J. G. Lees *et al.*, 2019 CATH: expanding the  
2 horizons of structure-based functional annotations for genome sequences. *Nucleic*  
3 *Acids Res.* 47: D280–D284.
- 4 Simão, F. A., R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva, and E. M. Zdobnov, 2015  
5 BUSCO: assessing genome assembly and annotation completeness with single-copy  
6 orthologs. *Bioinformatics* 31: 3210–3212.
- 7 Slater, G. S. C., and E. Birney, 2005 Automated generation of heuristics for biological  
8 sequence comparison. *BMC Bioinformatics* 6: 31.
- 9 Smit, A. F. A., R. Hubley, and P. Green, 2013 RepeatMasker Open-4.0.
- 10 Song, L., S. Sabunciyar, G. Yang, and L. Florea, 2019 A multi-sample approach increases  
11 the accuracy of transcript assembly. *Nat. Commun.* 10: 5000.
- 12 Stanke, M., M. Diekhans, R. Baertsch, and D. Haussler, 2008 Using native and syntenically  
13 mapped cDNA alignments to improve de novo gene finding. *Bioinformatics* 24: 637–  
14 644.
- 15 Stanke, M., O. Schöffmann, B. Morgenstern, and S. Waack, 2006 Gene prediction in  
16 eukaryotes with a generalized hidden Markov model that uses hints from external  
17 sources. *BMC Bioinformatics* 7: 62.
- 18 Tang, H., X. Zhang, C. Miao, J. Zhang, R. Ming *et al.*, 2015 ALLMAPS: robust scaffold  
19 ordering based on multiple maps. *Genome Biol.* 16: 3.
- 20 Taylor, J., and D. Butler, 2017 R Package ASMap: Efficient Genetic Linkage Map  
21 Construction and Diagnosis. *J. Stat. Softw.* 79: 1–29.
- 22 Torkamaneh, D., J. Laroche, M. Bastien, A. Abed, and F. Belzile, 2017 Fast-GBS: a new  
23 pipeline for the efficient and highly accurate calling of SNPs from genotyping-by-  
24 sequencing data. *BMC Bioinformatics* 18: 5.
- 25 Varet, H., L. Brillet-Guéguen, J.-Y. Coppée, and M.-A. Dillies, 2016 SARTools: A DESeq2-  
26 and EdgeR-Based R Pipeline for Comprehensive Differential Analysis of RNA-Seq  
27 Data. *PLoS One* 11: e0157022.



- 1 Vasimuddin, Md., S. Misra, H. Li, and S. Aluru, 2019 Efficient Architecture-Aware  
2 Acceleration of BWA-MEM for Multicore Systems, pp. 314–324 in *2019 IEEE*  
3 *International Parallel and Distributed Processing Symposium (IPDPS)*,.
- 4 Velasco, R., A. Zharkikh, M. Troggio, D. A. Cartwright, A. Cestaro *et al.*, 2007 A High Quality  
5 Draft Consensus Sequence of the Genome of a Heterozygous Grapevine Variety.  
6 PLOS ONE 2: e1326.
- 7 Vitulo, N., C. Forcato, E. C. Carpinelli, A. Telatin, D. Campagna *et al.*, 2014 A deep survey of  
8 alternative splicing in grape reveals changes in the splicing machinery related to  
9 tissue, stress condition and genotype. BMC Plant Biol. 14: 99.
- 10 Walker, B. J., T. Abeel, T. Shea, M. Priest, A. Abouelliel *et al.*, 2014 Pilon: An Integrated  
11 Tool for Comprehensive Microbial Variant Detection and Genome Assembly  
12 Improvement. PLOS ONE 9: e112963.
- 13 Wang, M., and L. Kong, 2019 pblat: a multithread blat algorithm speeding up aligning  
14 sequences to genomes. BMC Bioinformatics 20: 28.
- 15 Waterhouse, R. M., M. Seppey, F. A. Simão, M. Manni, P. Ioannidis *et al.*, 2018 BUSCO  
16 Applications from Quality Assessments to Gene Prediction and Phylogenomics. Mol.  
17 Biol. Evol. 35: 543–548.
- 18 Wilson, D., R. Pethica, Y. Zhou, C. Talbot, C. Vogel *et al.*, 2009 SUPERFAMILY--  
19 sophisticated comparative genomics, data mining, visualization and phylogeny.  
20 Nucleic Acids Res. 37: D380-386.
- 21 Wu, C. H., A. Nikolskaya, H. Huang, L.-S. L. Yeh, D. A. Natale *et al.*, 2004 PIRSF: family  
22 classification system at the Protein Information Resource. Nucleic Acids Res. 32:  
23 D112-114.
- 24 Wu, T. D., and C. K. Watanabe, 2005 GMAP: a genomic mapping and alignment program  
25 for mRNA and EST sequences. Bioinformatics 21: 1859–1875.
- 26 Xu, G.-C., T.-J. Xu, R. Zhu, Y. Zhang, S.-Q. Li *et al.*, 2019 LR\_Gapcloser: a tiling path-based  
27 gap closer that uses long reads to complete genome assembly. GigaScience 8:  
28 giy157.

1

2

3 **Table 1: Assembly statistics of the PN40024.v4 REF and PN12X.v2 genome assembly.** The  
4 table lists statistics of the PN40024.v4 REF and PN12X.v2 scaffolded and chromosome-  
5 anchored genome assemblies. ‘N’ denotes the number (No.) of unknown bases.

SCAFFOLDS	No. Scaf.	Min. size [bp]	Avg. size [kb]	Median size [kb]	L50	N50 [Mb]	Max. size [Mb]	Sum [Mb]	No. Ns	GC [%]
<b>PN12X.v2</b>	2,059	2,001	236	5	41	3.43	13.10	485.19	14,976,411	33.5
<b>PN40024.v4</b>	640	542	742	20	25	6.50	15.23	474.61	1,755,062	34.4
<b>Anchored PN12X.v2</b>	367	2,010	1,250	277	37	3.57	13.10	465.64	11,921,253	33.6
<b>Anchored PN40024.v4</b>	165	1,085	2,801	1,506	24	6.57	15.23	462.14	1,631,047	34.4

6

7 **Table 2: Assembly quality values of PN40024.v4 and 12X.v2.** Assembly quality values  
8 measured by Merqury for PN40024.v4 and 12X.v2 genome assemblies. QV denotes base  
9 level quality value.

	<b>12X.v2</b>	<b>PN40024.v4</b>
<b>QV</b>	37.4338	36.0171
<b>Error rate [%]</b>	0.000180559	0.000250199
<b>k-mer completeness [%]</b>	96.79	96.96

10

11 **Table 3: Haplotype block statistics of PN40024.v4 and 12X.v2.** Phasing accuracy  
12 estimation of Merqury for PN40024.v4 and 12X.v2 genome assembly. ALT denotes alternative  
13 heterozygous sequence parts of PN40024.v4.

	<b>12X.v2</b>	<b>PN40024.v4</b>	<b>ALT</b>

<b>Number of blocks</b>	2,575	1,454	289
<b>Total bases in blocks [bp]</b>	474,845,411	468,703,133	19,519,697
<b>Block N50 size [kb]</b>	1,762	2,050	250
<b>Switch error rate [%]</b>	0.766002	0.959042	4.75944

1

2 **Table 4: VCost.v3, PN40024.v4.1 REF haplotype and PN40024.v4.2 REF haplotype gene**  
3 **prediction overview.**

	<b>VCost.v3</b>		<b>PN40024.v4.1</b>		<b>PN40024.v4.2</b>	
	<b>Number</b>	<b>Mean length [bp]</b>	<b>Number</b>	<b>Mean length [bp]</b>	<b>Number</b>	<b>Mean length [bp]</b>
<b>Protein-coding genes</b>	41,182	4,485	35,197	4,742	35,230	4,735
<b>Transcripts</b>	47,363	1,383	41,160	1,433	41,173	1,440
<b>Exons</b>	239,165	273	208,581	282	208,719	283
<b>CDS</b>	225,869	220	199,956	231	200,059	232
<b>5' UTRs</b>	26,024	259	17,019	280	17,478	275
<b>3' UTRs</b>	26,994	327	17,873	440	18,344	433
<b>tRNAs</b>	19	74	725	75	725	75

4

5 **Table 5: Gene numbers of heterozygous sequence regions.** The abbreviation ALT denotes  
6 the alternative heterozygous sequence regions.

	<b>Bases [bp]</b>		<b>Number of genes</b>	
	<b>PN40024.v4</b>	<b>ALT</b>	<b>PN40024.v4.1</b>	<b>ALT</b>
<b>chr02</b>	1,610,271	1,886,900	190	214

<b>chr03</b>	288,001	287,774	14	13
<b>chr04</b>	1,049,642	929,781	123	122
<b>chr07</b>	11,422,405	10,851,409	830	792
<b>chr10</b>	5,475,057	5,100,371	650	623
<b>chr11</b>	733,078	630,772	43	41
<b>chr15</b>	60,730	52,641	5	4
<b>TOTAL</b>	<b>20,639,184</b>	<b>19,739,648</b>	<b>1,855</b>	<b>1,809</b>

1

2 **Figure 1: Assembly process for the PN40024.v4 genome sequence assembly.** A) Initial  
3 datasets: Sanger-based scaffolds of PN12X.v2 (red) with unknown bases ('N's'), genomic  
4 SMRT reads (green) and genomic short reads (blue). Erroneous bases are represented by a  
5 black line. B) Scaffold assembly steps. Same color code as A). Dark green regions represent  
6 newly incorporated SMRT sequencing regions. C) Pseudomolecule construction using the new  
7 scaffolds and genetic maps. Sequence regions of 12X.v2 are colored red and newly  
8 incorporated SMRT sequencing regions are colored dark green.

9 **Figure 2: Percentage of mapped genomic reads (A) and percentage of properly paired**  
10 **genomic reads (B) between PN40024.v4, PN12X.v2 and cv. 'Cabernet Sauvignon'**  
11 **(Massonnet et al., 2020) for 11 paired-end resequencing datasets of *V. vinifera* cultivars.**  
12 The x-axis denotes the source (cultivar) of the genomic reads and the y-axis the percentage [%]  
13 of mapped reads. Note that the PN40024 dataset was obtained with Illumina Genome Analyzer  
14 Iix sequencing and all other samples with Illumina HiSeq 4000. The PN40024 dataset is  
15 therefore of lower quality than the others.

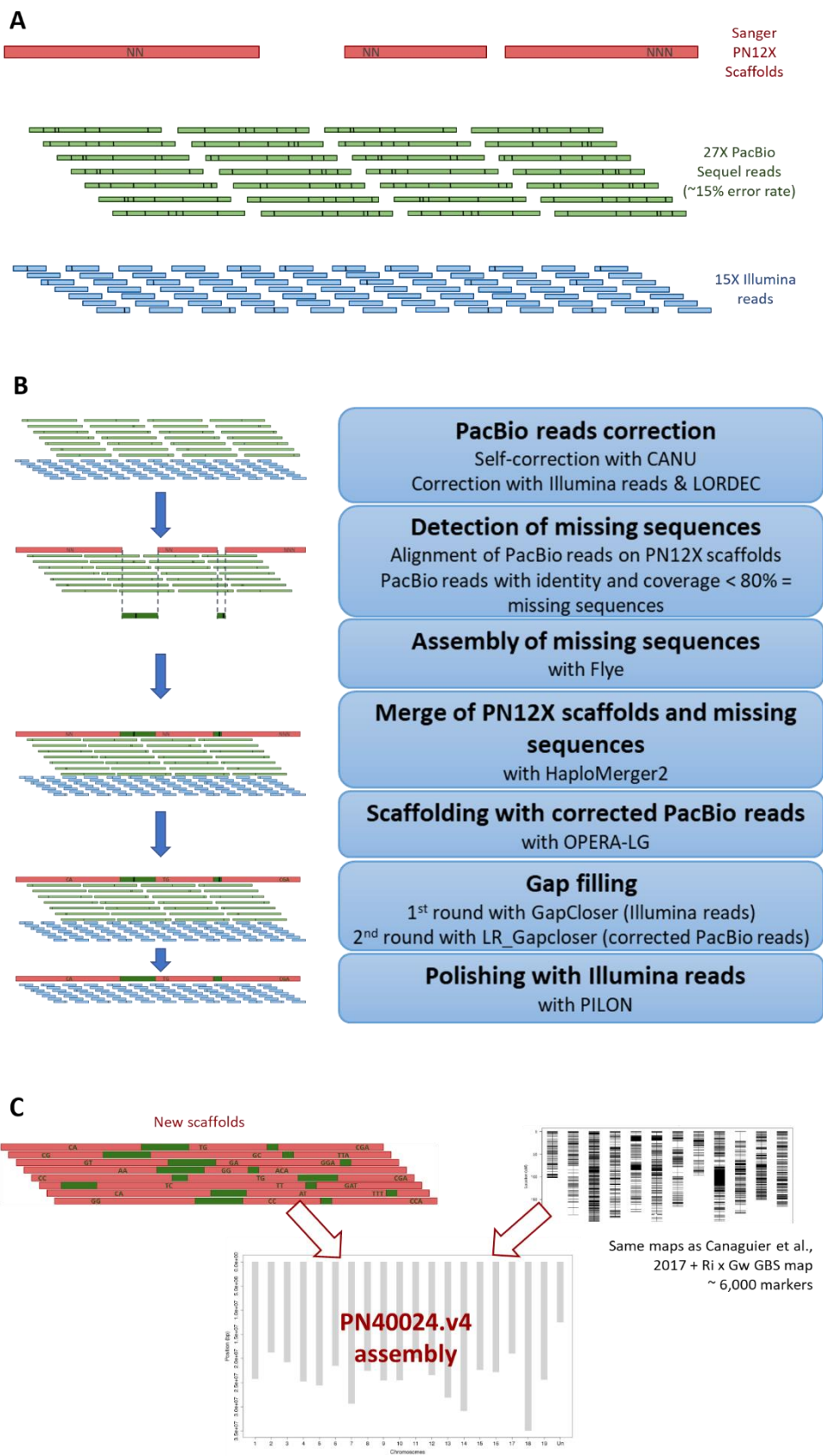
16 **Figure 3: Location of regions assembled using long reads, density of errors and of**  
17 **heterozygous SNPs in the PN40024.v4 genome sequence assembly.** The x-axis shows the  
18 19 main pseudo-chromosomes and the artificial chrUn ('Un'). The y-axis shows the base  
19 position in [bp]. 'Pacbio regions' refers to sequences derived from genomic SMRT reads. The  
20 seven heterozygous regions are squared in green.

1 **Figure 4: Density of ‘Pinot noir’ and ‘Schiava grossa’ homozygous SNPs compared to the**  
2 **PN40024.v4 genome assembly.** The x-axis shows the 19 main pseudochromosomes and the  
3 artificial chrUn (‘Un’). The y-axis shows the base position in [bp]. Where density of ‘Pinot  
4 noir’ SNPs is high, it means PN40024.v4 carries the ‘Schiava grossa’ haplotype and *vice versa*.  
5 The regions where both ‘Pinot noir’ and ‘Schiava grossa’ SNP density is low correspond to  
6 regions where both genomes share a common haplotype.

7  
8 **Figure 5: Density of ‘Helfensteiner’ homozygous SNPs compared to the PN40024.v4**  
9 **genome assembly.** The x-axis shows the 19 main pseudochromosomes and the artificial chrUn  
10 (‘Un’). The y-axis shows the base position in [bp]. The seven regions squared in green are the  
11 heterozygous regions.

12  
13  
14 **Figure 6: Plant core genes of the PN40024.v4 and PN12X.v2 genome assemblies and their**  
15 **annotations.** The 2,326 plant core genes of the database eudicots\_odb10 were determined in  
16 the PN40024.v4 genome assembly, in its annotation PN40024.v4.1, in the PN12X.v2 genome  
17 assembly and in the VCost.v3 gene annotation. ‘PN40024.v4.2’ is the PN40024.v4 gene  
18 annotation after manual curation of the fragmented and missing plant core genes.

19



1

2

*Figure 1*  
254x447 mm ( x DPI)

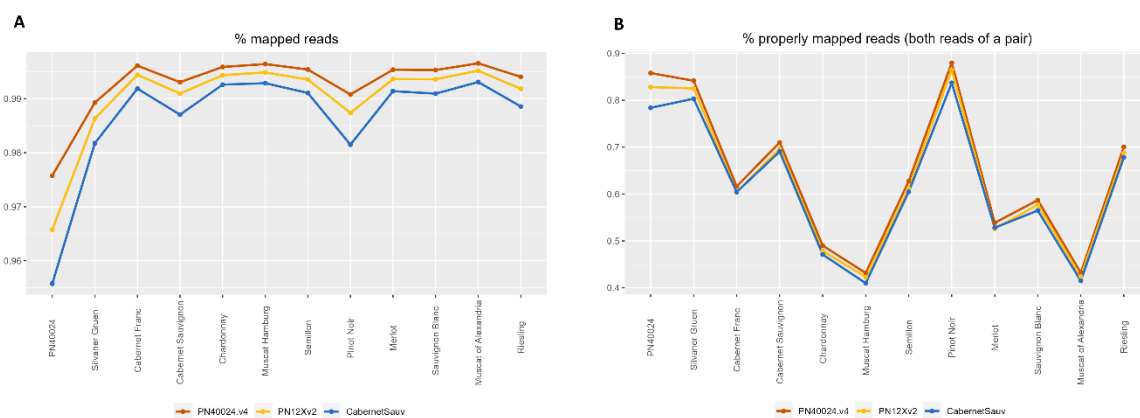


Figure 2  
406x152 mm ( x DPI)

1  
2  
3  
4  
5

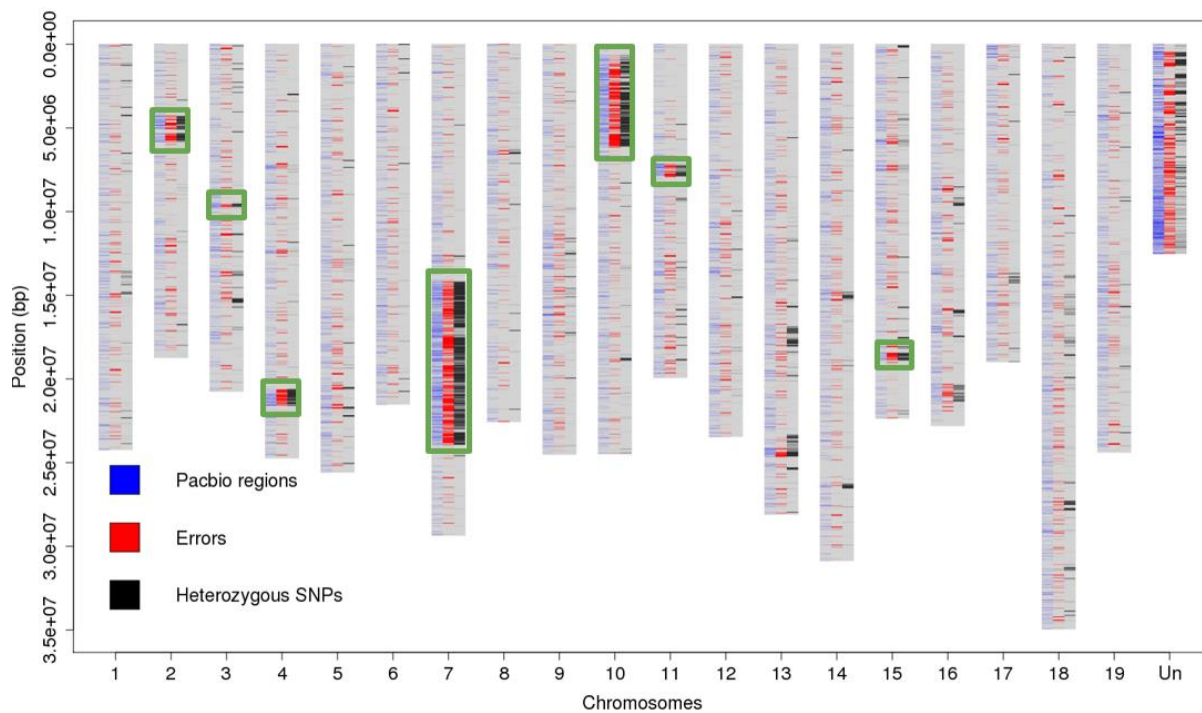


Figure 3  
184x108 mm ( x DPI)

6

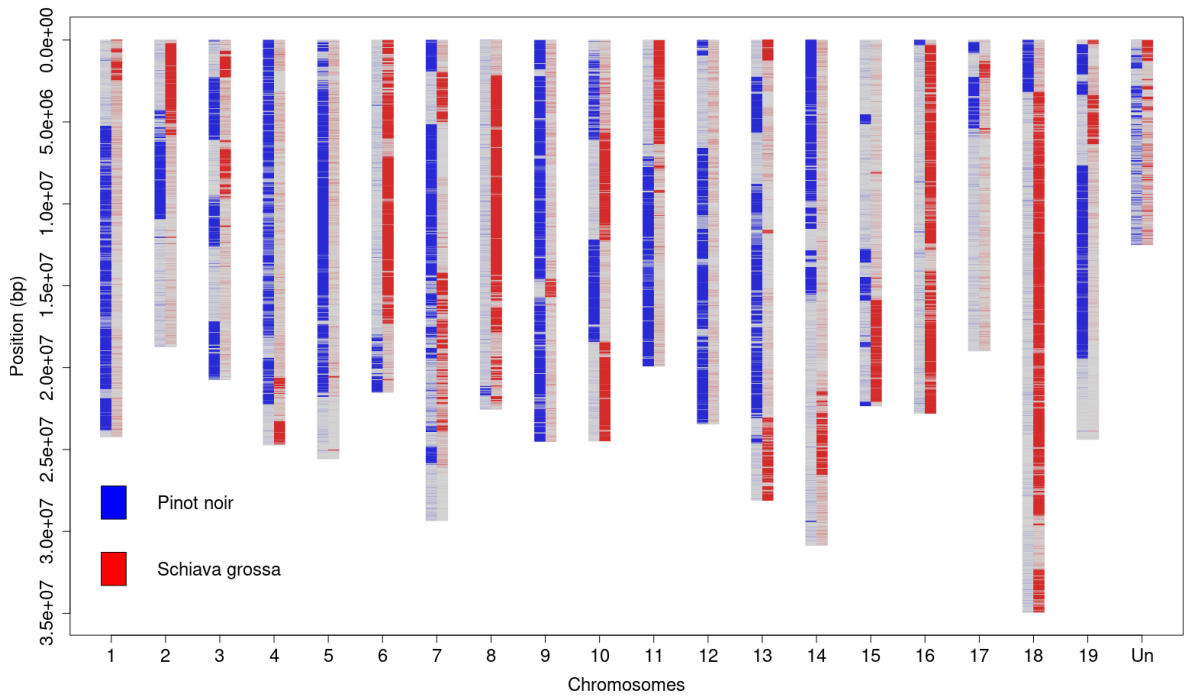


Figure 4  
418x261 mm ( x DPI)

- 1
- 2
- 3

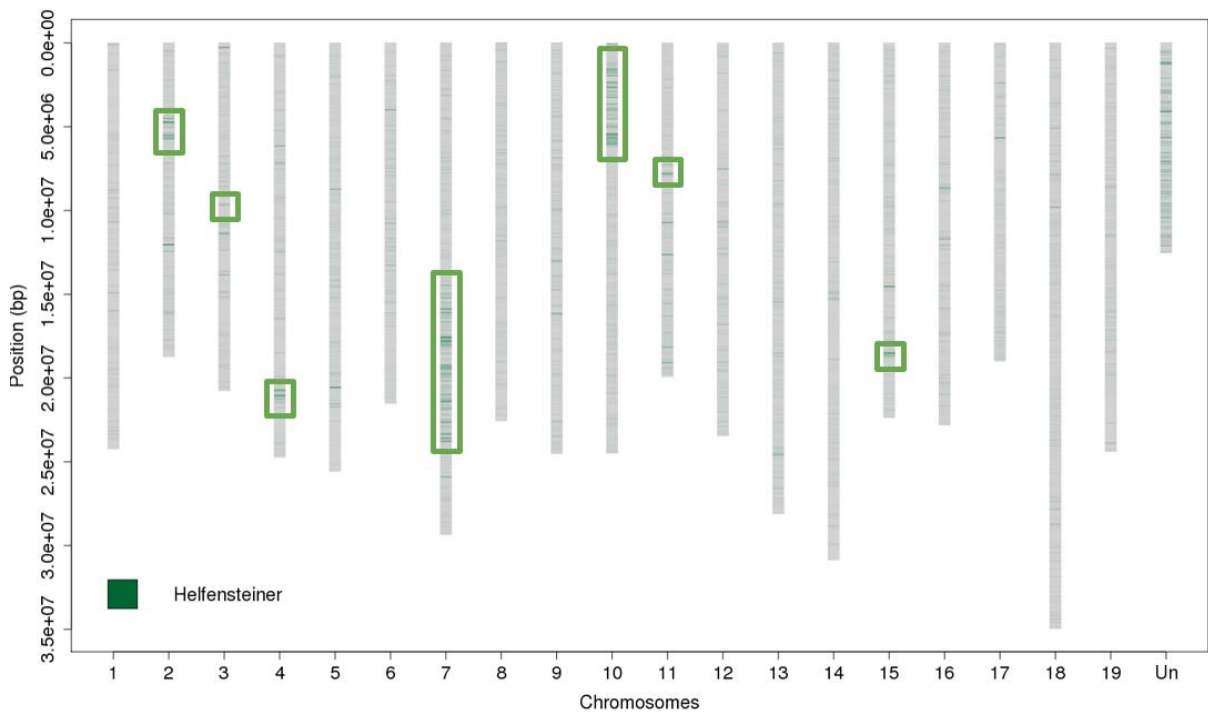


Figure 5  
184x108 mm ( x DPI)



1  
2  
3  
4

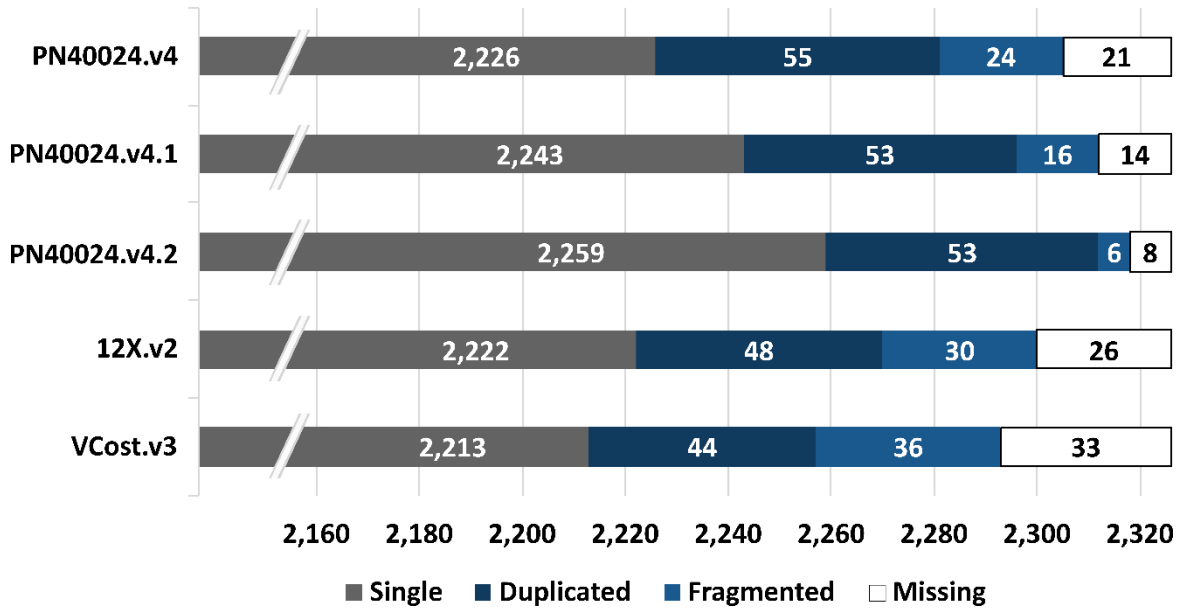


Figure 6  
154x79 mm ( x DPI)