A chromosome-level, haplotype-phased genome assembly for *Vanilla planifolia* highlights that partial endoreplication challenges accurate whole genome assembly

Q. Piet, G. Droc, W. Marande, G. Sarah, S. Bocs, C. Klopp, M. Bourge, S. Siljak-Yakovlev, O. Bouchez, C. Lopez-Roques, S. Lepers-Andrzejewski, L. Bourgois, J. Zucca, M. Dron, P. Besse, M. Grisoni, C. Jourda, C. Charron

PII: S2590-3462(22)00080-3

DOI: https://doi.org/10.1016/j.xplc.2022.100330

Reference: XPLC 100330

To appear in: PLANT COMMUNICATIONS

Received Date: 30 October 2021

Revised Date: 10 April 2022

Accepted Date: 27 April 2022

Please cite this article as: Piet, Q., Droc, G., Marande, W., Sarah, G., Bocs, S., Klopp, C., Bourge, M., Siljak-Yakovlev, S., Bouchez, O., Lopez-Roques, C., Lepers-Andrzejewski, S., Bourgois, L., Zucca, J., Dron, M., Besse, P., Grisoni, M., Jourda, C., Charron, C., A chromosome-level, haplotype-phased genome assembly for *Vanilla planifolia* highlights that partial endoreplication challenges accurate whole genome assembly, *PLANT COMMUNICATIONS* (2022), doi: https://doi.org/10.1016/j.xplc.2022.100330.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2022 The Author(s).



4	A abromogome lovel bonletune phaged general accomply for Varilla		
I	A chromosome-level, haplotype-phased genome assembly for vanua		
2	planifolia highlights that partial endoreplication challenges accurate whole		
3	genome assembly		
4			
5	Piet Q. ^{1,*} , Droc G. ^{2,3,4,*,#} , Marande W. ^{5,*} , Sarah G. ^{6,*} , Bocs S. ^{2,3,4,*} , Klopp C. ^{7,*} , Bourge M. ⁸ ,		
6	Siljak-Yakovlev S. ⁹ , Bouchez O. ¹⁰ , Lopez-Roques C. ¹⁰ , Lepers-Andrzejewski S. ¹¹ , Bourgois		
7	L. ¹² , Zucca J. ¹³ , Dron M. ¹⁴ , Besse P. ¹⁵ , Grisoni M. ^{16#} , Jourda C. ^{1,#} and Charron C. ^{1,*}		
8			
9	¹ CIRAD, UMR PVBMT, F-97410 Saint-Pierre, La Réunion, France. 📞		
10	² CIRAD, UMR AGAP Institut, F-34398 Montpellier, France.		
11	³ UMR AGAP Institut, Univ Montpellier, CIRAD, INRAE, Institut Agro, F-34398 Montpellier,		
12	France.		
13	⁴ French Institute of Bioinformatics (IFB) - South Green Bioinformatics Platform, Bioversity,		
14	CIRAD, INRAE, IRD, F-34398 Montpellier France.		
15	⁵ INRAE, CNRGV, Genotoul, 31326 Castanet-Tolosan, France.		

- 16 ⁶AGAP, Univ. Montpellier, CIRAD, INRAE, Montpellier SupAgro, Montpellier, France.
- 17 ⁷Plateforme Bioinformatique, Genotoul, BioinfoMics, UR875 Biométrie et Intelligence
- 18 Artificielle, INRAE, Castanet-Tolosan, France
- 19 ⁸Cytometry Facility, Imagerie-Gif, Université Paris-Saclay, CEA, CNRS, Institute for
- 20 Integrative Biology of the Cell (I2BC), 91198, Gif-sur-Yvette, France.
- 21 ⁹Université Paris-Saclay, Ecologie Systématique Evolution, Univ. Paris-Sud, CNRS,
- 22 AgroParisTech, Université Paris-Saclay, Orsay Cedex, France.
- 23 ¹⁰INRAE, GeT-PlaGe, Genotoul, 31326, Castanet-Tolosan, France.
- 24 ¹¹Etablissement Vanille de Tahiti, Uturoa, French Polynesia, France.
- 25 ¹²Eurovanille, Rue de Maresquel, 62870 Gouy Saint André, France.
- 26 ¹³V. Mane Fils, Département Biotechnologie, 06620 Le Bar Sur Loup, France.
- 27 ¹⁴Université Paris-Saclay, CNRS, INRAE, Univ Evry, Institute of Plant Sciences Paris-Saclay
- 28 (IPS2), 91405 Orsay, France.
- 29 ¹⁵Université de la Réunion, UMR PVBMT, Saint-Pierre, La Réunion, France.
- 30 ¹⁶CIRAD, UMR PVBMT, 501 Tamatave, Madagascar.
- 31
- 32 * These authors contributed equally to this work
- 33 [#]Correspondance : michel.grisoni@cirad.fr ; cyril.jourda@cirad.fr ; gaetan.droc@cirad.fr

34

35 SHORT SUMMARY

36

The genome of the orchid *Vanilla planifolia* (4.09 Gb, 16 pairs of chromosomes) is very prone to partial endoreplication (PE) which leads to a very unbalanced DNA content in the cells. We report here first molecular evidence of PE at chromosome scale through the assembly and annotation of an accurate haplotype-phased genome of *V. planifolia*. Distinct haplotypespecific sequencing depth variation patterns suggest complex molecular regulation of endoreplication along chromosomes. To facilitate post-genomics efforts, an integrated public and user-friendly web portal (the Vanilla Genome Hub) has been developed.

44 45

46 KEYWORDS

47

48 Vanilla, Whole genome sequencing, Optical mapping, Strict partial endoreplication, Genome49 Hub

50

51 ABSTRACT

52

53 Vanilla planifolia, the species cultivated to produce one of the world's most popular flavors, is 54 highly prone to partial genome endoreplication (PE) which leads to highly unbalanced DNA 55 content in cells. We report here first molecular evidence of PE at chromosome scale by the 56 assembly and annotation of an accurate haplotype-phased genome of V. planifolia. Cytogenetic 57 data demonstrated that the diploid genome size is 4.09 Gb, with 16 chromosome pairs although 58 aneuploid cells are frequently observed. Using PacBio HiFi and optical mapping, we assembled 59 and phased a diploid genome of 3.4 Gb with a scaffold N50 of 1.2 Mb and 59,128 predicted protein-coding genes. The atypical k-mers frequencies and the uneven sequencing depth 60 observed agreed with our expectation of unbalanced genome representation. Sixty-seven 61 percent of the genes were scattered over only 30% of the genome, putatively linking gene-rich 62 regions and the endoreplication phenomenon. On the contrary, low coverage regions (non-63 endoreplicated) were rich in repeated elements but also contained 33% of the annotated genes. 64 65 Furthermore, this assembly showed distinct haplotype-specific sequencing depth variation 66 patterns suggesting a complex molecular regulation of endoreplication along the chromosomes. This high-quality anchored assembly represented 83% of the estimated V. planifolia genome. 67

It provides a significant step towards the elucidation of this complex genome. To support postgenomics efforts, we developed the Vanilla Genome Hub, a user-friendly integrated web portal that allows centralized access to high-throughput genomic and other omics data, and interoperable use of bioinformatics tools.

- 72
- 73

74 INTRODUCTION

75

Endoreplication, characterized by a series of DNA replications in the nucleus without mitotic 76 77 cell division, is found in a large number of both animal and plant species (Lee et al., 2009). 78 During regular endoreplication, each step of this mechanism leads to a two-fold nuclear DNA 79 content increase in somatic cells (2C to 4C, 8C, 16C, etc.), where 1C corresponds to the DNA 80 content of the non-replicated holoploid chromosome set. Endoreplication is very common in 81 plants and related to various biological processes, such as plant development and growth, and 82 occurs in response to biotic and abiotic stresses (Bourdon et al., 2012; Lang and Schnittger, 2020). This phenomenon depends on the type of tissue and its stage of development, suggesting 83 84 involvement in cell differentiation and maintenance of the final stage of differentiation (Bhosale et al., 2018). The molecular mechanisms involved in regular endoreplication have 85 86 been particularly studied in Arabidopsis thaliana over the past few years. A downregulation of 87 mitotic activity caused by mitotic Cyclin-dependent kinase (CDK)-cyclin complexes has been 88 shown to be directly involved in the control of endoreplication (Lang and Schnittger, 2020).

89

90 In many orchid species, measurements of genomic content by flow cytometry (FCM) did not 91 agree with the commonly accepted model of complete endoreplication. In this case, nuclear 92 DNA content in endoreplicated cells was present at less than two-fold the 2C cells content. 93 Since this ratio was constant, whatever the cell ploidy level, for a given Vanilla species, it was 94 suggested that the nuclear DNA could be categorized into two parts: The P fraction, subject to 95 endoreplication, and the F fraction, not endoreplicated (Brown et al., 2017). These fractions 96 are constant in all cells undergoing PE, which suggests fine regulations of genome 97 rearrangements. The fact that the gametes are haploid also suggests the presence of molecular 98 mechanisms allowing the isolation of the holoploid genome. This type of endoreplication, that 99 appears to be specific to the Orchidaceae lineage in plants has been successively termed 100 "progressively partial endoreplication (PPE) (Bory et al. 2008; Trávníček et al., 2015; Hřibová 101 et al., 2016), strict partial endoreplication (SPE) (Brown et al., 2017), and more recently

"partial endoreplication (PE) (Chumová et. al, 2021; Trávníček et al., 2021). To be in line with
the latest works and in order to harmonize the terminology of this phenomenon, the term PE
will be used in this work. So far, PE has been observed in all species studied within the genus *Vanilla* (Bory et al., 2008; Brown et al., 2017; Lepers-Andrzejewski et al., 2011; Trávníček et
al., 2015).

107

108 Vanilla planifolia G. Jackson is an emblematic orchid cultivated for fruit (pod) fragrance. Pods 109 contain many aromatic compounds, particularly vanillin in high proportion (Perez-Silva et al., 110 2006). In this species, diploid nuclei (2C) are mainly found in nodal tissues (with PE up to 32E) 111 while the nuclei of mature leaf cells contain low 2C fraction and PE up to 64E (Brown et al., 112 2017). The F fraction was estimated to be 71.6% of the genome while the P fraction (28.4%) 113 could be duplicated up to 64E. In addition, the proportion of the non-endoreplicated (F) genome 114 varies greatly from species to species. It is very high in V. pompona (F=81%) but rather low in V. mexicana (F=17%) (Brown et al., 2017). Several studies on orchids have also shown that 115 116 species prone to PE have a larger genome than those prone to conventional endoreplication (Chumová et al., 2021; Trávníček et al., 2019; Trávníček et al., 2015). Nevertheless, the 117 118 molecular mechanisms involved in PE are not yet elucidated.

119

120 A chromosome-scaled, phased V. planifolia genome (Daphna cultivar) was recently reported, 121 highlighting haplotype difference and one ancestral whole genome duplication shared by all 122 sequenced orchids (Hasing et al., 2020). However, the 1.5 Gigabase (Gb) of the assembled 123 genome was far from the V. planifolia genome size estimated to be about 4 Gb using FCM 124 measurement (Bory et al., 2008; Lepers-Andrzejewski et al., 2011), which suggests a putative 125 high incomplete genome assembly concerning Daphna. As mentioned by Hasing et al. (2020), 126 the reason for the genome size discrepancy between flow cytometry and assembly results 127 remained to be elucidated. With about 65% of the V. planifolia genome missing in Daphna assembly, we hypothetize that the missing part of the genome corresponds mainly to the F 128 (non-endoreplicated) fraction (71.6%) (Brown et al., 2017), due to its lower representation and 129 130 therefore its lower sequencing depth.

131

Here we develop an approach combining FCM, cytogenetics and whole genome sequencing using most recently developed technologies (Supplemental Figure 1) and a 2C fraction enriched tissue (nodes), resulting in reduced P/F therefore greater proportion of the F fraction to unlock this issue. We demonstrate that size discrepancy was due to the occurrence of PE, for

which further knowledge at chromosome scale was gained from this study. We present the most complete version to date of a high-quality chromosome-level phased genome of *V*. *planifolia* using a traditional vanilla cultivar from the Indian Ocean region (CR0040). Our results are shared through a web portal facilitating data access, use and analyses by a wide community.

141

142

143 **RESULTS**

144

145 Genome size, ploidy level and chromosome content

146

147 The 2C genome size of V. planifolia CR0040, a traditional vanilla cultivar from La Reunion 148 island (Supplemental Note 1), was estimated in nodal tissues to be 4.18 ± 0.08 pg by FCM (Supplemental Note 1), corresponding to 4.09 Gb (Doležel et al., 2003). To estimate PE levels, 149 150 fluorescence ratio of DNA content between consecutive peaks of endoreplication levels was 151 estimated (Supplemental Note 1; Supplemental Figure 2 and Supplemental Table 1). Results 152 showed no significant difference (calculated t-value of 1.116, 1.900, 0.935, 0.365 compared to Student table t-value (a=0.05) of 2.131) between PE pattern of CR0040 and other V. planifolia 153 154 cultivars such as CR1110 ($2C = 4.16 \pm 0.04 \text{ pg}$) studied by Brown et al. (2017). Replicated fraction P was also calculated (P = $30.5\% \pm 3.2$). The equivalent amount 2p was then P x 2C = 155 156 1.275 pg, which meant that the absolute quantity p was 0.637 pg and the absolute quantity f of 157 fixed amount was 1.453 pg (Figures 1A, 1B; Supplemental Note 1). The karyotype of V. 158 planifolia obtained by cytogenetics approaches (Supplemental Note 1) appeared as of bimodal 159 type, namely composed by 16 both large and small chromosome pairs (Figures 2A, 2B, 2C), 160 although aneuploid cells were frequently observed, such as those with only 28 chromosomes (Figures 2D, 2E). V. planifolia chromosomes possess important portions of telomeric and 161 162 pericentromeric heterochromatin which made the determination of their morphology difficult. 163 In the interphase nuclei, this heterochromatin was present in the form of numerous 164 chromocenters clearly visible both after staining with orcein (Figure 2F) and DAPI (Figure 165 2G). This type of heterochromatin is unspecific while the heterochromatin linked to rRNA 166 genes is rich in G-C bases. Only one locus (two spots) of rDNA (18S-5.8S-26S) was present 167 in the genome of V. planifolia (Figure 2H, arrows), evidenced after chromomycin (CMA3) staining. Our results also revealed, after Hoechst 33258 staining, that in the V. planifolia 168

169 chromosomes AT-rich DNA regions were more represented than those GC-rich and that some170 chromosomes were entirely or almost entirely heterochromatinized (Figure 2I).

- 171
- 172

173 Whole genome assembly and k-mers analysis

174

175 CR0040 genome sequencing produced 69 Gb Pacific Biosciences (PacBio) HiFi long-reads, 176 147 Gb Oxford Nanopore Technology (ONT) long-reads and 200 Gb Illumina 10X Genomics 177 short-reads (Supplemental Note 1 and Supplemental Table 2). These DNAseq reads were 178 assembled using different bioinformatics pipelines. The best result was obtained using only 179 high quality HiFi long-reads (Supplemental Note 2 and Supplemental Tables 3 and 4). Contigs 180 from HiFi reads assembly were scaffolded with optical maps to obtain a final phased assembly 181 of 3.4 Gb (1.5 Gb for haplotype A and 1.9 Gb for haplotype B) representing around 83% of the 182 expected genome size. One third of the assembly could be anchored onto 14 chromosomes 183 using published Daphna chromosomes as references (Hasing et al., 2020). Unfortunately, no data could help to organize the remaining contigs into the two missing chromosomes. 184 185 Therefore, the remaining two thirds correspond to unanchored additional sequences that were compiled into two unknown random pseudomolecules A0 and B0. The final assembly 186 187 comprised 24,534 contigs with a contig N50 length of 924 kb. The lengths of the 14 188 chromosomes ranged from 73.5 Mb (Chr01) to 20 Mb (Chr14). Main genome assembly 189 statistics are synthesized in Table 1.

190

191 In order to understand how the PE affects the assembly, a k-mer analysis was produced. The 192 results should reflect the sequencing coverage of the different genome fractions present in our 193 raw data and assembly. Briefly, the reads were split into overlapping k-mers (47 mers in our 194 case), then the k-mers were sorted and the occurrences counted. These counts were then used 195 to produce a histogram. A spectra-cn plot was used to compare the k-mers found in the reads 196 versus the k-mers found in the assembly (Supplemental Figure 3). The X-axis gives the number 197 of times a given k-mer was found in all the reads, reflecting the coverage of the k-mer. The Y-198 axis gives a value representing the number of k-mers that were found a specified number of 199 times (X-axis value). Interestingly, two k-mers distributions were centered at 42X and 84X, 200 representing classical diploid distribution with the heterozygous and homozygous k-mer 201 content. We assumed that these peaks represented the k-mers of the endoreplicated fraction 202 with high sequencing depth due to higher representation. Remarkably, the graph also showed

an additional k-mers distribution centered around 10X (Supplemental Figure 3, red arrow).
This distribution can easily be mistaken with the erroneous k-mers distribution but we assumed
that it represented non-endoreplicated k-mers of the *V. planifolia* genome, with low-sequencing
depth due to lower representation.

207

208 To validate the assembly and compare it with the already published reference, we produced 209 four k-mers spectra-cn plots showing Daphna Illumina reads and CR0040 HiFi reads k-mer 210 distributions colored both with Daphna and CR0040 assemblies (Figure 3). A spectra-cn plot 211 enables to compare the k-mers found in the reads versus the k-mers found in the assembly. The 212 k-mer histogram from reads is colored by the number of times each k-mer is found in the 213 assembly. For a heterozygous diploid assembly, we expected to find two distributions, on the 214 left hand the heterozygous distribution which should be colored in red because each k-mer is 215 only found once in the assembly and on the right hand the homozygous distribution which is purple because the corresponding k-mers are found twice in the assembly. The black area at 216 217 the far left of the diagram corresponds to k-mers including sequencing errors which are found 218 a limited number of times in the reads and never in the assembly. Daphna Illumina sequencing 219 being deeper resulted in a better separation between the homozygous (80X sequencing depth) 220 and heterozygous (160X depth) k-mers fractions in the spectra-cn graph, compared to CR0040 221 (Figures 3A and 3B against 3C and 3D). The same pattern occurred for CR0040 HiFi data 222 around 45X and 90X (Figure 3C). The differences between Figures 3A and 3C come from the 223 sequencing depth and the type of tissue used: mature leaves with a higher proportion of the P 224 fraction for Daphna and nodal tissues with a lower P/F ratio for CR0040. The k-mer distribution 225 of the non-endoreplicated fraction (low coverage) were not found in the Daphna assembly 226 (black area left of Figures 3B and 3D) but are mostly present in the CR0040 assembly. 227 Regarding the completeness of the Daphna reference assembly, the spectra-cn plots (Figures 228 3B and 3D) showed that part of the heterozygous fraction was missing (orange arrows) and 229 some k-mers were in over-represented copies (>2X) in both heterozygous and homozygous fractions (Figure 3B, black arrows). The spectra-cn diagram also showed heterozygous content 230 231 present twice or more times instead of once in this assembly (Figure 3, black arrows) which 232 could indicate spurious duplications. As a whole, our genome assembly of CR0040 is close in 233 genome size to FCM estimation and has the expected k-mer diploid profile with a well-234 represented non-endoreplicated fraction (Figures 1C and 1D).

- 235
- 236

238

237 Genes and transposable elements annotation

- 239 The assembled genome supplemented with transcriptomic data from nine distinct tissues made 240 it possible to identify 59,128 protein coding genes (26,392 for haplotype A and 32,736 for 241 haplotype B) among which 90.31% could be associated with a function (Supplemental Note 3, 242 Supplemental Tables 5 to 10). Sixty-seven percent of predicted genes were anchored onto the 243 14 chromosome pairs, and the remaining 33% onto the two random mosaic chromosomes that 244 were constructed from the unanchored scaffolds and contigs (Figure 4A, blue distributions). 245 We estimated the annotation completeness at 93.2% with the BUSCO (Benchmarking 246 Universal Single-Copy Orthologs) approach using the Viridiplantae database. In total, 72% of 247 the assembly consisted of repeats including single sequence repeats (SSR, 15.4%) and 9.7% 248 other low complexity regions (Supplemental Note 3, Supplemental Table 11). High content of 249 retrotransposons was found (41.5%), while content in DNA transposons was low (1.4%). The 250 Long Terminal Repeats (LTR) retrotransposon content was richer in Gypsy (9.7%, Figure 4A, 251 purple distributions) than in Copia (6.1%, Figure 4A, orange distribution), although a number 252 of annotated retrotransposons (12.5%) were not more precisely classified. The two random 253 mosaic chromosomes were enriched in repeats and showed low gene density and low 254 sequencing depth (Figure 4A, green distributions and Supplemental Table 12). Indeed, 255 compared to the 14 chromosome sequences, the unanchored regions showed higher proportions of Long Interspersed Nuclear Elements (LINEs) sequences (8% and 14.05%, respectively), and 256 257 this for both haplotypes. On the contrary, DNA transposons (3.14% and 0.93%) as well as 258 Short Interspersed Nuclear Elements (SINEs, 0.12% and 0.05%) and LTRs (21.67% and 259 15.57%) represented a larger part of the 14 chromosome sequences than of unanchored regions. 260 The biggest difference in favor of unanchored regions was observed for unclassified 261 retrotransposons which represents 16.76% of unanchored sequences against 3.28% of the 14 262 chromosome sequences. Main genome annotation statistics are synthesized in Table 1.
- 263

264

265 *V. planifolia* pangenomics and whole genome duplication

266

The comparison of the four mosaic haplotypes from the two *V. planifolia* cultivars, CR0040 and Daphna (Supplemental Tables 12 and 13) showed that the 14 pseudomolecules of CR0040 were shorter and contained less genes than Daphna and that a large number of regions in the CR0040 pseudomolecules (haplotype A or B) were not located in the Daphna pseudomolecules

271 (Supplemental Figure 4). Pan-genomic analysis of the orthogroups from proteomes derived 272 from the 14 chromosomes only (Supplemental Figure 5; Supplemental Table 14; Supplemental 273 Note 4) indicated that the core genome was composed of 14,210 families and 77,692 genes 274 (35,972 CR0040 and 41,720 Daphna). The dispensable genome of CR0040 contains 1,266 275 families and 3,613 genes specific to CR0040. The dispensable genome of Daphna contains 276 3,997 Daphna specific families and 13,645 genes. Finally, we looked at the expansion or 277 reduction of gene families in relation to six proteomes (CR0040, Daphna, Phalaenopsis 278 equestris, Phalaenopsis aphrodite, A. thaliana, Oryza sativa; Supplemental Figure 6). From an 279 orchid perspective, the expansion number for the orchid node is rather low (+36) while the 280 Daphna specific number is rather high (A +1841 and B +1943) compared to CR0040 (A +418 281 and B +826).

To identify whole genome duplications (WGD), pairwise genome synteny analyses between CR0040, Daphna, and *P. aphrodite*, and between themselves were carried out (Supplemental Figures 7 and 8, and Supplemental Notes 4 and 6). The CR0040 haplotype A dotplot validated at least one pan-orchid WGD (α° , the origin of the paleo-allotetraploid) previously found by Hasing et al. (2020). An additional dotplot diagonals and dS peak suggested a second WGD, possibly the tau of Monocots (τ^{m}).

- 288 289
- 290 Detection of non-endoreplicated regions

291

292 PE induce highly unbalanced DNA representation with a P/F DNA ratio ranging from 3 to 10 293 according to tissues. This was reflected in our assembly by highly variable sequencing depth 294 (Figure 4A, green lane). The two random mosaic chromosomes which showed a low 295 sequencing depth at most loci, may therefore contain a large part of the non-endoreplicated F 296 fraction of the genome. It is likely that a large part of unanchored sequences originates from 297 the two fully heterochromatinized chromosomes observed in the interphase nuclei (Figure 2I), 298 possibly chromosome pairs 15 or 16). The remaining part of unanchored sequences should 299 correspond to missing fractions in the anchored chromosomes. Interestingly, the sequencing 300 reads mapping on CR0040 and Daphna assemblies also showed intra-chromosomal sequencing 301 depth variations (Figure 4B). These patterns were consistent whatever the technology used. To 302 observe this phenomenon globally on all chromosomes and genomes with all technologies 303 (HiFi, ONT, Illumina), sequencing depth analysis tools were used and manual validation 304 performed (Supplemental Note 5; Supplemental Tables 15 and 16). Two patterns of sequencing

305 depth variations have been identified along all chromosomes. The first one (indicated with a dotted box labeled "1" in Figure 4B and Supplemental Figure 9) corresponded to a sharp 306 307 decrease of sequencing depth for both cultivars, with all sequencing technologies which 308 dropped down from 45X-120X to less than 20X. Surprisingly, this pattern occurred 309 independently on the two haplotypes. A total of 37 very low coverage regions (from 0.4 to 6 310 Mb in length) with this pattern were identified along the chromosomes (24 in haplotype A and 311 13 in haplotype B) for a cumulative size of 60.1 Mb. For a large part of these regions, we found 312 low gene density and high repeat density. This pattern could correspond to non-endoreplicated 313 regions present in both Daphna and CR0040 genomes. The fact that these patterns are 314 systematically located on junctions between super-scaffolds is consistent with the decrease in 315 sequencing depth inherent from non-endoreplication which impaired the assembly of the 316 endoreplicated regions located on either side. The second pattern (indicated with a dotted box 317 labeled "2" in Figure 4B and Supplemental Figure 9) corresponded to 36 regions (from 1.2 Mb to 20 Mb in length, cumulative size of 207.2 Mb) with segmental sequencing depth variation 318 319 present in CR0040 (with HiFi, ONT and Illumina) but not in Daphna (with ONT and Illumina). 320 Furthermore, these variations were syntenic along the two haplotypes but the direction of 321 variation was inverted between the two phases. Their respective level of sequencing depth 322 differs by a factor of about 3 in CR0040. The cause of these apparently coordinated sequencing 323 depth inversions between CR0040 haplotypes remains unclear. After analyzing the locations 324 of these k-mers in CR0040 assembly, it appeared that these low depth k-mers (between 5X and 15X) were mostly present in the unanchored part of the genome comparatively to the 325 326 chromosome sequences (Figure 5) with median ratios values equal to 0.27 and 0.036327 respectively, showing significant difference (Wilcoxon-Mann-Whitney test; p-value = 4e-13). 328 However, chromosomes 7A and 6B were outliers, showing also a high proportion of low depth 329 k-mers. In addition, the distribution of these k-mers along the genome was globally consistent 330 with the areas identified except for some discrepancy (Supplemental Figure 10). Chromosomes 331 6B and 7A showed strong signals in terms of low depth k-mers proportions as already pointed out by Figure 5. Indeed, on chromosome 6B, k-mers of this type were positioned on nearly all 332 333 the assembled sequence while they were localized on approximately half of the chromosome 334 assembled 7A sequence.

- 335
- 336

337 Cell cycle regulator genes orthologs involved in A. thaliana endoreplication

338

339 The search for orthologs of the Cyclin-dependent Kinases (CDK) and cyclins (Cyc) families 340 of A. thaliana, involved in regular endoreplication mechanism, showed that representatives of 341 these two families were indeed found in the proteomes of CR0040 and P. aphrodite 342 (Supplemental Table 17). However, the number of genes coding for CDKs and cyclins found 343 via the orthogroups approach were lower for these two species. For example, the gene coding for Cyc-D3-1 in A. thaliana (At4g34160) was part of a species-specific orthogroup containing 344 345 some other D-types cyclins genes. Regarding the genes encoding the regulatory proteins of the 346 CDK/Cyc complexes, all of them had orthologs in both orchids. However, it appeared that 347 these multigenic families were slightly underrepresented in CR0040 gene annotation compared 348 to A. thaliana and P. aphrodite. Finally, an imbalance between the A and B haplotypes was 349 observed for Fizzy-related proteins (Fzr) and cyclin-dependent kinase inhibitor (Krp) 350 orthologs.

- 351
- 352

353 Vanilla genome hub

354

To support post-genomics efforts, the Vanilla Genome Hub (VGH, https://vanilla-genome-355 hub.cirad.fr) has been developed. It centralizes vanilla genomic information with a set of user-356 friendly interconnected modules and interfaces for analyzing and visualizing genomic data. 357 358 From the main menu of VGH (Supplemental Note 6 and Figure 6A), the search for genes of 359 interest to biologists is simplified using the interoperable system by the identification of 360 paralogous genes by keywords and sequence homology (Figures 6B and 6C) and information 361 report with gene name, gene localization and polypeptide function (Figure 6D). The genome 362 browser was built to offer tracks of supplemental information such as: GC content, gene 363 structure, gene expression, DNA sequencing depth and composition in repeats to support the identification of new genes of interest (Figure 6E, Supplemental Figure 11). A metabolic 364 365 pathways reconstruction and visualization tool allows to identify the annotated genes involved 366 in pathways (Figure 6F). A GO enrichment tool allows to test and visualize the enrichment 367 according to a GO category of a gene group (Figure 6G). At last, comparative analysis at 368 genome-scale is supported by an interactive multiscale synteny visualization (Figure 6H).

- 369
- 370
- 371 DISCUSSION
- 372

373 Flow cytometry and cytogenetic data validate the genome size and chromosome content 374 375 Genome size, ploidy level and chromosome content of the V. planifolia CR0040 cultivar were 376 validated by FCM and cytogenetic analyses. The estimated size at 4.09 Gb indicated a ploidy 377 level similar to other traditional diploid V. planifolia cultivars (Bory et al., 2008; Lepers-378 Andrzejewski et al., 2011). Estimation of endoreplication levels confirmed PE as previously 379 described in V. planifolia (Brown et al., 2017). This species was shown to harbor diploidized 380 meiotic chromosome pairing with 16 bivalents (Bory, 2007). This demonstrates the complete 381 diploidization of this supposed segmental paleo-allotetraploid (Nair and Ravindran, 1994; 382 Ravindran, 1979). The same meiotic observation has been also performed for $V_{\cdot} \times tahitensis$ 383 by Lepers-Andrzejewski et al. (2011). Aneuploid chromosome numbers were frequently 384 observed in mitotic metaphases of V. planifolia (Bory et al., 2008; Nair and Ravindran, 1994), 385 possibly due to the observed mitotic associations which could lead to unequal anaphase 386 separation. This may lead to errors in the evaluation of the basic chromosome number, as it is 387 the case in a recent paper where the authors considered that the basic number was x=14 (Hasing 388 et al., 2020). The phenomenon of an euploidy apparently occurs only in somatic cells while 389 meiosis appears to be regular with a stable number of chromosomes (Bory, 2007). Although the CR0040 assembly is more complete than that of Daphna, only 14 pseudomolecules were 390 391 obtained because CR0040 scaffolds were anchored on the 14 Daphna pseudomolecules. 392 Chromosomes 15 and 16 are probably not endoreplicated and present in the unanchored part 393 (CR0040 A0 and CR0040 B0).

394

395

396 PE hinders whole genome assembly

397

398 Given the CR0040 diploid genome size estimated to 4.09 Gb using FCM, our genome assembly 399 represented around 83% of the expected genome size which was twice the size of the Daphna 400 genome previously published (Hasing et al., 2020). This difference could be explained by the 401 fact that the assembly for CR0040 was done from HiFi reads which allowed to assemble in 402 different contigs the repeated regions despite their low sequencing depth. We were thus able to 403 assemble a greater number of repeated sequences that might correspond to a large fraction of 404 the non-endoreplicated genome, missing from the Daphna genome assembly. The biological 405 reality of this hypothesis is reinforced by the consistency of k-mer depth profiles and 406 sequencing depth patterns resulting from the mapping of reads from different sequencing

407 technologies (HiFi, 10X and ONT) tested in this study for CR0040 (Supplemental Figure 9). A k-mer spectra-cn diagram is an efficient tool to visually compare reads and assembly k-mer 408 409 compositions. They are used to validate diploid or haploid assembly quality (Yen et al., 2020). 410 The k-mer spectra-cn diagram shows clearly a diploid general pattern with a heterozygous 411 distribution containing only k-mers in single copy in the assembly and a homozygous distribution harboring, as expected, only k-mers present twice in the assembly. Unexpectedly 412 413 for a diploid genome, this figure includes a third distribution which is located in the low coverage area of the diagram. The color pattern shows clearly that these k-mers present in low 414 415 frequencies (5 to 15 times) are also present in our assembly. These k-mers represent the non-416 repeated fraction of low coverage sections of the assembly which are mainly located in the 417 unanchored sequences but are also present in low coverage sections of the other chromosomes. 418 Even if the unanchored sequences are mainly built of repeats they also harbor genes and other 419 non-repeated blocks and these parts are large enough in terms of k-mers to generate this unexpected k-mer distribution in the spectra-cn plot. These k-mers are not present in the public 420 421 V. planifolia Daphna assembly and therefore the corresponding distribution is black in Figure 422 3B.

- 423
- 424

425 Molecular signatures of partial endoreplication

426

The abundance of interspersed repeats detected in CR0040 was consistent with already 427 428 mentioned data in other orchids such as *Phalaenopsis equestris* (Cai et al., 2015) and *P*. 429 aphrodite (Chao et al., 2018) and in other lineages, like the Oryza genus (Stein et al., 2018). 430 High content of retrotransposons and low content in DNA transposons were in the range of 431 what has been found for different orchids (Cai et al., 2015; Chao et al., 2018). High repeat 432 content was found in candidate non-endoreplicated regions, which is in agreement with what 433 has already been described in other orchids (Chumová et al., 2021). Furthermore, some types of repeats might be preferentially found in non-endoreplicated regions as shown by repeats 434 435 proportions differences between assembled chromosomes and unanchored sequences, and in 436 particular retrotransposons proportions. Thus, LINEs for example, occupy a larger portion of 437 the unanchored regions than that of the 14 chromosomes even though these regions are 438 overrepresented in the present assembly. However, the lack of a more detailed annotation of 439 the retrotransposons class hampers the search for a potential preferential distribution of repeat 440 families between endoreplicated and non-endoreplicated regions. It is therefore crucial to better annotate these repeats in order to determine exactly which kinds are preferentially found in the
two fractions of the genome. On the other hand, the distribution of genes in the genome shows
the opposite trend with approximately two-third of the protein coding sequences localized in
the anchored region.

445

The distinct sequencing depth profiles observed between CR0040 and Daphna likely reflected 446 447 a tissue-specific endoreplication pattern. Indeed, the nodes used for sequencing CR0040 genome are growing and differentiating tissues while the leaves used for sequencing Daphna 448 449 genome are made of fully differentiated cells. The irregular haplotype-specific endoreplication 450 pattern (segmental or not) observed in CR0040 could thus result from a peculiar physiological 451 activity. Whatever the reason, this intriguing pattern suggests a complex and fine regulation of 452 PE at chromosome level which deserves further studies. While no previous study has 453 demonstrated the mechanisms underlying PE in orchids, many works have focused on the 454 regulation of regular endoreplication found in a large number of plant species, and well 455 analyzed on tomato and Arabidopsis (Lang and Schnittger, 2020). The common mechanism to trigger endoreplication is a downregulation of mitotic CDK activity to suppress the mitosis and 456 457 a fine regulation of this activity throughout the induced endocycle with an alternance between 458 high and low activity levels at specific checkpoints in order to maintain the replication process (De Veylder et al., 2011; Shimotohno et al., 2021). CDK controls cell cycle progression and 459 460 mitosis entry via its phosphorylation activity, which is activated by association with CYC 461 proteins. Recently, Inada et al. (2021) demonstrated the involvement of actin and actin-binding 462 protein in the regulation of A. thaliana endoreplication. The whole genome analysis made it possible to identify orthologous CDK, cyclins, CDK-activators/repressors and Actin 463 464 Depolymerizing Factors (ADF) in V. planifolia CR0040. A first step in understanding orchid 465 PE would therefore be to further analyze these molecular regulators. Indeed, the recognition of 466 orthologs and paralogs in large gene families, such as CDK-cyclin complex, is challenging and 467 requires deepening by a high-quality manual annotation of the genes of interest (Vaattovaara et al., 2019). 468

Even though PE seems specific to orchids in plants (Trávníček et al., 2015), such phenomenon
of under-represented genomic regions is well known in metazoan. Ciliates such as *Paramecium tetraurelia* or *Tetrahymena thermophile* show programmed DNA elimination following
endoreplication in their MAC nucleus, involving chromosome fragmentation and elimination
of specific sequences called IES (Internal Eliminated Sequences) (Sellis, 2021; Bracht et al.,
2013). However, FCM approaches on *Ludisia discolor*, an orchid subject to PE, have ruled out

the hypothesis of such DNA elimination and favor that of under-replication (Hřibová et al., 2016). Under-replication has also been studied in several organisms such as Drosophila for which it has been proposed that a reduction in expression of genes involved in DNA replication may lead to a slower mitosis S-phase, and an incomplete replication of genomic regions during late S-phase (Lilly and Spradling, 1996). Molecular mechanisms described in Drosophila highlighted an inhibition of the replication fork progression involving Rif1 protein, which interacts with the SUUR protein (Munden et al., 2018; Armstrong et al., 2019).

482

483 Finally, cytogenetic studies using in situ hybridization techniques (FISH, fluorescence in situ 484 hybridization, and GISH, genomic in situ hybridization) could also be used to increase our 485 knowledge of the molecular signatures of PE (Younis et al., 2015). Recent advance in FISH is 486 the development of probes based on synthetic oligonucleotides specific to repetitive sequences 487 or to particular chromosome regions (Jiang, 2019). This new-generation of FISH probes in plants was applied to species with a sequenced genome, for example for Zea and Cucumis 488 489 species (Braz et al., 2020; Han et al., 2015; Martins et al., 2019; Zhang et al., 2021a). 490 Endoreplicated versus non-endoreplicated genomic regions could serve to synthesize oligo-491 based FISH probes specific to each fraction for precisely locating these PE signatures on 492 chromosomes. GISH technique uses the total genomic DNA of a species, in contrast to FISH. We hypothesize that hybridizing the total DNA of highly endoreplicated nuclei (16E, 32E) to 493 494 CR0040 chromosomes would induce a more intense hybridization signal in endoreplicated 495 regions, which would allow us to identify non-endoreplicated areas with little hybridization.

- 496
- 497

498 Impact of technologies on whole genome evolution analysis

499

500 The strategy of combining optical mapping and HiFi long reads sequencing for CR0040 501 genome assembly resulted in a haplotype A with 14 pseudomolecules of better quality and with less scaffolding errors than Daphna haplotype A, which was built with Hi-C and ONT 502 503 technologies (Hasing et al., 2020). Indeed, comparisons between Daphna and CR0040 504 haplotypes A revealed a dual haplotype conservation problem in the Daphna phased assembly, 505 which is reflected in the Daphna Hi-C scaffolding. The use of HiFi long-reads and optical maps 506 allowed for more accurate haplotype separation as shown in previous works (Du et al., 2020; 507 Matthews et al., 2018). In the case of CR0040 not only did HiFi allow better assembly of nonendoreplicated regions, but HiFiasm allowed better separation of haplotypes. These 508

509 improvements were therefore necessary to better solve the sequencing of the complex vanilla 510 genome, with a high rate of heterozygosity (Ho V. planifolia cultivars = 0.362; Favre et al., 511 2022) and subjected to PE (Brown et al., 2017). However, this dual haplotype conservation 512 problem observed in Daphna, and not in CR0040, impacted comparative pan-genomics 513 analyses and distorted results obtained. Thus, the differences observed between the two genomes of V. planifolia (number of paralogs, numbers of gene families with 514 515 expansions/contractions, complete and duplicated BUSCO scores) could be explained by these 516 mosaic assembly problems and therefore an incorrect separation between haplotypes A and B. 517 Monocot genome evolution analyzes were carried out using the high-quality haplotype A 518 sequence of CR0040. Dotplot results were in agreement with the fact that V. planifolia is a 519 diploidized paleo-polyploid species with primary basic chromosome number x=8 and 520 secondary basic number x=16, as described for the whole Vanilla genus (Felix and Guerra, 521 2005). Moreover, only one locus (two spots) of rDNA (18S-5.8S-26S) was identified in the 522 genome of V. planifolia by cytogenetic approaches, which brings one additional evidence of 523 an ancient diploidization of this supposed segmental paleo-allotetraploid. Finally, two WGDs, possibly corresponding to α° and τ^{m} , were highlighted, as also described for the *Dendrobium* 524 525 chrysotoxum chromosome-scale genome assembly (Zhang et al., 2021b).

526 527

528 Efficiency of the integrative approach combining cytogenetics and high-quality whole 529 genome sequencing

530

531 In this study, we confirmed the size and structure of the V. planifolia genome using both cytogenetics and nuclear DNA sequencing methods. The particular phenomenon of PE at play 532 533 in many orchids has been explored at chromosome level, for the first time in plants to our 534 knowledge. Our data showed that the non-endoreplicated sequences are very predominantly 535 made up of repeated sequences. This confirmed, at the genomic level, previous findings in 536 orchids by Chumová et al. (2021), based on Phylogenetics Generalised Least Squares (PGLS) 537 model, and by Brown et al. (2017) who demonstrated in Vanilla by nuclei imaging that, on the 538 other hand, the endoreplicated part was transcribed. We nevertheless revealed that 33% of the 539 59,128 protein coding genes annotated were present in the two random mosaic chromosomes 540 corresponding mainly to the non-endoreplicated part as shown by low-sequencing depth. In 541 addition, a thorough examination of the sequencing depths along anchored chromosomes with 542 three different technologies has revealed 73 regions that appear with different endoreplication

543 levels that vary with haploid phase. Half of which may be linked to tissue type (leaves vs nodes). This last conclusion has to be confirmed with DNA sequencing from different tissues 544 545 of the same cultivar. This work constitutes a considerable progress in the knowledge of V. 546 *planifolia* genomics and sheds light on the most relevant methodologies for further deciphering 547 this complex genome and the PE phenomenon. The Vanilla Genome Hub was built to help the 548 community to address major unresolved questions about vanilla such as PE, biosynthesis of 549 aromatic compounds, and resistance to pathogens. We are working on a new version of the 550 vanilla nuclear genome sequence that will be improved in terms of haplotype separation, 551 chromosome reconstruction, gene and repeat element annotation in order to further investigate the molecular mechanisms of PE with appropriate plant material, bio-technologies and 552 553 bioinformatics tools.

554

555

556 MATERIALS AND METHODS

557

558 Cytometry, cytogenetics and DNA sequencing

559

A traditional vanilla cultivar (CR0040) from Reunion island was used in this study (Supplemental Note 1). FCM and cytogenetics studies were performed using protocols described in (Supplemental Note 1). High molecular weight (HMW) DNA and ultra-HMW DNA was extracted from node tissues and sequenced using PacBio HiFi, ONT and Illumina technologies (Supplemental Note 1). Optical genome maps were produced using Bionano Genomics® protocol and the Saphyr G1 System (Supplemental Note 1).

566

567

568 Genome assembly and analysis

569

570 HiFi reads were assembled into contigs using hifiasm 0.13 with default parameters (Cheng et 571 al., 2021). The hybrid scaffolding was performed between DNA contigs and optical genome 572 maps using the hybrid Scaffold pipeline of Bionano Genomics® with default parameters. These 573 scaffolds were phased into two haplotypes using in house scripts and the not-scaffolded contigs (https://github.com/dfguan/purge_dups). 574 were phased using purge dups Then, 575 pseudomolecules were reconstructed using alignments of the phased-assembly on Daphna 576 chromosomes (Hasing et al., 2020, Supplemental Note 2). The assembly quality was estimated

577 with quast 5.1.0 (Gurevich et al., 2013) and using the approach of Benchmarking Universal Single-Copy Orthologs (BUSCO, version 5.0.0) (Simao et al., 2015) (Supplemental Note 2). 578 579 The k-mers analysis was performed with kat 2.4.2 using the comp tool (Mapleson et al., 2017). 580 The plot script was slightly modified to project on the Y axis the number of distinct k-mers 581 multiplied by the k-mers multiplicity instead of just the number of distinct k-mers. In parallel, 582 the k-mers of size 47 having a depth between 5 and 15 were extracted within PacBio sequences 583 using Jellyfish 2.3.0 (Marcais and Kingsford, 2011). These k-mers were repositioned on our 584 reference using the tool 'query per sequence' 585 (https://github.com/gmarcais/Jellyfish/tree/master/examples/query_per_sequence) and the ratio of these k-mers was computed among each sequence of our genome. These sequences 586 587 were splitted between chromosomes and unanchored sequences and the repartition of the k-588 mer ratio was drawn using python seaborn library (https://seaborn.pydata.org/).

- 589
- 590

591 Structural and functional genome annotation

592

593 Automatic gene prediction was performed on CR0040 contigs with the EuGene Eukaryotic 594 Pipeline (EGNEP version 1.5) (Sallet et al., 2019) (Supplemental Note 3). Transcriptomic data 595 from CR0040 were produced using RNA sequencing of nine organs with Illumina technology 596 (Supplemental Note 3). In addition, gene expression profiles and putative novel isoforms were 597 identified with Stringtie v. 2.0.3 (Kim et al., 2019) (Supplemental Note 3). Transcriptomic data 598 from V. planifolia cultivars (CR0040, Daphna (ncbi Bioprojects: PRJNA668740 and 599 PRJNA633886), and an unspecified cultivar (ncbi GEO: GSE134155)), proteomic data from V. planifolia Daphna (Hasing et al., 2020), Phalaenopsis equestris (ncbi Bioproject: 600 601 PRJNA382149) and the Liliopsida class (Swissprot 2020_06), as well as a custom orchids 602 specific statistical model for splice site detection were used for this analysis (Supplemental 603 Note 3). Functions were assigned through InterProScan domain searches as well as similarity 604 searches against Uniprot/Swissprot and Uniprot/TrEMBL databases (BlastP). Gene Ontology 605 (GO) terms were assigned through InterProScan (Jones et al., 2014) results while Enzyme 606 Classification (EC) numbers were predicted combining both tools PRIAM (Claudel-Renard et 607 al., 2003) and BlastKOALA (Kanehisa et al., 2016).

Repeats were first identified using RepeatModeler v2.0.1 (Flynn et al., 2020), RepeatScout
v1.0.5 and transposable element genes predicted from EGNEP annotation and then classified
with REPET v3.0 (Flutre et al., 2011) and PASTEC v2.0 (Hoede et al., 2014) according to the

611 Wicker's TE classification (Wicker et al., 2007). After cleaning steps (see details in note S9.2),

repeats were clustered with CD-HIT v4.8.1 (Fu et al., 2012) to produce two banks of repeats.

613 The CR0040 genome was then annotated for repeats using previous banks, RepeatMasker

614 v4.1.1 (Tarailo-Graovac and Chen, 2009) and bedtools intersect v2.29.2 (Quinlan and Hall,

615 2010).

- 616
- 617

618 Genomic comparisons and gene family's reconstruction

619

In order to compare the 14 haplotype A chromosomes of both vanilla cultivars, to check the 620 621 completeness of the Vanilla genome and to study the pan-orchid α° WGD, a series of analyses 622 were performed with the CoGe Synmap pipeline as described in Supplemental Note 4. Gene 623 family reconstruction was performed using Orthofinder2 (v.2.4.0) (Emms and Kelly, 2019). 624 Genes known to be involved in cell cycle control in A. thaliana such as Cyclins (Cyc), Cyclin-625 Dependent Kinases (CDKs) and known regulators of these genes were searched in CR0040 626 and *P. aphrodite* proteomes with a combination of Blastp searches and orthogroups. This 627 analysis was applied on CDK-A and -B types as well as Cyc-A -B -D types. Regulators of these 628 genes included cyclin-dependent kinase inhibitor (KRP), transcriptional repressor ILP1, 629 WEE1, Actin-Depolymerizing Factor (ADF), Fizzy-Related proteins (FZR).

630

631

632 Detection of non-endoreplicated genomic regions

633

634 Reads from each sequencing technology used in this study (HiFi, ONT and Illumina reads from 635 CR0040), as well as ONT and Illumina reads from Daphna were mapped on the CR0040 636 assembly. Illumina short-reads and long-reads (HiFi and ONT) were mapped on CR0040 637 assembly using bwa-mem2 (Vasimuddin et al., 2019) and Minimap2 (Li, 2018), respectively. Sequencing depths were averaged for genomic windows of 20 kb. To detect sequencing depth 638 639 bias and to limit the risk to detect false positives, the mean of sequencing depth for every 20 640 successive windows of 20 kb was computed using Illumina reads for Daphna and using long-641 reads (HiFi and ONT) for CR0040. Identified regions were manually validated and refined by 642 visualization of sequencing depth drops for each CR0040 chromosome and for all available 643 sequencing datasets (see details in Supplemental Note 5).

644

645			
646	6 Vanilla Genome Hub (VGH)		
647			
648	The VGH was constructed using the Tripal system, a specific toolkit for the construction of		
649	online community genomic databases, by integrating the GMOD Chado database schema and		
650	the Drupal open source platform (<u>https://www.drupal.org/</u>). The VGH implements a set of		
651	interconnected modules and user-friendly interfaces (details in Supplemental Note 6).		
652			
653	Data availability		
654	The chromosome assembly and accompanying data received the following identifiers in NCBI:		
655	BioProject (with SRA database) ID PRJNA753216, (haplotype A) and PRJNA754028		
656	(haplotype B) BioSample (node) SAMN20691751. Data can be accessed at the vanilla genome		
657	hub site (see below).		
658	RNA-Seq data are readily accessible on the NCBI portal : BioSamples SAMN20691786 (fruit),		
659	SAMN20691787 (leaf), SAMN20691788 (flower), SAMN20691789 (stem), SAMN20691790		
660	(soil root), SAMN20691791 (aerial root), SAMN20691792 (bud), SAMN20691793 (flower		
661	bud), SAMN20691794 (ovary), SAMN20691795 (mixed tissues); SRA : SRR15411867		
662	(mixed tissues), SRR15411868 (ovary), SRR15411869 (flower bud), SRR15411870 (bud),		
663	SRR15411871 (aerial root), SRR15411872 (soil root), SRR15411873 (stem), SRR15411874		
664	(flower), SRR15411875 (leaf), SRR15411876 (fruit)		
665	In addition, these data and various exploration tools are accessible at "Vanilla Genome Hub"		
666	(https://vanilla-genome-hub.cirad.fr/).		
667			
668	8 AUTHORS CONTRIBUTIONS		
669			
670	C.J., M.D., M.G., P.B. contributed to conceptualization of the study.		
671	C.C., C.J., G.S., M.B., M.D., O.B., S.B., W.M. designed the experiments.		
672	M.B., S.S-Y. performed flow cytometry and cytogenetic experiments and analyses.		
673	L.B., J.Z. contributed to the funding of the research, monitored the progress of the work and		
674	supported the researchers throughout the project		
675	C.C., C.L.R., O.B., W.M. performed nucleic acid preparation and sequencing.		
676	C.C., C.J., C.K., G.D., G.S., Q.P., S.B., W.M. performed sequences analyses and assemblies.		
677	C.C., G.D., Q.P., S.B., S.L.A performed genome annotation and built the genome hub		

678 C.C., C.J., C.K, M.G., Q.P. and W.M. outlined the manuscript and wrote first drafts.

- 679 C.C., C.J., C.K., C.L.R., G.D., G.S., M.B., M.D., M.G., P.B., Q.P., S.B., S.S-Y., W.M.
 680 contributed with inputs and revisions in the manuscripts.
- 681

682 The authors declare that they have no conflict of interests.

- 683
- 684

685 ACKNOWLEDGMENTS

686

We are grateful to Jean Bernard Dijoux and Katia Jade for preparing the plant material, and to the Plant Protection Platform (3P, IBISA) for lab facilities and plant resources access (BRC Vatel). We acknowledge the SouthGreen Bioinformatics Platform (<u>http://www.southgreen.fr/</u>) for computational resources access and the GeT-PlaGe platform (INRAE, Toulouse, France) for the use of sequencing facilities. Finally, the authors would like to thank the reviewers for their suggestions which helped to improve the manuscript.

- 693
- 694

695 FUNDING

696

This work was supported by grants from Eurovanille and V. Mane Fils companies. The research was co-funded by the Centre de Coopération Internationale en Recherche Agronomique pour le Développement (CIRAD), the Université de la Réunion (UR), the Institut National de Recherche pour l'Agriculture l'Alimentation et l'Environnement (INRAE), the Centre National de la Recherche Scientifique (CNRS), and the Etablissement Vanille de Tahiti (EVT). This work was also granted by the European Regional Development Fund (ERDF), the Conseil Régional de la Réunion, and the Conseil Départemental de la Réunion.

- 704
- 705

706 **REFERENCES**

707

708 Armstrong, R. L., Penke, T., Chao, S. K., Gentile, G. M., Strahl, B. D., Matera, A. G.,

709 McKay, D. J., & Duronio, R. J. (2019). H3K9 Promotes Under-Replication of

- 710 Pericentromeric Heterochromatin in Drosophila Salivary Gland Polytene Chromosomes.
- 711 Genes, 10(2), 93. 10.3390/genes10020093

- 712 Bhosale, R., Boudolf, V., Cuevas, F., Lu, R., Eekhout, T., Hu, Z.B., Van Isterdael, G.,
- 713 Lambert, G.M., Xu, F., Nowack, M.K., et al. (2018). A Spatiotemporal DNA Endoploidy
- 714 Map of the Arabidopsis Root Reveals Roles for the Endocycle in Root Development and Stress
- 715 Adaptation. Plant Cell **30**:2330-2351. 10.1105/tpc.17.00983.
- 716 Bory, S. (2007). Diversity of *Vanilla planifolia* in the Indian ocean and its related species :
- 717 Genetics, cytogenetics and epigenetics aspect. Université de La Réunion, France.
- 718 Bory, S., Catrice, O., Brown, S., Leitch, I.J., Gigant, R., Chiroleu, F., Grisoni, M., Duval,
- 719 M.F., and Besse, P. (2008). Natural polyploidy in *Vanilla planifolia* (Orchidaceae). Genome
- **720 51**:816-826. 10.1139/G08-068.
- 721 Bourdon, M., Pirrello, J., Cheniclet, C., Coriton, O., Bourge, M., Brown, S., Moise, A.,
- 722 Peypelut, M., Rouyere, V., Renaudin, J.P., et al. (2012). Evidence for karyoplasmic
- homeostasis during endoreduplication and a ploidy-dependent increase in gene transcription
- 724 during tomato fruit growth. Development **139**:3817-3826. 10.1242/dev.084053.
- 725 Bracht, J. R., Fang, W., Goldman, A. D., Dolzhenko, E., Stein, E. M., & Landweber, L.
- **F.** (2013). Genomes on the edge: programmed genome instability in ciliates. Cell, 152(3), 406–
- 727 416. https://doi.org/10.1016/j.cell.2013.01.005
- 728 Braz, G.T., Martins, L.D., Zhang, T., Albert, P.S., Birchler, J.A., and Jiang, J.M. (2020).

A universal chromosome identification system for maize and wild Zea species. Chromosome

730 Research **28**:183-194. 10.1007/s10577-020-09630-5.

- 731 Brown, S.C., Bourge, M., Maunoury, N., Wong, M., Bianchi, M.W., Lepers-
- 732 Andrzejewski, S., Besse, P., Siljak-Yakovlev, S., Dron, M., and Satiat-Jeunematre, B.
- 733 (2017). DNA Remodeling by Strict Partial Endoreplication in Orchids, an Original Process in
- the Plant Kingdom. Genome Biology and Evolution 9:1051-1071. 10.1093/gbe/evx063.
- 735 Cai, J., Liu, X., Vanneste, K., Proost, S., Tsai, W.C., Liu, K.W., Chen, L.J., He, Y., Xu,
- 736 Q., Bian, C., et al. (2015). The genome sequence of the orchid *Phalaenopsis equestris*. Nat
- 737 Genet 47:65-72. 10.1038/ng.3149.
- 738 Chao, Y.T., Chen, W.C., Chen, C.Y., Ho, H.Y., Yeh, C.H., Kuo, Y.T., Su, C.L., Yen, S.H.,
- 739 Hsueh, H.Y., Yeh, J.H., et al. (2018). Chromosome-level assembly, genetic and physical
- 740 mapping of *Phalaenopsis aphrodite* genome provides new insights into species adaptation and
- resources for orchid breeding. Plant Biotechnol J **16**:2027-2041. 10.1111/pbi.12936.
- 742 Cheng, H., Concepcion, G.T., Feng, X., Zhang, H., and Li, H. (2021). Haplotype-resolved
- de novo assembly using phased assembly graphs with hifiasm. Nat Methods 18:170-175.
- 744 10.1038/s41592-020-01056-5.

- 745 Chumová, Z., Záveská, E., Hloušková, P., Ponert, J., Schmidt, P. A., Čertner, M.,
- 746 Mandáková, T., & Trávníček, P. (2021). Repeat proliferation and partial endoreplication
- 747 jointly shape the patterns of genome size evolution in orchids. The Plant Journal : for cell and
- 748 molecular biology, **107**(2):511–524. <u>10.1111/tpj.15306</u>
- Claudel-Renard, C., Chevalet, C., Faraut, T., and Kahn, D. (2003). Enzyme-specific
 profiles for genome annotation: PRIAM. Nucleic acids research, 31(22), 6633–6639.
 10.1093/nar/gkg847
- 752 De Veylder, L., Larkin, J.C., and Schnittger, A. (2011). Molecular control and function of
 753 endoreplication in development and physiology. Trends in Plant Science 16:624-634.
 754 10.1016/j.tplants.2011.07.001.
- 755 Doležel J, Bartoš J, Voglmayr H, Greilhuber J. (2003). Nuclear DNA content and genome
- size of trout and human. Cytom. Part A. **51**A:127–128. doi: 10.1002/cyto.a.10013.
- 757 Du, K., Stock, M., Kneitz, S., Klopp, C., Woltering, J.M., Adolfi, M.C., Feron, R.,
- 758 **Prokopov, D., Makunin, A., Kichigin, I., et al.** (2020). The sterlet sturgeon genome sequence
- and the mechanisms of segmental rediploidization. Nat Ecol Evol 4:841-852. 10.1038/s41559020-1166-x.
- 761 Emms, D.M., and Kelly, S. (2019). OrthoFinder: phylogenetic orthology inference for
 762 comparative genomics. Genome Biology 20Artn 238. 10.1186/S13059-019-1832-Y.
- 763 Favre, F., Jourda, C., Grisoni, M., Piet, Q., Rivallan, R., Dijoux, J.B., Hascoat, J., Lepers-
- 764 Andrzejewski, S., Besse, P., and Charron, C. (2022). A genome-wide assessment of the
- 765 genetic diversity, evolution and relationships with allied species of the clonally propagated
- 766 crop *Vanilla planifolia* Jacks. ex Andrews. Genet Resour Crop Ev. 10.1007/s10722-022767 01362-1.
- 768 Felix, L.P., and Guerra, M. (2005). Basic chromosome numbers of terrestrial orchids. Plant
- 769 Syst Evol **254**:131-148. 10.1007/s00606-004-0200-9.
- 770 Flutre, T., Duprat, E., Feuillet, C., and Quesneville, H. (2011). Considering Transposable
- Element Diversification in De Novo Annotation Approaches. Plos One **6**ARTN e16526
- 772 10.1371/journal.pone.0016526.
- 773 Flynn, J.M., Hubley, R., Goubert, C., Rosen, J., Clark, A.G., Feschotte, C., and Smit, A.F.
- 774 (2020). RepeatModeler2 for automated genomic discovery of transposable element families.
- 775 Proceedings of the National Academy of Sciences of the United States of America 117:9451-
- 776 9457. 10.1073/pnas.1921046117.

- Fu, L.M., Niu, B.F., Zhu, Z.W., Wu, S.T., and Li, W.Z. (2012). CD-HIT: accelerated for
 clustering the next-generation sequencing data. Bioinformatics 28:3150-3152.
 10.1093/bioinformatics/bts565.
- Gurevich, A., Saveliev, V., Vyahhi, N., and Tesler, G. (2013). QUAST: quality assessment
 tool for genome assemblies. Bioinformatics 29:1072-1075. 10.1093/bioinformatics/btt086.
- Han, Y.H., Zhang, T., Thammapichai, P., Weng, Y.Q., and Jiang, J.M. (2015).
 Chromosome-Specific Painting in Cucumis Species Using Bulked Oligonucleotides. Genetics
- **200**:771-779. 10.1534/genetics.115.177642.
- 785 Hasing, T., Tang, H.B., Brym, M., Khazi, F., Huang, T.F., and Chambers, A.H. (2020). A
- phased *Vanilla planifolia* genome enables genetic improvement of flavour and production. Nat
 Food 1:811-819. 10.1038/s43016-020-00197-2.
- 788 Hoede, C., Arnoux, S., Moisset, M., Chaumier, T., Inizan, O., Jamilloux, V., and
- 789 Quesneville, H. (2014). PASTEC: An Automatic Transposable Element Classification Tool.
- 790 Plos One 9ARTN e91929.10.1371/journal.pone.0091929.
- 791 Hřibová, E., Holušová, K., Trávníček, P., Petrovská, B., Ponert, J., Šimková, H.,
- 792 Kubátová, B., Jersáková, J., Čurn, V., Suda, J., Doležel, J., & Vrána, J. (2016). The
- 793 Enigma of Progressively Partial Endoreplication: New Insights Provided by Flow Cytometry
- and Next-Generation Sequencing. Genome biology and evolution, 8(6):1996–2005.
 https://doi.org/10.1093/gbe/evw141.
- 796 Inada, N., Takahashi, N., and Umeda, M. (2021). Arabidopsis thaliana subclass I ACTIN
- 797 DEPOLYMERIZING FACTORs and vegetative ACTIN2/8 are novel regulators of
 798 endoreplication. Journal of Plant Research 134:1291-1300. 10.1007/s10265-021-01333-0.
- 799 Jiang, J.M. (2019). Fluorescence *in situ* hybridization in plants: recent developments and
- 800 future applications. Chromosome Research 27:153-165. 10.1007/s10577-019-09607-z.
- Jones, P., Binns, D., Chang, H. Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H.,
- 802 Maslen, J., Mitchell, A., Nuka, G., Pesseat, S., Quinn, A. F., Sangrador-Vegas, A.,
- 803 Scheremetjew, M., Yong, S. Y., Lopez, R., & Hunter, S. (2014). InterProScan 5: genome-
- scale protein function classification. Bioinformatics (Oxford, England), 30(9), 1236–1240.
- 805 <u>10.1093/bioinformatics/btu031</u>
- 806 Kanehisa, M., Sato, Y., & Morishima, K. (2016). BlastKOALA and GhostKOALA: KEGG
- 807 Tools for Functional Characterization of Genome and Metagenome Sequences. Journal of
- 808 molecular biology, 428(4), 726–731. https://doi.org/10.1016/j.jmb.2015.11.006

- 809 Kim, D., Paggi, J.M., Park, C., Bennett, C., and Salzberg, S.L. (2019). Graph-based genome
- alignment and genotyping with HISAT2 and HISAT-genotype. Nature Biotechnology **37**:907-
- 811 +. 10.1038/s41587-019-0201-4.
- 812 Lang, L., and Schnittger, A. (2020). Endoreplication a means to an end in cell growth and
- 813 stress response. Current Opinion in Plant Biology **54**:85-92. 10.1016/j.pbi.2020.02.006.
- Lee, H.O., Davidson, J.M., and Duronio, R.J. (2009). Endoreplication: polyploidy with
 purpose. Genes Dev 23:2461-2477. 10.1101/gad.1829209.
- 816 Lepers-Andrzejewski, S., Siljak-Yakovlev, S., Brown, S.C., Wong, M., and Dron, M.
- 817 (2011). Diversity and Dynamics of Plant Genome Size: An Example of Polysomaty from a
- 818 Cytogenetic Study of Tahitian Vanilla (Vanilla X Tahitensis, Orchidaceae). American Journal
- 819 of Botany **98**:986-997. 10.3732/ajb.1000415.
- Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics
 34:3094-3100.
- 822 Lilly, M. A., & Spradling, A. C. (1996). The Drosophila endocycle is controlled by Cyclin E
- and lacks a checkpoint ensuring S-phase completion. Genes & development, 10(19), 2514–
 2526. 10.1101/gad.10.19.2514
- 825 Mapleson, D., Accinelli, G.G., Kettleborough, G., Wright, J., and Clavijo, B.J. (2017).
- KAT: a K-mer analysis toolkit to quality control NGS datasets and genome assemblies.
 Bioinformatics 33:574-576. 10.1093/bioinformatics/btw663.
- Marcais, G., and Kingsford, C. (2011). A fast, lock-free approach for efficient parallel
 counting of occurrences of k-mers. Bioinformatics 27:764-770.
 10.1093/bioinformatics/btr011.
- 831 Martins, L.D., Yu, F., Zhao, H.N., Dennison, T., Lauter, N., Wang, H.Y., Deng, Z.H.,
- 832 Thompson, A., Semrau, K., Rouillard, J.M., et al. (2019). Meiotic crossovers characterized
- by haplotype-specific chromosome painting in maize. Nature Communications **10**Artn 4604
- 834 10.1038/S41467-019-12646-Z.
- 835 Matthews, B.J., Dudchenko, O., Kingan, S.B., Koren, S., Antoshechkin, I., Crawford,
- J.E., Glassford, W.J., Herre, M., Redmond, S.N., Rose, N.H., et al. (2018). Improved
 reference genome of *Aedes aegypti* informs arbovirus vector control. Nature 563:501-+.
 10.1038/s41586-018-0692-z.
- 839 Munden, A., Rong, Z., Sun, A., Gangula, R., Mallal, S., & Nordman, J. T. (2018). Rif1
- 840 inhibits replication fork progression and controls DNA copy number in Drosophila. eLife, 7,
- e39140. 10.7554/eLife.39140

- Nair, R.R., and Ravindran, P.N. (1994). Somatic Association of Chromosomes and Other
 Mitotic Abnormalities in *Vanilla planifolia* (Andrews). Caryologia 47:65-73. Doi
- 844 10.1080/00087114.1994.10797284.
- 845 Perez-Silva, A., Odoux, E., Brat, P., Ribeyre, F., Rodriguez-Jimenes, G., Robles-Olvera,
- 846 V., Garcia-Alvarado, M.A., and Gunata, Z. (2006). GC-MS and GC-olfactometry analysis
- 847 of aroma compounds in a representative organic aroma extract from cured vanilla (*Vanilla*
- 848 *planifoli*a G. Jackson) beans. Food chemistry **99**:728-735. 10.1016/j.foodchem.2005.08.050.
- Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing
 genomic features. Bioinformatics 26:841-842.
- **Ravindran, P.N.** (1979). Nuclear Behavior in the Sterile Pollen of *Vanilla planifolia*(Andrews). Cytologia 44:391-396. DOI 10.1508/cytologia.44.391.
- 853 Sallet, E., Gouzy, J., and Schiex, T. (2019). EuGene: An Automated Integrative Gene Finder
- 854 for Eukaryotes and Prokaryotes. Gene Prediction: Methods and Protocols 1962:97-120.
 855 10.1007/978-1-4939-9173-0_6.
- - 856 Sellis, D., Guérin, F., Arnaiz, O., Pett, W., Lerat, E., Boggetto, N., Krenek, S., Berendonk,
 - 857 T., Couloux, A., Aury, J. M., Labadie, K., Malinsky, S., Bhullar, S., Meyer, E., Sperling,
- L., Duret, L., & Duharcourt, S. (2021). Massive colonization of protein-coding exons by
 selfish genetic elements in Paramecium germline genomes. PLoS biology, 19(7), e3001309.
 10.1371/journal.pbio.3001309
- 861 Shimotohno, A., Aki, S.S., Takahashi, N., and Umeda, M. (2021). Regulation of the Plant
- 862 Cell Cycle in Response to Hormones and the Environment. Annual Review of Plant Biology,
- 863 Vol 72, 2021 72:273-296. 10.1146/annurev-arplant-080720-103739.
- 864 Simao, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V., and Zdobnov, E.M.
- 865 (2015). BUSCO: assessing genome assembly and annotation completeness with single-copy
 866 orthologs. Bioinformatics 31:3210-3212. 10.1093/bioinformatics/btv351.
- 867 Stein, J.C., Yu, Y., Copetti, D., Zwickl, D.J., Zhang, L., Zhang, C., Chougule, K., Gao, D.,
- 868 Iwata, A., Goicoechea, J.L., et al. (2018). Genomes of 13 domesticated and wild rice relatives
- 869 highlight genetic conservation, turnover and innovation across the genus *Oryza*. Nat Genet
- **50**:285-296. 10.1038/s41588-018-0040-0.
- 871 Tarailo-Graovac, M., and Chen, N. (2009). Using RepeatMasker to identify repetitive
- 872 elements in genomic sequences. . In Current protocols in bioinformatics A.D. Baxevanis, ed.
- 873 Trávníček, P., Čertner, M., Ponert, J., Chumová, Z., Jersáková, J., & Suda, J. (2019).
- 874 Diversity in genome size and GC content shows adaptive potential in orchids and is closely

- 875 linked to partial endoreplication, plant life-history traits and climatic conditions. The New
 876 phytologist, 224(4):1642–1656. https://doi.org/10.1111/nph.15996.
- 877 Trávníček, P., Ponert, J., Urfus, T., Jersáková, J., Vrána, J., Hřibová, E., Doležel, J., &
- 878 Suda, J. (2015). Challenges of flow-cytometric estimation of nuclear genome size in orchids,
- a plant group with both whole-genome and progressively partial endoreplication. Cytometry.
- 880 Part A : the journal of the International Society for Analytical Cytology, 87(10):958–966.
- 881 https://doi.org/10.1002/cyto.a.22681.
- Vaattovaara, A., Leppälä, J., Salojärvi, J., & Wrzaczek, M. (2019). High-throughput
 sequencing data and the impact of plant gene annotation quality. Journal of experimental
 botany, 70(4):1069–1076. https://doi.org/10.1093/jxb/ery434.

885 Vasimuddin, M., Misra, S., Li, H., and Aluru, S. (2019). Efficient Architecture-Aware
886 Acceleration of BWA-MEM for Multicore Systems. Int Parall Distrib P:314-324.

- 887 10.1109/Ipdps.2019.00041.
- 888 Wicker, T., Sabot, F., Hua-Van, A., Bennetzen, J.L., Capy, P., Chalhoub, B., Flavell, A.,
- Leroy, P., Morgante, M., Panaud, O., et al. (2007). A unified classification system for
 eukaryotic transposable elements. Nat Rev Genet 8:973-982. nrg2165 [pii]
- 891 10.1038/nrg2165.
- 892 Yen, E.C., McCarthy, S.A., Galarza, J.A., Generalovic, T.N., Pelan, S., Nguyen, P., Meier,
- **J.I., Warren, I.A., Mappes, J., Durbin, R., et al.** (2020). A haplotype-resolved, de novo
- genome assembly for the wood tiger moth (*Arctia plantaginis*) through trio binning.
 GigaScience 9ARTN giaa088.10.1093/gigascience/giaa088.
- Younis, A., Ramzan, F., Hwang, Y.J., and Lim, K.B. (2015). FISH and GISH: molecular
 cytogenetic tools and their applications in ornamental plants. Plant Cell Reports 34:1477-1488.
 10.1007/s00299-015-1828-3.
- **Zhang, T., Liu, G.Q., Zhao, H.N., Braz, G.T., and Jiang, J.M.** (2021a). Chorus2: design of
- 900 genome-scale oligonucleotide-based probes for fluorescence *in situ* hybridization. Plant
 901 Biotechnology Journal 10.1111/pbi.13610.
- 902 Zhang, Y.X., Zhang, G.Q., Zhang, D.Y., Liu, X.D., Xu, X.Y., Sun, W.H., Yu, X., Zhu,
- 903 X.E., Wang, Z.W., Zhao, X., et al. (2021b). Chromosome-scale assembly of the *Dendrobium*
- 904 *chrysotoxum* genome enhances the understanding of orchid evolution. Hortic Res-England
- 905 **8**ARTN 183.10.1038/s41438-021-00621-z.
- 906
- 907

908 FIGURE LEGENDS

909

910 Figure 1. Endoreplicated and non-endoreplicated fractions in the CR0040 Vanilla 911 *planifolia* genome. A. The histogram represents the distribution of nuclei in V. *planifolia* nodal 912 tissues according to the strict partial endoreplication state of cells, from 2C (green), to 4E 913 (blue), 8E (vellow), 16E (orange) and 32E (grey). The disks below represent the endoreplicated (coloured) and non-endoreplicated (black) DNA content for each class of nuclei, proportionally 914 915 to their mass (pg). The lower cases f and p denote the respective DNA quantities of F (fixed 916 proportion of the haploid genome which cannot endoreplicate) and P (part participating in 917 endoreplication) fractions. The mean and the standard deviation (sd) of interpeak ratio has been 918 indicated below the dotted arrows. **B.** F and P fractions and P/F ratio values obtained by flow cytometry, and detailed for P fraction for each nuclear class (2C, green; 4E, blue; 8E, yellow, 919 920 16E, orange and 32E, grey). C. Theoretical F and P fractions expected from HiFi sequencing 921 and from flow cytometry obtained data. **D.** Theoretical (dotted) and experimental k-mer 922 coverages for F (black) and P (hatched) fractions.

923

924 Figure 2. Cytogenetic analysis of Vanilla planifolia CR0040. A-D. Orcein staining: A. and 925 **B.** Mitotic metaphases with 2n=32 chromosomes; **C.** Karyotype corresponding to figure B; **D.** 926 Hypoaneuploid mitotic metaphase with 2n=28 chromosomes; E. Karyotype corresponding to 927 figure D; F. Interphase nuclei showing heterochromatic chromocenters; G. DAPI stained 928 interphase nucleus showing unspecific heterochromatin; H. Chromomycin fluorochrome 929 staining with two CMA⁺ regions (arrows) corresponding to rDNA sites; I. Hoechst stained AT-930 rich DNA in metaphase and interphase nucleus (IN), with two fully heterochromatinized 931 chromosomes (arrows). Bar = $10 \mu m$.

932

933 Figure 3. Assemblies k-mers content comparison between CR0040 PacBio HiFi long reads 934 and Daphna Illumina short reads using spectra-cn graph. X-axis represents k-mers 935 multiplicity (counts) and Y-axis the number of distinct k-mers multiplied by their counts. Because of different sequencing depths between read sets, the Y-axis upper values are 10^9 for 936 figures **A** and **B** and 10⁸ for figures **C** and **D**. The area colors indicate the number of k-mer 937 copies (black: 0x or missing k-mers, red:1x, purple: 2x, green: 3x, blue: 4x and orange: 5x) 938 939 found in the assembly. Four spectra-cn plots are presented: A. Daphna reads vs CR0040 940 assembly, **B.** Daphna reads vs Daphna assembly, **C.** CR0040 reads vs CR0040 assembly and 941 **D.** CR0040 reads vs Daphna assembly. The red arrows point towards a low coverage k-mer 942 distribution not expected in a diploid genome assembly spectra-cn graph. The black arrows

point towards the heterozygous (on the left) and homozygous (on the right) k-mer distributions
expected in a diploïd genome assembly. The orange arrows point towards missing k-mers in
the heterozygous k-mer distribution. The lower the black distribution at this location the less
k-mers are missing in the assembly.

947

Figure 4. Overview of the assembled vanilla genome. A. Circos-plot of the genomic content 948 949 along V. planifolia haplotypes A and B and the relationship between them. All tracks are 950 divided in 500Kb genomic windows. From the outside to the inside of the circular 951 representation: ideograms of 28 chromosomes and 2 random mosaic chromosomes that contain 952 the unanchored scaffolds. Gene density (blue), interspersed repeats RepeatMasker hits density 953 (black: retroelements; orange: LTR/Copia ; purple : LTR/Gypsy). Sequencing depth obtained 954 by mapping CR0040 PacBio Hifi reads on the assembly (green) and N density (grey). Syntenic 955 blocks across haplotypes are connected by lines in the innermost part of the figure. **B.** Sequencing depth along the CR0040 A03 and B03 chromosomes (red rectangles) obtained by 956 957 mapping Daphna Illumina (yellow) and ONT (pink) reads, CR0040 PacBio Hifi (blue), 958 Nanopore (green) and Illumina (grey) reads on the CR0040 assembly. Synteny between 959 homologous chromosomes are represented by red boxes. Gaps (N stretches) explaining sudden 960 drops in sequencing depth are shown with white blocks. (1) Low level of sequencing depth for 961 all data. (2) Inverted level of sequencing depth for CR0040 between haplotypes A and B, and 962 constant level of sequencing depth for both Daphna haplotypes. Gene and retrotransposons 963 distributions along the chromosomes are represented by a blue line chart and a stacked 964 histogram (copia: red, gypsy: purple, other retrotransposons: black) respectively.

965

Figure 5. Ratio of k-mers within unanchored and anchored CR0040 genome. This boxplot
shows the ratio of k-mers with a depth below 15 in our HiFi reads, within unanchored
sequences (blue) and within chromosomes (orange).

969

970 Figure 6. Overview (screen shots) of some interoperable vanilla genome analysis tools
971 integrated into the vanilla genome hub. A. Main menu, B. Gene search (Tripal Megasearch),
972 C. Sequence homology search (Blast), D. Gene report (Tripal), E. Genome Browser (JBrowse),
973 F. Metabolic pathway visualization (Pathways Tools), G. GO Enrichment (DIANE), H.
974 Comparison of genomic sequences (Synvisio).

- 975
- 976

TABLE

Table 1. HiFi assembly and annotation statistics of the diploid CR0040 genome.

Total assembly size (Gb)	3.4
Total contig number	24,534
Contig N50 length (Mb)	0.924
Maximum contig length (Mb)	31
GC content (%)	31.6
Number of protein coding genes	59,128
BUSCO completeness (%)	93.2
Total of interspersed repeats (%)	47.0
•	ole Q.







Reads





В





