

## Using RNA-Seq to assemble a rose transcriptome with more than 13,000 full-length expressed genes and to develop the WagRhSNP 68k Axiom SNP array for rose (*Rosa L.*)

Carole F S Koning-Boucoiran, G Danny Esselink, Mirjana Vukosavljev, Wendy P C Van\_t\_Westende, Virginia W Gitonga, Frans Andries Krens, Roeland E Voorrips, W Eric Van\_de\_Weg, Dietmar Schulz, Thomas Debener, Paul Arens and Marinus J M Smulders

Journal Name:	Frontiers in Plant Science
ISSN:	1664-462X
Article type:	Original Research Article
Received on:	28 Jan 2015
Accepted on:	27 Mar 2015
Provisional PDF published on:	27 Mar 2015
Frontiers website link:	<a href="http://www.frontiersin.org">www.frontiersin.org</a>
Citation:	Koning-boucoiran CF, Esselink GD, Vukosavljev M, Van_t_westende WP, Gitonga VW, Krens FA, Voorrips RE, Van_de_weg WE, Schulz D, Debener T, Arens P and Smulders M(2015) Using RNA-Seq to assemble a rose transcriptome with more than 13,000 full-length expressed genes and to develop the WagRhSNP 68k Axiom SNP array for rose ( <i>Rosa L.</i> ). <i>Front. Plant Sci.</i> 6:249. doi:10.3389/fpls.2015.00249
Copyright statement:	© 2015 Koning-boucoiran, Esselink, Vukosavljev, Van_t_westende, Gitonga, Krens, Voorrips, Van_de_weg, Schulz, Debener, Arens and Smulders. This is an open-access article distributed under the terms of the <a href="https://creativecommons.org/licenses/by/4.0/">Creative Commons Attribution License (CC BY)</a> . The use, distribution and reproduction in other forums is permitted, provided the original author(s) or licensor are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

This Provisional PDF corresponds to the article as it appeared upon acceptance, after rigorous peer-review. Fully formatted PDF and full text (HTML) versions will be made available soon.

1 **Frontiers in Plant Science**

2

3 **Using RNA-Seq to assemble a rose transcriptome with more than 13,000 full-length**  
4 **expressed genes and to develop the WagRhSNP 68k Axiom SNP array for rose (*Rosa L.*)**

5

6 Carole F.S. Koning-Boucoiran<sup>1,2</sup>, G. Danny Esselink<sup>1</sup>, Mirjana Vukosavljev<sup>1,3</sup>, Wendy P.C.  
7 van 't Westende<sup>1</sup>, Virginia W. Gitonga<sup>1,4</sup>, Frans A. Krens<sup>1</sup>, R.E. Voorrips, W. Eric van de  
8 Weg<sup>1</sup>, Dietmar Schulz<sup>4</sup>, Thomas Debener<sup>4</sup>, Chris Maliepaard<sup>1</sup>, Paul Arens<sup>1</sup>, Marinus J.M.  
9 Smulders<sup>1\*</sup>

10

11 <sup>1</sup> Wageningen UR Plant Breeding, Wageningen University and Research centre, P.O. Box  
12 386, 6700 AJ Wageningen, The Netherlands.

13 <sup>4</sup> Institute for Plant Genetics, Leibnitz University Hannover, Germany

14 Current addresses:

15 <sup>2</sup> Hogeschool van Arnhem en Nijmegen (HAN), Nijmegen, The Netherlands

16 <sup>3</sup> Pheno Geno Roses D.O.O. Bulevar oslobođenja 65a, 21000 Novi Sad, Serbia

17 <sup>4</sup> Lex+ East Africa, PO Box 1739, Naivasha-20117, Kenya

18

19 \* Corresponding author

20 email addresses: Boucoiran, Carole <Carole.KoningBoucoiran@han.nl>; Esselink, Danny  
21 <danny.esselink@wur.nl>; Mirjana Vukosavljev <mirjana.vukosavljev@wur.nl>; Westende,  
22 Wendy van t <wendy.vantwestende@wur.nl>; Virginia W. Gitonga  
23 <V.Gitonga@DNAGreenGroup.com>; Frans Krens <frans.krens@wur.nl>; Roeland  
24 Voorrips <roeland.voorrips@wur.nl>; Weg, Eric van de <eric.vandeweg@wur.nl>; Dietmar  
25 Schulz <schulz@genetik.uni-hannover.de>; Thomas Debener <debener@genetik.uni-  
26 hannover.de>; Chris Maliepaard <chris.maliepaard@wur.nl>; Paul Arens  
27 <paul.arens@wur.nl>; René Smulders <rene.smulders@wur.nl>

28

29 Running title: Using RNA-Seq to assemble a rose transcriptome and to develop a 68k Axiom  
30 SNP array for rose

31 **Abstract**

32 **In order to develop a versatile and large SNP array for rose, we set out to mine ESTs from**  
33 **diverse sets of rose germplasm. For this** RNA-Seq libraries containing about 700 million  
34 reads were generated from tetraploid cut and garden roses using Illumina paired-end  
35 sequencing, and from diploid *Rosa multiflora* using 454 sequencing. Separate *de novo*  
36 assemblies were performed in order to identify single nucleotide polymorphisms (SNPs)  
37 within and between rose varieties. SNPs among tetraploid roses were selected for  
38 constructing a genotyping array that can be employed for genetic mapping and marker-trait  
39 association discovery in breeding programs based on tetraploid germplasm, both from cut  
40 roses and from garden roses. In total 68,893 SNPs were included on the WagRhSNP Axiom  
41 array.

42 **Next**, an orthology-guided assembly was performed for the construction of a non-redundant  
43 rose transcriptome database. A total of 21,740 transcripts had significant hits with  
44 orthologous genes in the strawberry (*Fragaria vesca* L.) genome. Of these 13,390 appeared  
45 to contain the full-length coding regions. This newly established transcriptome resource adds  
46 considerably to the currently available sequence resources for the *Rosaceae* family in general  
47 and the genus *Rosa* in particular.

48

49 Key-words: Rosa, transcriptomics, EST, SNP array, genotyping, assembly

50

## 51 **Introduction**

52 Whereas cut rose is economically the most important ornamental crop worldwide (761  
53 million euros in The Netherlands in 2011), the rose genome **sequence** has not been completed  
54 yet. In fact, a long history of interspecific hybridization and selection (Debener and Linde,  
55 2009; Smulders et al., 2011; Zhang et al., 2012; Vukosavljev et al., 2013) has led to a  
56 complicated taxonomy that is not fully resolved (Koopman et al., 2008; **Fougère-Danezan et**  
57 **al., 2015**). Commercial rose cultivars are mostly tetraploid and highly heterozygous.  
58 Therefore, inheritance patterns of quantitative traits may be complex. Phenotypic traits such  
59 as flower stem production, flower shape, flower color or disease resistance are of economic  
60 importance for breeders and growers (Debener and Linde, 2009; Smulders et al., 2011), and  
61 need to be better understood genetically in order to be able to apply marker-assisted selection  
62 in breeding programs.

63 EST sequences are an efficient source of various markers for the construction of dense  
64 genetic linkage maps and the identification of QTLs (Vukosavljev et al., 2015). Based on  
65 next-generation sequencing technologies, **four** EST studies on roses have been published so  
66 far. Dubois *et al.* (2012) produced ESTs from 13 rose tissues and studied the expression of  
67 genes involved in flowering and scent biosynthesis. Yan *et al.* (2014) analyzed gene  
68 expression during flower blooming, whereas Kim *et al.* (2012) focused on miRNAs related to  
69 color genes in four rose cultivars. **Yan et al. (2015) analysed ascorbate biosynthesis genes and**  
70 **transcription factors in *Rosa roxburghii* fruits**. These studies used a *de novo* assembly  
71 pipeline. An orthology-guided reference transcriptome assembly was developed by Ruttink *et*  
72 *al.* (2013) in *Lolium*, in order to be able to improve the assembly sequences from a highly  
73 heterozygous species in the absence of a reference genome sequence. The Genome Database  
74 for Rosaceae (GDR, Jung *et al.*, 2014) contains, besides about 500,000 rose ESTs (2013),  
75 data of the fully sequenced woodland strawberry (*Fragaria vesca* L.) genome that may be  
76 used as a reference genome for comparative genomics in rose. Strawberry is the closest  
77 related species for which a genome sequence is available.

78 SNP array platforms have been developed for various *Rosaceae* crops including apple  
79 (Chagné *et al.*, 2012; Bianco *et al.*, 2014), cherry (Peace *et al.*, 2012) and peach (Verde *et al.*,  
80 2012). They have been shown to be an important resource facilitating the production of dense  
81 genetic maps and subsequent QTL mapping of important traits, genome-wide association  
82 analysis, pedigree-based analysis and genomic selection (Bianco et al. 2014). The required  
83 density depends on the type of application and on the degree of LD on the germplasm –

84 which in outcrossing species is often small (a few centimorgan). Dense maps are necessary to  
85 be able to identify haplotypes and study patterns of introgression with sufficient resolution  
86 power (Zhang et al. 2013). Tightly linked markers improve the efficiency of marker-assisted  
87 selection in breeding (Jänsch et al. 2015). In order to generate a versatile and large SNP array  
88 for rose, we set out to mine ESTs from diverse sets of germplasm comprising tetraploid cut  
89 and garden roses as well as a diploid garden rose, so that the SNPs on the array would be  
90 polymorphic in a wide range of rose genetic backgrounds.

91 Here, we present the results of (i) three transcriptome *de novo* assemblies based on tetraploid  
92 cut and garden rose genotypes and a diploid rose, (ii) the development of a 68K genotyping  
93 SNP marker array on the Axiom platform, and the (iii) construction and (iv) annotation of a  
94 non-redundant rose transcriptome for tetraploid roses using the genome sequence of diploid  
95 strawberry. This study generates valuable genomic resources (i.e. EST library with  
96 annotations and genetic diversity between different genotypes) and resulted in the  
97 construction of a large SNP array that can serve as a genotyping platform for future studies in  
98 diploid and tetraploid roses.

99

## 100 **Materials & Methods**

### 101 *Plant material*

102 Three sources of material were used. From the parents of the K5 segregating population of  
103 tetraploid cut roses (*Rosa hybrida*, Yan *et al.*, 2005; Koning-Boucoiran et al., 2012), P540  
104 (mother) and P867 (father), petals were harvested in three stages of flower development (S1,  
105 S2, S3, Fig. 1A). This material is designated as K5. For garden roses (designated as GR)  
106 whole flowers at three flowering stages (closed, half-opened and fully-opened) as well as  
107 young leaves were harvested of twelve European and Canadian garden rose cultivars: Morden  
108 Fireglow, Adelaide Hoodless, Prairie Joy, Morden Blush, Diamond Border, Nipper, J.P.  
109 Connell, Princess of Wales, Heritage, Graham Thomas, Morden Centennial (MC) and Red  
110 New Dawn (RND) (Fig. 1A; Vukosavljev et al., 2013). For the diploid *Rosa multiflora* hybrid  
111 88/124-46 (Biber *et al.*, 2010; designated as Rh88) leaves were harvested from plants grown  
112 under optimal growing conditions (control) and from plants 1-6 days after inoculation with  
113 black spot, *Diplocarpon rosae* (Debener *et al.*, 1998), powdery mildew, *Podosphaera*  
114 *pannosa* (Linde & Debener, 2003), or downy mildew, *Peronospora sparsa* (Schulz *et al.*,

115 2009), one hour after wounding or after 1 hour of 40°C heat stress. Sampled tissues were  
116 immediately frozen in liquid nitrogen and stored at -80°C until RNA extraction.

117

### 118 *RNA preparation*

119 Total RNA was isolated from P540 and P867 by using the RNeasy Plant Mini Kit (QIAGEN,  
120 Westburg, The Netherlands) according to the manufacturer's instructions. The K5 RNA  
121 samples were prepared by pooling equal amounts of total RNA isolated from petals of the  
122 three flower stages as mentioned above. For GR total RNA was isolated from the 12 garden  
123 roses according to Chang *et al.* (1993) with modifications described in Supplementary file  
124 ESM1. Total RNA samples were prepared by pooling equal amounts of RNA from the four  
125 samples (three flower stages and leaves). For Rh88 total RNA from treated and untreated  
126 leaves was isolated using the Invitex RNA extraction kit (STRATEC Molecular GmbH,  
127 Berlin, Germany) according to the manufacturer's instructions. Remaining DNA was  
128 removed by digestion with RNase-free DNase as specified in the extraction kit.

129

### 130 *Sequencing of cDNA*

131 RNA samples of the K5 parents were sequenced by ServiceXS (Leiden, The Netherlands)  
132 using Illumina's standard operation protocols (2 x 75 bp paired end) on a Genome Analyser  
133 II. The RNA samples from the twelve garden rose cultivars (GR) were sent to GATC Biotech  
134 (Constance, Germany) where twelve cDNA libraries were sequenced using 2 x 100 bp paired  
135 end sequencing on a HiSeq 2000. For the pools of stressed and unstressed Rh88 leaves two  
136 random primed normalised cDNA libraries were constructed and sequenced in two 454 FLX  
137 Titanium runs at the Roche 454 sequencing center in Branford (USA). The first pool  
138 consisted of eight independent RNA isolations of untreated leaves and the second pool  
139 consisted of equal amounts of RNA from four independent extractions of all stress-treated  
140 leaves.

141

### 142 *Pre-processing of the sequences*

143 Illumina reads were pre-processed using Prinseq-lite (vs 0.20.3) which included the trimming  
144 of nucleotides having a phred score lower than 25, the trimming of 10 nucleotides from the 5'  
145 end to remove the bias of the nucleotide content of the reads due to the random hexamer

146 priming (Hansen *et al.*, 2010), the trimming of poly A/T tails, the removal of duplicate reads,  
147 of low complexity reads (DUST approach), of reads shorter than 50 nucleotides and of reads  
148 with more than one ambiguous nucleotide. Next, the paired-end reads were processed for  
149 overlapping sequences using COPE (Liu *et al.*, 2012). All unconnected reads were used as  
150 normal paired-end reads, all connected read pairs (i.e. merged read pairs) and single reads  
151 were used as single-end reads.

152 The 454 reads were pre-processed using the FASTX toolkit (v0.0.13) with the same trimming  
153 and filtering as the Illumina reads except that reads smaller than 100 bp were discarded.  
154 Duplicate reads were removed using USEARCH v5.2.32 (Edgar 2010). The reads were  
155 corrected for homopolymer nucleotide tracks using Acacia v1.52 (Bragg *et al.*, 2012).

156

### 157 *Transcriptome de novo assembly*

158 *De novo* transcriptomes were assembled per sample set using Trinity (min\_kmer\_cov 2,  
159 Grabherr *et al.*, 2011) for the Illumina datasets for transcriptome assembly and SNP calling.  
160 For the 454 dataset the MIRA/CAP3 assemblers in the iAssembler (v1.3.2) pipeline (Zheng  
161 *et al.*, 2011) were used for the transcriptome assembly, and the CLC assembler for SNP  
162 detection, in which sequences of the control and stressed leaves were combined to increase  
163 the number of reads. To select for relevant biological transcripts within each data set, RSEM  
164 (RNA-Seq by Expectation-Maximization) was used for transcript abundance estimation. If  
165 less than 1% of the total reads of a component (IsoPct) matched with a specific transcript, the  
166 transcript was not included in the subsequent steps of the analysis. Of the others the most  
167 abundant transcript/isoform was selected. To this end, all reads were mapped against all  
168 transcripts. Rh88 transcripts were also filtered for fungal sequences by blasting against  
169 available fungal sequences of *Marssonina brunnea* (Zhu *et al.*, 2012).

170

### 171 *SNP mining*

172 SNPs were identified within subsets of the sequences: (1) **between** K5 parents P540 and  
173 P867, (2a) **among** twelve garden rose cultivars, (2b) **in a subset of those, namely between** the  
174 **two** garden roses MC and RND (Vukosavljev *et al.* in preparation), and (3) within Rh88 (Fig  
175 1C). For each of the sets the individual transcriptomes were assembled with CAP3 (default  
176 settings with -p 97) to generate a reference transcriptome per dataset. All reads of each  
177 dataset were mapped to their specific reference transcriptome with Bowtie 2 (Langmead **and**



178 Salzberg, 2012) with modified settings (--very-sensitive --rfg 5, 10) and filtered for map  
179 quality (>25) using SAMtools.

180 The resulting SAM file was used for SNP calling using QualitySNPng (Nijveen *et al.*, 2013)  
181 with modified settings (the minimal similarity score per polymorphic site,  
182 similarityAllPolymorphicSites: 0.8, the minimal number of reads per allele set at 5). Using  
183 the filtering options of QualitySNPng, SNPs found in transcripts displaying a larger number  
184 of haplotypes than theoretically expected were discarded. For the tetraploid parents of each  
185 mapping population (P540xP867 and MCxRND), a maximum of 8 haplotypes (four  
186 haplotypes per parent) per transcript was expected. The resulting SNP markers with 35  
187 flanking nucleotides on both sites without additional SNPs or InDels were marked as  
188 potential markers. To filter against paralogous markers, transcripts containing selected SNP  
189 markers were searched against their own reference transcripts (BLASTn e-value 1-30).  
190 Transcripts with two or more hits were discarded since this indicates that they may be present  
191 several times in the genome.

192 To prevent interference with chloroplast DNA during the array hybridization all sequences  
193 around SNP markers were screened against the chloroplast genome of *Fragaria vesca*  
194 (<http://www.rosaceae.org>). The genome sequence of *Fragaria vesca* was used to identify and  
195 remove sequences with potential splice junction sites. For this, markers were searched against  
196 the *Fragaria* sequences (BLASTn e-value 1-5) and discarded using custom perl scripts if their  
197 sequences matched with fewer than 68 bp (95%). All A/T and C/G polymorphisms were also  
198 excluded since genotyping these SNPs requires twice the number of probes using the Axiom  
199 platform. As a last step the mined SNP markers of the three sets were compared to remove  
200 redundancy across sets.

201

202 *Axiom genotyping array: WagRhSNP*

203 The selected SNPs were submitted to Affymetrix (Santa Clara, CA, USA) for a final analysis  
204 to determine whether the probes could be synthesized. Affymetrix discarded SNPs that  
205 shared similar sequences, as this could interfere with hybridization. We decided to include  
206 two probes for each SNP on the array, each probe targeting one of the strands (coded as  
207 AX\_set\_ID), as this allows an additional quality check during genotype calling with dosage  
208 scoring (Smulders *et al.*, 2015). The array also includes 3000 non-polymorphic control probes  
209 (coded as DQC-sample name, DQC is a measure of the extent to which the distribution of

210 signal values is separated from background values). The SNPs on the array are listed in  
211 Supplementary File ESM2, mentioning for each SNP its SNP\_ID, the flanking sequences, the  
212 alleles of the SNP, the Rh-Fv ortholog transcript code, the *Fragaria vesca* protein code, and  
213 the annotation in *Fragaria* whenever available.

214

#### 215 *Orthology-guided reference transcriptome assembly*

216 An orthology-guided assembly procedure according to Ruttink *et al.* (2013) was followed for  
217 the construction of a non-redundant rose transcriptome sequence. The set of non-redundant  
218 proteins of *Fragaria vesca* was downloaded from PLAZA2.5 (Van Bel *et al.* 2012) to guide  
219 the assembly. tBLASTn, carried out in 2013, (e-value cut-off 1e-5 and up to 250 hits  
220 allowed) was used to search for protein hits (e-value 1e-10) with all retained rose transcripts  
221 of this study. The transcripts with a significant tBLASTn hit with a *F. vesca* protein were  
222 grouped and assembled using CAP3 with default settings. The assembled transcripts were  
223 compared to the *F. vesca* proteins using BLASTx, and if the highest-scoring protein returned  
224 the *Fragaria* gene originally used as the highest scorer, then the two genes were considered  
225 as putative orthologs, and the transcript was selected and tentatively named after the most  
226 likely orthologous gene in *F. vesca* (Rh-Fv transcripts). Next, the longest ORF of each  
227 selected rose sequence was determined (Trinity package, Grabherr *et al.* 2011), the 3' and 5'  
228 UTR sequences trimmed off and the remaining sequences were reassembled using CAP3  
229 with default settings to select the final set of orthologous sequences of rose.

230 Functional domains predicted in *Fragaria* and available in GDR (Jung *et al.* 2014;  
231 <http://www.rosaceae.org>) were mined for our rose Rh-Fv transcripts. They were also scanned  
232 for protein signatures from superfamilies reported in various databases such as Smart, Tigr,  
233 Panther, Pfam, FPrint, Profilescan, ProDom and Gene3D (Zdobnov and Apweiler, 2001).

234

## 235 **Results**

### 236 *Individual de novo transcriptome assemblies*

237 RNA-Seq results of K5 (two cut rose genotypes) and GR (twelve garden rose cultivars) were  
238 obtained by using Illumina paired-end sequencing (Fig. 1A, Table 1). To generate the RNA-  
239 Seq data of Rh88 (garden rose Rh88 hybrid, Fig. 1A) 454 sequencing was used.

240 Sequencing of cDNAs from petals of the tetraploid cut rose K5 parents gave 44 million 75 bp  
241 paired-end reads. After quality trimming and merging, 81% of the reads remained either as  
242 paired-end (12 million) or as single end reads (23 million). In total ~78,000 transcripts were  
243 obtained, but ~18,000 (24%) of these were subsequently discarded since they were less  
244 abundant and possibly splice variants. For SNP calling it is better to avoid these. On one hand  
245 they could represent paralogous genes, which would lead to nucleotide differences between  
246 paralogs rather than between alleles of the same locus. On the other hand, if they are from  
247 splice variants they may map in multiple contigs, which could mean that all of these would  
248 unnecessarily be excluded from the SNP calling. Ultimately, approximately 30,000  
249 transcripts per genotype (Table 1) were identified.

250 Similarly, the sequencing of cDNAs from petals and leaves of GR resulted in 632 million 100  
251 bp paired-end reads. After quality trimming and merging, 71% of the reads remained either as  
252 paired-end or single end reads. After filtering of the most abundant transcripts during the  
253 assembly, 48% were discarded as possible splice variants. Ultimately, between 35,000 and  
254 47,000 transcripts per cultivar were identified (Table 1).

255 The sequences of Rh88 of healthy and stressed leaves were combined and resulted in 2.3  
256 million reads with an average length of 360 bp. After quality trimming, duplicate removal  
257 and homopolymer correction 38% remained. Filtering of the most abundant transcripts  
258 yielded 93,974 transcripts (Table 1). This large number was partly due to the presence, in the  
259 black spot and powdery mildew infected leaves, of fungal genes. The sequences were  
260 therefore blasted (blastx) against available fungal sequences of *M. brunnea* (Zhu *et al.* 2012),  
261 based on which 12,705 sequences with an average homology of 53.4% were discarded. The  
262 remainder was used in further analysis.

263

#### 264 *Development of the WagRhSNP array*

265 SNPs were mined in transcripts containing at least one reliable SNP (Tang *et al.* 2006) for the  
266 three rose datasets separately (Table 2). The smallest transcript was 135 bp long, assembled  
267 from 10 reads and showed 1 reliable SNP while the largest transcript was 14,270 bp long,  
268 from 15,286 reads, showing 77 reliable SNPs. Some transcripts contained up to 123 reliable  
269 SNPs. A small majority of the SNPs (62.3% for K5, 59.8% for GR, and 57.7% for Rh88)  
270 were transitions (C/T or G/A), as is common (in almond: 51%, Wu *et al.* 2008; in wheat and  
271 maize: 45% and 55% respectively, Edward *et al.* 2008; in cassava, up to 65%, Lopez *et al.*

272 2005). The average SNP density varied between 0.4 and 0.6 per 100 bp among the three  
273 sample sets but the variation between transcripts within each sample was large, as pointed out  
274 by the standard deviations (Table 2). The distribution of reliable SNPs per transcript is given  
275 in Fig. 2 (note the log scale). Half of the transcripts had 1-3 SNPs per transcript.

276 The WagRhSNP array includes a total of 68,893 reliable SNPs. Of these, 26,354 SNPs were  
277 identified between the cut rose parents (labelled: RhK5\_transcript number\_SNP number),  
278 26,364 SNPs among the 12 garden rose cultivars (named: Rh12GR\_transcript number\_SNP  
279 number), 14,293 SNPs between MC and RND (named: RhMCRND\_transcript number\_SNP  
280 number) and 1,882 SNPs between alleles of Rh88 (named Rh88\_transcript number\_SNP  
281 number). Probes for both strands are on a genotyping array which we named the WagRhSNP  
282 rose array totaling 137,786 probes (described in ESM2). This Axiom<sup>®</sup> array is commercially  
283 available for genetic studies in rose. The availability of a signal from two independent probes  
284 for each SNP enables additional quality control during scoring of the signal dosage, which is  
285 important for accurate genotyping in tetraploids (Vukosavljev et al., in prep; Arens et al., in  
286 prep).

287

### 288 *Orthology-guided assembly of the rose transcriptome*

289 The orthology-guided assembly procedure, developed for *Lolium* by Ruttink *et al.* (2011),  
290 was applied. The non-redundant protein coding sequences from strawberry (34,748 unigenes)  
291 were mapped against the 628,240 rose transcripts identified during the *de novo* assembly  
292 (Fig. 1). In total 381,621 (60.7%) transcripts were mapped against 28437 strawberry  
293 unigenes. They could be assembled into 21740 orthologous sequences, whereby singletons  
294 (i.e., strawberry unigenes with only a hit against single rose transcript) were discarded. Of  
295 these 13,390 sequences (61.6%) corresponded to complete unique ORFs (ESM3).

296 A single *Fragaria* protein could map against up to 25 rose transcripts, which partly  
297 overlapped. For example, when looking at FV0G46670, the orthology-guided assembly  
298 identified 25 rose transcripts out of 250 transcripts that had been identified during the *de novo*  
299 assembly.

300

### 301 *Annotation of the rose transcripts*

302 The annotated rose transcripts with orthologs in *Fragaria* (named Rh-Fv transcripts) were  
303 investigated for functions and involvement in processes such as disease resistance and  
304 defense mechanisms, flower development and flower color (Supplementary file ESM4). The  
305 transcripts were further grouped into GO classes and functional domains, and were mined  
306 based on the InterPro Scan prediction of *Fragaria* (Zdobnov & Apweiler, 2001). A total of  
307 2,498 different protein domains were predicted to be present (ESM5). The putative  
308 annotations of the *Fragaria vesca* genome were linked to the protein domains identified in  
309 our dataset. Within those domains, 8,090 proteins were identified based on 17,726 transcripts  
310 (ESM5). Figure 3 illustrates the distribution of four protein functions (resistance/defense,  
311 flower color, flowering, and cold tolerance) and the number of rose candidate genes  
312 identified within each functional class. For instance, 300,000 transcripts were annotated as  
313 involved in resistance/defense mechanisms, with homology to 1000 *Fragaria* proteins.  
314 Among them six different putative mlo-like genes (up to 25 transcripts) were identified out of  
315 the eight mlo-like genes present in the *Fragaria* genome. These included the full sequence of  
316 RhMLO1, RhMLO3 and RhMLO4, and a partial sequence of RhMLO2 (Kaufman et al.  
317 2012). The unique putative TMV resistance protein N matched with 270 rose genes.

318 GO terms were assigned to the annotated transcripts (ESM4). Around 20,000 genes (38%)  
319 were assigned to both molecular functions (such as transport, signal transduction and  
320 structural molecules) and to biological processes (such as flower development, protein  
321 metabolism and response to stress). Around 10,000 genes (21%) were assigned to the  
322 category cellular compounds (e.g. organelle synthesis/regulation). Table 3 shows, for  
323 instance, that 978 transcripts belong to the GO class (GO:0004674) protein serine/threonine  
324 kinase activity, which includes not only the above identified protein (LRR receptor-like  
325 serine/threonine-protein kinase) but also other types of kinases involved in other processes  
326 (ESM4).

327 The Rh-Fv transcripts (i.e., those with orthologous sequences in the *Fragaria* genome) of our  
328 dataset were compared to the ROSAseq database (Dubois *et al.*, 2012;  
329 <http://iant.toulouse.inra.fr/R.chinensis>). This database contains 80,714 rose EST clusters  
330 (based on 454 sequencing) longer than 100 nucleotides (average length of  $444 \pm 209.4$  bp),  
331 annotated with the *Fragaria vesca* genome. Of these, 56,899 EST clusters had a BLASTn hit  
332 to 14,302 Rh-Fv transcripts of our study with a mean nucleotide identity of 96.2%. In general  
333 multiple ROSAseq EST clusters mapped to a single Rh-Fv transcript: 95% of these Rh-Fv  
334 transcripts matched with up to 10 EST clusters from the ROSAseq database, while ca. 85% of

335 the ROSAseq EST clusters matched a single Rh-Fv transcript (ESM6). For instance, three not  
336 annotated ROSAseq EST clusters (RC013751, RC050162, and RC061808) were similar to  
337 Rh-Fv transcript FV1G02570.m1, which was annotated as a putative mlo-like protein 1.  
338 Furthermore, five EST clusters from the ROSAseq database (RC016326, RC022993,  
339 RC028093, RC040307, RC072319), four of which not annotated, showed similarity to one  
340 Rh-Fv transcript (FV1G02570.m1) annotated as putative mlo-like protein 6. On the other  
341 hand, ROSAseq clusters annotated as putative mlo-like proteins 2, 11 and 14 did not have  
342 any similarities with the Rh-Fv transcripts.

343

## 344 **Discussion**

345 In this study, three sets of transcript sequence data were combined and analyzed in order to  
346 develop a SNP genotyping array. The generated markers may be used to produce dense  
347 genetic linkage maps at tetraploid and diploid level, improve QTL and gene function  
348 analyses, and the study of synteny with *Fragaria*. The genotyping array with 68,893 SNPs is  
349 commercially available for the *Rosaceae* community.

350

### 351 *The WagRhSNP Axiom array*

352 We chose to use the Axiom<sup>®</sup> array system of Affymetrix (Santa Clara, CA, USA).  
353 Advantages of the Axiom array system for SNP detection include the large number of probes  
354 that fit on the array, the small size of conserved probe sequences (so that additional SNPs do  
355 not interfere so often in sequences with a high density of SNPs), and that the array can be  
356 produced for 480 samples (5 microtiter plates) onwards using photolithographic templates (so  
357 that arrays ordered later will be identical to the original ones). Axiom arrays are being  
358 developed for other rosaceous crops as well, notably a 90K array for octoploid strawberry  
359 within the RosBREED project (Smulders et al., 2015).

360 By analyzing three data sets with two different sequencing platforms, more than 68k SNPs  
361 were identified and included on the WagRhSNP Axiom array. Per transcript the SNP  
362 frequencies of the K5 and the GR samples were 1 SNP every 167 bp and 200 bp,  
363 respectively, which is higher than the SNP frequency of 1 SNP/288 bp found in the highly  
364 heterozygous genome of apple (Chagné et al. 2012). SNPs on this array originate for  
365 approximately 40% from tetraploid cut roses and 60% from tetraploid garden roses, but it can  
366 be expected that they will be useful in all tetraploid germplasm, as cut roses represent a

367 subset of the germplasm of garden roses (Vukosavljev *et al.*, 2013). We included around  
368 1000 SNPs identified in diploid *R. multiflora*. As tetraploid roses are the result of extensive  
369 hybridisation between (diploid) species and are probably segmental allopolyploids, many of  
370 the SNPs on the array that have been identified as polymorphism within and between  
371 tetraploid cultivars, may be polymorphic in diploid germplasm as well. Zhang *et al.* (2013)  
372 observed that SNP haplotypes at a SNPSTR locus were shared between tetraploid and diploid  
373 *Rosa* species.

374 The array will also be very useful as the SNPs reside in coding regions of genes that are being  
375 expressed. Genetic map positions of the SNP markers can thus be linked to transcript  
376 sequences and, if available, gene annotation. The Supplementary files contain the keys for the  
377 connections to the genes predicted in the *Fragaria vesca* genome sequence. Similarly, gene  
378 annotations can be screened for candidate genes, which then can be examined for the  
379 presence of SNPs.

380

#### 381 *Assembly issues*

382 The individual assemblies were performed without a reference genome and are therefore  
383 difficult to validate but they highlighted the diversity among the samples and they were used  
384 to identify reliable SNPs. Ruttink *et al.* (2013) indicated that *de novo* assemblies in highly  
385 heterozygous species typically yield more transcripts than the actual number of genes  
386 expressed, and that was the case here as well (Table 1). We proceeded to construct a common  
387 transcriptome from these samples by using orthology-guided assembly from Ruttink *et al.*  
388 (2013). The three different rose datasets had been produced with two different sequencing  
389 platforms (Illumina paired-end on a GAI and a HiSeq, and 454). The Illumina paired-end  
390 reads were short (110 bp to 200 bp) but we obtained up to 15,000 reads per transcript. The  
391 454 reads of the diploid *R. multiflora* cultivar were longer (350 bp) but transcript depth was  
392 limited (max. 400 reads/transcript). Finseth & Harrison (2014) concluded that using Illumina  
393 reads alone one can produce a high quality transcriptome appropriate for RNA-Seq gene  
394 expression analyses, but that utilizing both 454 and Illumina is preferred. Hodgins *et al.*  
395 (2014) came to a similar conclusion.

396

#### 397 *Annotation*

398 The orthology-guided assembly based on the *Fragaria* genome (ESM3) produced 21,740  
399 orthologs, of which 13,390 appeared to be full-length coding regions with an average length  
400 of 1089 bp. In this way, we could identify over 1/3 of all the genes estimated to be present in  
401 the rose genome, a good result as tissues such as root, stem, fruit and seed were not included  
402 in this study. They provide an additional resource to the 14,252 peptides (EST clusters with  
403 an average length of 444 bp) orthologous to *Fragaria* identified by Dubois *et al.* (2012) in  
404 their ROSAseq database, which was produced from diploid roses.

405

## 406 **Conclusion**

407 Our data provides the most comprehensive transcriptome resource currently available for rose  
408 with 13,390 expressed full-length genes identified. This resource adds significantly to the  
409 currently available genomics and bioinformatics resources for the genus *Rosa*. SNPs in many  
410 of these genes are present on the 68k WagRhSNP Axiom array, which will support candidate  
411 gene identification. The dense SNP array with 68,893 SNPs will enable producing dense  
412 genetic maps that are useful in genetic research and marker-assisted breeding.

413

## 414 **Acknowledgements**

415 This research was partly supported by the TTI Green Genetics projects ‘Hyperrose’ and  
416 ‘Polyploids’ and by the TKI-U ‘Polyploids’ project (BO-26.03-002-001). The Roche 454  
417 sequencing centre in Branford (USA) is acknowledged for sequencing the Rh88 samples. L.  
418 Bellon, F. Brew, M. Mittmann, A. Pirani and T. Webster of Affymetrix are thanked for  
419 constructive discussions on the design of the array.

420

## 421 **Supplementary material**

422 ESM1. Adjusted protocol for RNA extraction from rose flowers based on Chang *et al.*  
423 (1993). (docx)

424 ESM2. [The complete description](#) of all SNPs on the WagRhSNP array. For each SNP the  
425 following is listed: SNP\_ID, flanking sequences, and alleles of the SNP (A/B). In addition,  
426 the Rh-Fv ortholog transcript code, the *Fragaria vesca* protein code, and the annotation in  
427 *Fragaria* are given whenever available. [This file provides all sequence information about the](#)



428 SNPs on the array and their flanking regions, but it can also be used as a key to the most  
429 similar sequence in the *Fragaria vesca* (Fv) genome, if available, and vice versa. (csv)  
430 ESM3. Orthologous ORF sequences of roses, constructed using the orthology-guided  
431 assembly procedure. The sequence header contains the Rh\_Fv orthologue name and the  
432 related *Fragaria* (Fv) transcript annotation. (csv)  
433 ESM4. InterPro predictions for the orthologous transcripts of cut and garden roses. The  
434 predictions originate from the related *Fragaria* (Fv) transcript. (csv)  
435 ESM5. GO classes and GO annotations of the Rh-Fv transcripts. The annotations originate  
436 from the related *Fragaria* (Fv) transcript. (csv)  
437 ESM6. Key from ROSASeq transcripts to Rh\_Fv\_orthologs. The file contains ROSASeq  
438 transcripts with high similarity with Rh\_Fv orthologous ORF sequences. The annotations  
439 originate from the related *Fragaria* (Fv) transcript. (csv)

440

441 *ESM3 will be uploaded to GDR and ROSAseq.*

442

#### 443 **Author and Contributors**

444 Conception of the study: CFSK-B, VWG, FAK, WEVDW, TD, CM, PA, MJMS; Collection  
445 of material: CFSK-B, MV, WPCVTW, VWG, DS, PA; Production of the data: CFSK-B,  
446 MV, WPCVTW, DS, DT, MJMS; Assembly: DE; SNP selection and array design: GDE,  
447 REV, WEVDW, CM, PA, MJMS; Annotation of the rose transcripts: CFSK-B, GDE;  
448 Writing of the manuscript: CFSK-B, PA and MJMS. All authors read and approved the final  
449 version of the manuscript.

450

#### 451 **References**

452 Bianco, L., Cestaro, A., Sargent, D.J., Banchi, E., Derdak, S., Di Guardo, M., Salvi, S., Viola,  
453 R., Gut, I., Laurens, F., Chagné, D., Velasco, R., Van de Weg, E., Troggio, M. (2014).  
454 Development and Validation of a 20K Single Nucleotide Polymorphism (SNP) Whole  
455 Genome Genotyping Array for Apple (*Malus × domestica* Borkh). *PLoS ONE* 9(10):  
456 e110377. doi: 10.1371/journal.pone.0110377

457 Biber, A., Kaufmann, H., Linde, M., Spiller, M., Terefe, D., Debener, T. (2010). Molecular  
458 markers from a BAC transcript spanning the *Rdr1* locus: a tool for marker-assisted selection  
459 in roses. *Theoretical and Applied Genetics* 120, 765-773.

460 Bragg, L., Stone, G., Imelfort, M., Hugenholtz, P., Tyson, G. W. (2012). Fast, accurate error-  
461 correction of amplicon pyrosequences using Acacia. *Nature Methods* 9, 425–426.

462 Chagné, D., Crowhurst, R. N., Troggio, M., Davey, M. W., Gilmore, B., Lawley, C.,  
463 Vanderzande, S., Hellens, R. P., Kumar, S., Cestaro, A., Velasco, R., Main, D., Rees, J. D.,  
464 Iezzoni, A., Mockler, T., Wilhelm, L., Van de Weg, E., Gardiner, S. E., Bassil, N., Peace, C.  
465 (2012) Genome-Wide SNP Detection, Validation, and Development of an 8K SNP Array for  
466 Apple. *PLOS ONE* 7, e31745.

467 Chang, S., Puryear, J., Cairney, J. (1993). A simple and efficient method for isolating RNA  
468 from pine trees. *Plant Molecular Biology Reporter* 11, 113-116.

469 Clark, M., Schmitz, C. A., Rosyara, U. R., Luby, J. J., Bradeen, J. M. (2014). A consensus  
470 ‘Honeycrisp’ apple (*Malus × domestica*) genetic linkage map from three full-sib progeny  
471 populations. *Tree Genetics & Genomes* 10, 627-639.

472 Debener, T., Drewes-Alvarez, R., Rockstroh, K. (1998). Identification of five physiological  
473 races of blackspot, *Diplocarpon rosae*, Wolf on roses. *Plant Breeding* 117, 267-270.

474 Debener, T., Linde, M. (2009). Exploring Complex Ornamental Genomes: The Rose as a  
475 Model Plant. *Critical Reviews in Plant Sciences* 28(4), 267-280. doi:  
476 10.1080/07352680903035481

477 Dubois, A., Carrere, S., Raymond, O., Pouvreau, B., Cottret, L., Roccia, A., Onesto, J. P.,  
478 Sakr, S., Atanassova, R., Baudino, S., Foucher, F., Bris, M. L., Gouzy, J., Bendahmane, M.  
479 (2012). Transcriptome database resource and gene expression atlas for the rose. *BMC*  
480 *Genomics* 13, 638.

481 Edgar, R. C. (2010). Search and clustering orders of magnitude faster than  
482 BLAST. *Bioinformatics* 26, 2460-2461.

483 Edward, K. J., Poole, R. L., Barker, G. L. (2008). SNP Discovery in Plants. In *Plant*  
484 *Genotyping II: SNP Technology*, 2, 1. Ed: Henry, R. J. 272 pp.

485 Edwards, D., Forster, J. W., Chagné, D., Batley, J. (2007). What Are SNPs? In *Association*  
486 *Mapping in Plants*. Ed: Oraguzie, N. C., Rikkerink, E. H. A., Gardiner, S. E., De Silva, H. N.  
487 Springer New York: 41-52.

488 Finseth, F. R., Harrison, R. G. (2014). A Comparison of Next Generation Sequencing  
489 Technologies for Transcriptome Assembly and Utility for RNA-Seq in a Non-Model Bird.  
490 *PLOS ONE*, 9, e108550.

491 Fougère-Danezan, M., Joly, S., Bruneau, A., Gao, X. F., Zhang, L. B. (2015). Phylogeny and  
492 biogeography of wild roses with specific attention to polyploids. *Annals of Botany* 115, 275-  
493 291. doi: 10.1093/aob/mcu245.

494 Grabherr, M. G., Haas, B. J., Yassour, M., Levin, J. Z., Thompson, D. A., Amit, I., Adiconis,  
495 X., Fan, L., Raychowdhury, R., Zeng, Q., Chen, Z., Mauceli, E., Hacohen, N., Gnirke, A.,  
496 Rhind, N., di Palma, F., Birren, B. W., Nusbaum, C., Lindblad-Toh, K., Friedman, N., Regev,  
497 A. (2011). Full-length transcriptome assembly from RNA-Seq data without a reference  
498 genome. *Nature Biotechnology* 29, 644-652.

499 Hansen, K. D., Brenner, S. E., Dudoit, S. (2010). Biases in Illumina transcriptome sequencing  
500 caused by random hexamer priming. *Nucleic Acids Research* 38, e131-e131.

501 Hodgins, K. A., Lai, Z., Oliveira, L. O., Still, D. W., Scascitelli, M., Barker, M. S., Kane, N.  
502 C., Dempewolf, H., Kozik, A., Kesseli, R. V., Burke, J. M., Michelmore, R. W., Rieseberg,  
503 L. H. (2014). Genomics of Compositae crops: reference transcriptome assemblies and  
504 evidence of hybridization with wild relatives. *Molecular Ecology Resources* 14, 166-177. doi:  
505 10.1111/1755-0998.12163.

506 Jansch, M., Broggin, G. A. L., Weger, J., Bus, V. G. M., Gardiner, S. E., Bassett, H.,  
507 Patocchi, A. (2015). Identification of SNPs linked to eight apple disease resistance loci.  
508 *Molecular Breeding* 35, 45. doi: 10.1007/s11032-015-0242-4

509 Jung, S., Ficklin, S. P., Lee, T., Cheng, C.-H., Blenda, A., Zheng, P., Yu, J., Bombarely, A.,  
510 Cho, I., Ru, S., Evans, K., Peace, C., Abbott, A. G., Mueller, L. A., Olmstead, M. A., Main,  
511 D. (2014). The Genome Database for Rosaceae (GDR): year 10 update. *Nucleic Acids*  
512 *Research* 42 (D1), D1237-D1244. doi:10.1093/nar/gkt1012.

513 Kaufmann, H., Qiu, X., Wehmeyer, J., Debener, T. (2012). Isolation, molecular  
514 characterization, and mapping of four rose MLO orthologs. *Frontiers in Plant Science* 3, 244.  
515 doi: 10.3389/fpls.2012.00244.

516 Kim, J., Park, J. H., Lim, C. J., Lim, J. Y., Ryu, J. Y., Lee, B.-W., Choi, J.-P., Kim, W. B.,  
517 Lee, H. Y., Choi, Y., Kim, D., Hur, C.-G., Kim, S., Noh, Y.-S., Shin, C., Kwon, S.-Y. (2012).

518 Small RNA and transcriptome deep sequencing proffers insight into floral gene regulation in  
519 Rosa cultivars. *BMC Genomics* 13, 657.

520 Koning-Boucoiran, C. F. S., Gitonga, V. W., Yan, Z., Dolstra, O., van der Linden, C. G., van  
521 der Schoot, J., Uenk, G. E., Verlinden, K., Smulders, M. J. M., Krens, F. A., Maliepaard, C.  
522 (2012). The mode of inheritance in tetraploid cut roses. *Theoretical and Applied Genetics*  
523 125, 591-607.

524 Koopman, W. J. M., Wisseman, V., De Cock, K., Van Huylbroeck, J., De Riek, J.,  
525 Sabatino, G. J. H., Visser, D., Vosman, B., Ritz, C. M., Maes, B., Werlemark, G., Nybom, H.,  
526 Debener, T., Linde, M., Smulders, M. J. M. (2008). AFLP markers as a tool to reconstruct  
527 complex relationships: A case study in Rosa (Rosaceae). *American Journal of Botany* 95,  
528 353-366.

529 Langmead, B. and Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature*  
530 *Methods* 9, 357-359.

531 Liu, B., Yuan, J., Yiu, S.-M., Li, Z., Xie, Y., Chen, Y., Shi, Y., Zhang, H., Li, Y., Lam, T.-  
532 W., Luo, R. (2012). COPE: an accurate k-mer-based pair-end reads connection tool to  
533 facilitate genome assembly. *Bioinformatics* 28, 2870-2874.

534 Linde, M., Debener, T. (2003). Isolation and identification of eight races of powdery mildew  
535 of roses (*Podosphaera pannosa*) (Wallr.: Fr.) de Bary and the genetic analysis of the  
536 resistance gene *Rpp1*. *Theoretical and Applied Genetics* 107, 256-262.

537 Lopez, C., Piegu, B., Cooke, R., Delseny, M., Tohme, J., Verdier, V. (2005). Using cDNA  
538 and genomic sequences as tools to develop SNP strategies in cassava (*Manihot esculenta*  
539 Crantz). *Theoretical and Applied Genetics* 110, 425-431.

540 Nijveen, H., van Kaauwen, M., Esselink, D. G., Hoegen, B., Vosman, B. (2013).  
541 QualitySNPng: a user-friendly SNP detection and visualization tool. *Nucleic Acids Research*  
542 41(Web Server issue), W587–W590. doi: 10.1093/nar/gkt333.

543 Peace, C., Bassil, N., Main, D., Ficklin, S., Rosyara, U. R., Stegmeir, T., Sebolt, A., Gilmore,  
544 B., Lawley, C., Mockler, T. C., Bryant, D. W., Wilhelm, L., Iezzoni, A. (2012). Development  
545 and Evaluation of a Genome-Wide 6K SNP Array for Diploid Sweet Cherry and Tetraploid  
546 Sour Cherry. *PLoS ONE* 7, e48305.

547 Ruttink, T., Sterck, L., Rohde, A., Bendixen, C., Rouzé, P., Asp, T., Van de Peer, Y., Roldan-  
548 Ruiz, I. (2013). Orthology Guided Assembly in highly heterozygous crops: creating a

549 reference transcriptome to uncover genetic diversity in *Lolium perenne*. *Plant Biotechnology*  
550 *Journal* 11, 605-617.

551 Schulz, D. F., Linde, M., Blechert, O., Debener, T. (2009). Evaluation of genus *Rosa*  
552 germplasm for resistance to black spot, downy mildew and powdery mildew. *European*  
553 *Journal of Horticultural Science* 74, 1-9.

554 Shulaev, V., Sargent, D. J., Crowhurst, R. N., Mockler, T. C., Folkerts, O., Delcher, A. L.,  
555 Jaiswal, P., Mockaitis, K., Liston, A., Mane, S.P., Burns, P., Davis, T. M., Slovin, J. P.,  
556 Bassil, N., Hellens, R. P., Evans, C., Harkins, T., Kodira, C., Desany, B., Crasta, O. R.,  
557 Jensen, R. V., Allan, A. C., Michael, T.P., Setubal, J. C., Celton, J. M., Rees, D. J., Williams,  
558 K. P., Holt, S. H., Ruiz Rojas, J. J., Chatterjee, M., Liu, B., Silva, H., Meisel, L., Adato, A.,  
559 Filichkin, S. A., Troggio, M., Viola, R., Ashman, T. L., Wang, H., Dharmawardhana, P.,  
560 Elser, J., Raja, R., Priest, H. D., Bryant, D. W., Fox, S. E., Givan, S. A., Wilhelm, L. J.,  
561 Naithani, S., Christoffels, A., Salama, D. Y., Carter, J., Lopez Girona, E., Zdepski, A., Wang,  
562 W., Kerstetter, R. A., Schwab, W., Korban, S. S., Davik, J., Monfort, A., Denoyes-Rothan,  
563 B., Arus, P., Mittler, R., Flinn, B., Aharoni, A., Bennetzen, J. L., Salzberg, S. L., Dickerman,  
564 A. W., Velasco, R., Borodovsky, M., Veilleux, R. E., Folta, K. M. (2011). The genome of  
565 woodland strawberry (*Fragaria vesca*). *Nature Genetics* 43: 109-116.

566 Smulders M. J. M., Arens, P., Koning-Boucoiran, C. F. S., Gitonga, V. W., Krens, F.,  
567 Atanassov, A., Atanassov, I., Rusanov, K.E., Bendahmane, M., Dubois, A., Raymond, O.,  
568 Caissard, J. C., Baudino, S., Crespel, L., Gudin, S., Ricci, S. C., Kovatcheva, N., Van  
569 Huylenbroeck, J., Leus, L., Wissemann, V., Zimmermann, H., Hensen, I., Werlemark, G.,  
570 Nybom, H. (2011). *Rosa*. Chapter 12 in C. Kole (ed.), *Wild Crop Relatives: Genomic and*  
571 *Breeding Resources Plantation and Ornamental Crops*, Springer-Verlag Berlin Heidelberg  
572 2011. Pp 243-275. doi: 10.1007/978-3-642-21201-7\_12

573 Smulders M. J. M., Voorrips, R. E., Esselink, G. D., Santos Leonardo, T. M., Van 't  
574 Westende, W. P. C., Vukosavljev, M., Koning-Boucoiran, C. F. S., Van de Weg, W. E.,  
575 Arens, P., Schulz, D., Debener, T., Bellon, L., Mittmann, M., Pirani, A., Webster, T., Brew,  
576 F., Cox, P., Maliepaard, C. (2015). Development of the WagRhSNP Axiom SNP array based  
577 on sequences from tetraploid cut roses and garden roses. *Acta Horticulturae* 1064: 177-184.  
578 [http://www.actahort.org/books/1064/1064\\_20.htm](http://www.actahort.org/books/1064/1064_20.htm)

579 Tang, J., Vosman, B., Voorrips, R. E., van der Linden, C. G., Leunissen, J. A. M. (2006).  
580 QualitySNP: a pipeline for detecting single nucleotide polymorphisms and

581 insertions/deletions in EST data from diploid and polyploid species. *BMC Bioinformatics* 7,  
582 438.

583 Van Bel, M., Proost, S., Wischnitzki, E., Movahedi, S., Scheerlinck, C., Van de Peer, Y.,  
584 Vandepoele, K. (2012). Dissecting plant genomes with the PLAZA comparative genomics  
585 platform. *Plant Physiology* 158, 590-600.

586 Verde, I., Bassil, N., Scalabrin, S., Gilmore, B., Lawley, C. T., Gasic, K., Micheletti, D.,  
587 Rosyara, U. R., Cattonaro, F., Vendramin, E., Main, D., Aramini, V., Blas, A. L., Mockler, T.  
588 C., Bryant, D. W., Wilhelm, L., Troglio, M., Sosinski, B., Aranzana, M. J., Arús, P., Iezzoni,  
589 A., Morgante, M., Peace, C. (2012). Development and Evaluation of a 9K SNP Array for  
590 Peach by Internationally Coordinated SNP Detection and Validation in Breeding Germplasm.  
591 *PLoS ONE* 7, e35668.

592 Vukosavljev, M., Zhang, J., Esselink, G. D., van 't Westende, W. P. C., Cox, P., Visser, R. G.  
593 F., Arens, P., Smulders, M. J. M. (2013). Genetic diversity and differentiation in roses: A  
594 garden rose perspective. *Scientia Horticulturae* 162, 320-332.

595 Vukosavljev, M., Esselink, G. D., Van 't Westende, W. P. C., Cox, P., Visser, R. G. F.,  
596 Arens, P., Smulders, M. J. M. (2015). Efficient development of highly polymorphic  
597 microsatellite markers based on polymorphic repeats in transcriptome sequences of multiple  
598 individuals. *Molecular Ecology Resources* 15, 17–27. doi: 10.1111/1755-0998.12289

599 Ward, J. A., Ponnala, L., Weber, C. A. (2012). Strategies for transcriptome analysis in  
600 nonmodel plants. *American Journal of Botany* 99, 267-276.

601 Wu, S.-B., Wirthensohn, M. G., Hunt, P., Gibson, J. P., Sedgley, M. (2008). High resolution  
602 melting analysis of almond SNPs derived from ESTs. *Theoretical and Applied Genetics* 118,  
603 1-14.

604 Yan, H., Zhang, H., Chen, M., Jian, H., Baudino, S., Caissard, J.-C., Bendahmane, M.,  
605 Shubin, L., Zhang, T., Zhou, N., Qiu, X., Wang, Q., Tang, T. (2014). Transcriptome and gene  
606 expression analysis during flower blooming in *Rosa chinensis* 'Pallida'. *Gene* 540, 96-103.

607 Yan, X., Zhang, X., Lu, M., He, Y., An, H. (2015). De novo sequencing analysis of the *Rosa*  
608 *roxburghii* fruit transcriptome reveals putative ascorbate biosynthetic genes and EST-SSR  
609 markers. *Gene* 561, 54-62. doi: 10.1016/j.gene.2015.02.054

610 Yan, Z., Denneboom, C., Hattendorf, A., Dolstra, O., Debener, T., Stam, P., Visser, P.B.  
611 (2005). Construction of an integrated map of rose with AFLP, SSR, PK, RGA, RFLP, SCAR

612 and morphological markers. *Theoretical and Applied Genetics* 110, 766-777. doi:  
613 10.1007/s00122-004-1903-6

614 Zdobnov, E. M., Apweiler, R. (2001). InterProScan—an integration platform for the  
615 signature-recognition methods in InterPro. *Bioinformatics* 17, 847–848.

616 Zhang, J., Esselink, G. D., Che, D., Fougère-Danezan, M., Arens, P., Smulders, M. J. M.  
617 (2013). The diploid origins of allopolyploid rose species studied using single nucleotide  
618 polymorphism haplotypes flanking a microsatellite repeat. *Journal of Horticultural Science*  
619 *and Biotechnology* 88, 85-92. [http://www.jhortscib.org/Vol88/88\\_1/11.htm](http://www.jhortscib.org/Vol88/88_1/11.htm)

620 Zheng, Y., Zhao, L., Gao, J., Fei, Z. (2011). iAssembler: a package for de novo assembly of  
621 Roche-454/Sanger transcriptome sequences. *BMC Bioinformatics* 12, 453.

622 Zhu, S., Cao, Y.-Z., Jiang, C., Tan, B.-Y., Wang, Z., Feng, S., Zhang, L., Su, X.-H., Brejova,  
623 B., Vinar, T., Xu, M., Wang, M.-X., Zhang, S.-G., Huang, M.-R., Wu, R., Zhou, Y. (2012).  
624 Sequencing the genome of *Marssonina brunnea* reveals fungus-poplar co-evolution. *BMC*  
625 *Genomics* 13, 382.

626

627

628 Figure legends.

629

630 Figure 1. Overview of the strategy used to assemble the rose EST data, to mine SNPs, and to  
631 develop the WagRhSNP Axiom SNP array. (A) Pictures of the plant material used to isolate  
632 RNA. (B) Sequencing of the three sets. (C) Data analysis: (i) SNP mining in the three sample  
633 sets (GR: garden roses, RND: MC: Morden Centennial, Red New Dawn), (ii) development of  
634 the WagRhSNP Axiom® genotyping array, (iii) transcriptome assembly and (iv)  
635 identification and annotation of sequences orthologous with *F. vesca*.

636

637 Figure 2. Distribution of reliable SNPs per transcript of K5 (blue bars), GR (red bars) and  
638 Rh88 (green bars).

639

640 Figure 3. Number of proteins (blue bars) and their respective transcripts (red dots) in which  
641 they were identified for relevant protein functions in the sampled tissues.

642



**Table 1.** Sequencing results and analysis.

	<b>K5</b>		<b>GR</b>										<b>Rh88</b>		
	P540	P867	Morden Fireglow	Adelaide Hoodless	Prairy Joy	Morden Blush	Diamond Border	Nipper	J.P. Connell	Princess of Wales	Heritage	Graham Thomas	Morden Centennial	Red New Dawn	
Raw sequences	21,618,027	22,429,015	60,647,695	41,913,639	48,707,935	41,848,533	43,835,093	50,571,712	75,422,211	51,936,149	45,693,660	64,714,441	56,337,877	51,161,500	2,321,272
Pre-processing															
Nb of paired end reads*	5,007,927	7,373,422	26,291,879	16,666,189	21,243,155	18,504,871	19,233,108	17,583,539	29,474,176	22,940,332	19,677,516	27,523,415	23,823,735	22,280,757	
Nb of single end reads*	10,947,580	12,280,512	17,443,525	12,039,348	14,115,627	12,061,017	12,681,671	14,427,176	21,993,145	14,860,946	13,137,663	18,819,850	16,129,980	14,498,963	885,661
De novo assembly															
Transcripts	39,354	38,905	77,119	71,751	77,995	69,404	70,791	74,434	110,332	79,095	66,307	82,446	84,115	82,244	112,023
Discarded variants	9,331	8,894	30,196	39,677	32,298	29,675	35,294	37,457	67,667	41,332	31,281	42,848	45,570	43,066	0
Remaining transcripts	30,023	30,011	41,555	38,318	37,102	39,725	35,497	36,976	42,664	37,763	35,025	39,597	38,545	39,176	93,974

\*Number of paired-end reads and single reads per sample after quality trimming and merging during pre-processing of the data.

Table 2. Number of SNPs mined in the three rose datasets. STD=standard deviation

Sample set	Number of transcripts <sup>1</sup>	Average		Average		Number of reliable SNPs per transcript		Average density	
		transcript length (bp)	STD	transcript coverage <sup>2</sup>	STD	SNPs per transcript	SNP/100 bp	STD	
K5	19080	1134.5	952.2	643.0	2017.6	1-123	0.6	0.41	
GR <sup>3</sup>	51106	1342.9	994.0	787.9	2700.3	1-96	0.5	0.40	
Rh88	5493	1296.2	740.2	48.5	34.7	1-57	0.4	0.53	

<sup>1</sup>Transcripts containing at least one reliable SNP.

<sup>2</sup>Average of the number of reads per transcript.

<sup>3</sup> GR: All 12 garden rose cultivars included.

656 Table 3. List of the top 20 gene ontology (GO) terms the most represented among  
 657 annotated transcripts from cut and garden roses, for each of three main GO categories

<b>GO terms</b>	<b>Number of transcripts</b>
GO:0005515 protein binding	1871
GO:0003677 DNA binding	1077
GO:0004674 protein serine/threonine kinase activity	978
GO:0003676 nucleic acid binding	881
GO:0003700 transcription factor activity	828
GO:0005488 binding	613
GO:0003824 catalytic activity	588
GO:0016491 oxidoreductase activity	553
GO:0003723 RNA binding	532
GO:0005524 ATP binding	464
GO:0000166 nucleotide binding	410
GO:0009055 RNA binding Electron carrier activity	376
GO:0043565 Sequence-specific DNA binding	316
GO:0004888 Transmembrane receptor activity	265
GO:0016301 kinase activity	262
GO:0004497 monooxygenase activity	251
GO:0016758 transferase activity	242
GO:0020037 heme binding	241
GO:0004553 hydrolase activity	218
GO:0005215 transporter activity	205
GO:0006468 protein amino acid phosphorylation	1027
GO:0055114 oxidation reduction	1012
GO:0006355 regulation of transcription, DNA-dependent	644
GO:0008152 metabolic process	553
GO:0055085 transmembrane transport	534
GO:0006508 proteolysis	466
GO:0006915 apoptosis	423
GO:0007165 signal transduction	396
GO:0006952 defense response	340
GO:0005975 carbohydrate metabolic process	336
GO:0006412 translation	329
GO:0045087 innate immune response	260
GO:0006457 protein folding	258
GO:0009651 response to salt stress	244
GO:0046686 response to cadmium ion	232
GO:0006629 lipid metabolic process	186
GO:0009793 embryo development ending in seed dormancy	176
GO:0009737 response to abscisic acid stimulus	170
GO:0006886 intracellular protein transport	165
GO:0006810 transport	162
GO:0005634 nucleus	1401

			<b>compound 20.8 %</b>
GO:0016020	membrane	1175	
GO:0005886	plasma membrane	1172	
GO:0016021	integral to membrane	634	
GO:0009507	chloroplast	566	
GO:0005737	cytoplasm	462	
GO:0005622	intracellular	452	
GO:0005739	mitochondrion	333	
GO:0005773	vacuole	306	
GO:0005840	ribosome	281	
GO:0031224	intrinsic to membrane	259	
GO:0005829	cytosol	243	
GO:0005783	endoplasmic reticulum	205	
GO:0009941	chloroplast envelope	199	
GO:0009570	chloroplast stroma	188	
GO:0005618	cell wall	177	
GO:0009505	plant-type cell wall	173	
GO:0005730	nucleolus	112	
GO:0005794	Golgi apparatus	98	
GO:0009535	Chloroplast thylakoid membrane	96	
			<b>Others 3%</b>

Figure 1.TIFF

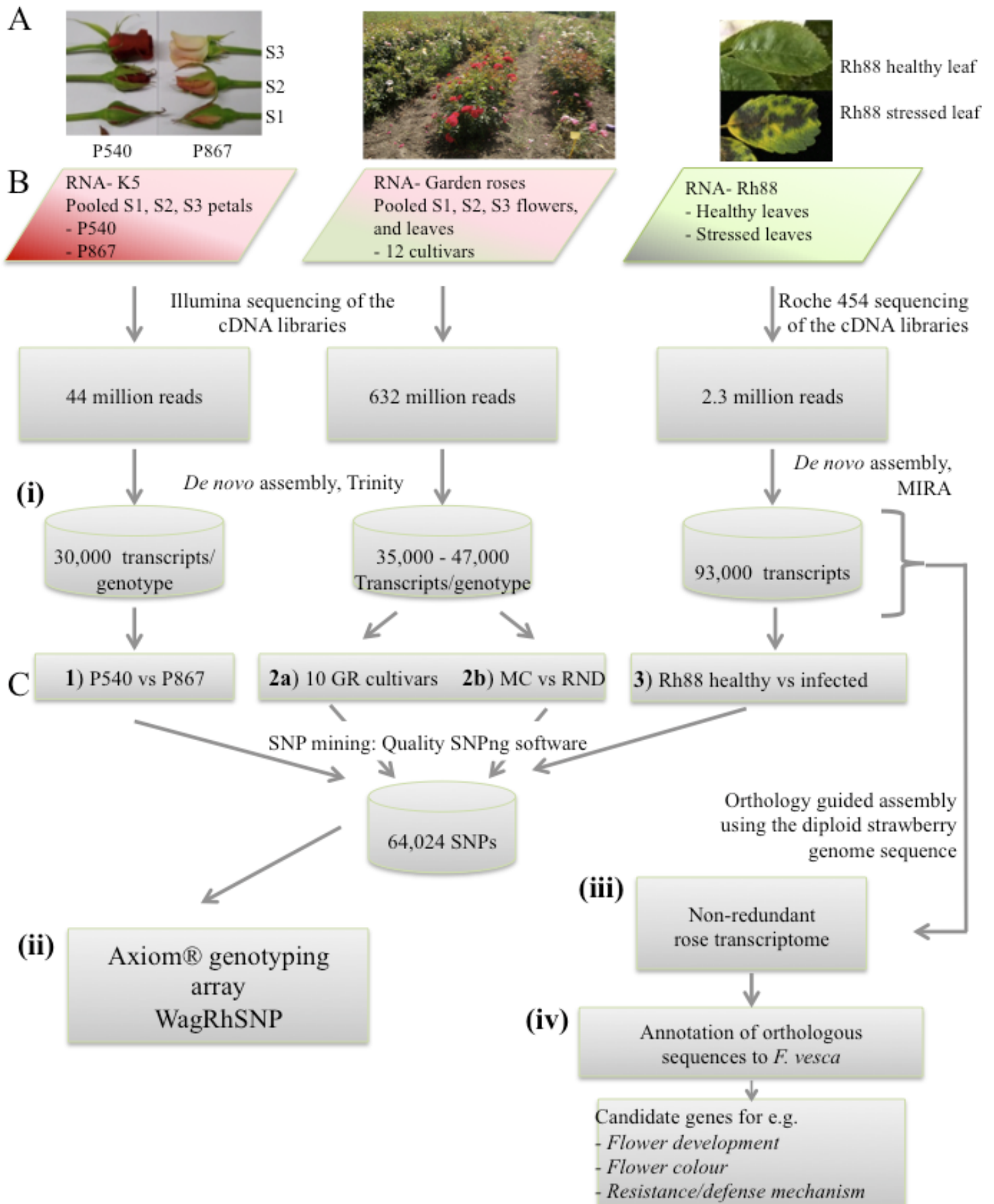


Figure 2.TIFF

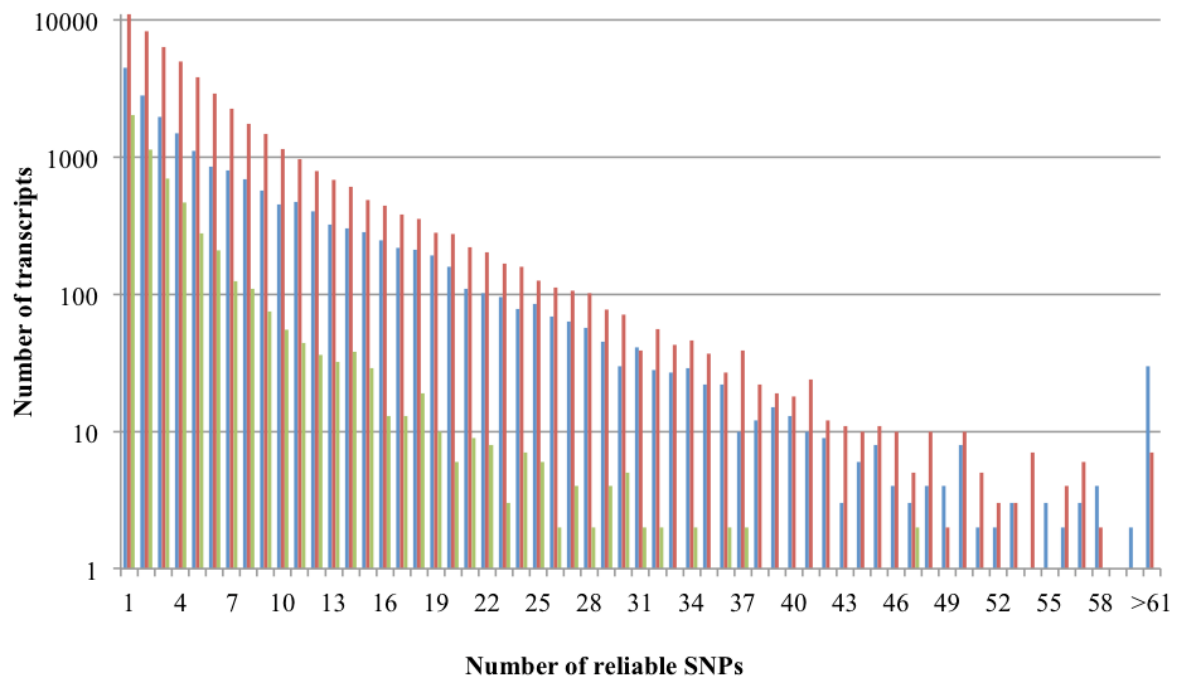


Figure 3.TIFF

