



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

ARCHIVIO ISTITUZIONALE
DELLA RICERCA

Alma Mater Studiorum Università di Bologna Archivio istituzionale della ricerca

Characterization of 25 full-length S-RNase alleles, including flanking regions, from a pool of resequenced apple cultivars

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

Published Version:

De Franceschi, P., Bianco, L., Cestaro, A., Dondini, L., Velasco, R. (2018). Characterization of 25 full-length S-RNase alleles, including flanking regions, from a pool of resequenced apple cultivars. *PLANT MOLECULAR BIOLOGY*, 97(3), 279-296 [10.1007/s11103-018-0741-x].

Availability:

This version is available at: <https://hdl.handle.net/11585/636519> since: 2018-10-11

Published:

DOI: <http://doi.org/10.1007/s11103-018-0741-x>

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).
When citing, please refer to the published version.

(Article begins on next page)

This is a post-peer-review, pre-copyedit version of an article published in “Plant Molecular Biology”. The final authenticated version is available online at: <https://doi.org/10.1007/s11103-018-0741-x>

This version is subjected to Springer Nature terms for reuse that can be found at: <https://www.springer.com/gp/open-access/authors-rights/aam-terms-v1>

1 **Characterization of 25 full-length *S-RNase* alleles, including flanking regions, from a pool of resequenced**
2 **apple cultivars**

3 Paolo De Franceschi¹ (corresponding author), Luca Bianco², Alessandro Cestaro², Luca Dondini¹, Riccardo
4 Velasco^{2,3}

5

6 ¹ Dipartimento di Scienze e Tecnologie Agroalimentari (DISTAL), Università degli Studi di Bologna, Bologna,
7 Italy

8 ² Research and Innovation Centre, Fondazione Edmund Mach, San Michele all'Adige, Trento, Italy

9 ³ Present address: Centro di Ricerca in Viticoltura ed Enologia del Consiglio per la Ricerca in Agricoltura e
10 l'Analisi dell'Economia Agraria (CREA-VE), Conegliano, Treviso, Italy

11 e-mail: paolo.defranceschi2@unibo.it - Telephone number: +39-0512096422

12

13 **Key message**

14 Data obtained from Illumina resequencing of 63 apple cultivars were used to obtain full-length *S-RNase*
15 sequences using a strategy based on both alignment and *de novo* assembly of reads.

16 **Abstract**

17 The reproductive biology of apple is regulated by the *S-RNase*-based gametophytic self-incompatibility
18 system, that is genetically controlled by the single, multi-genic and multi-allelic *S* locus. Resequencing of
19 apple cultivars provided a huge amount of genetic data, that can be aligned to the reference genome in
20 order to characterize variation to a genome-wide level. However, this approach is not immediately
21 adaptable to the *S*-locus, due to some peculiar features such as the high degree of polymorphism, lack of
22 colinearity between haplotypes and extensive presence of repetitive elements. In this study we describe a
23 dedicated procedure aimed at characterizing *S-RNase* alleles from resequenced cultivars. The *S*-genotype of
24 63 apple accessions is reported; the full length coding sequence was determined for the 25 *S-RNase* alleles
25 present in the 63 resequenced cultivars; these included 10 previously incomplete sequences (*S*₅, *S*_{6a}, *S*_{6b}, *S*₈,
26 *S*₁₁, *S*₂₃, *S*₃₉, *S*₄₆, *S*₅₀ and *S*₅₈). Moreover, sequence divergence clearly suggests that alleles *S*_{6a} and *S*_{6b},
27 proposed to be neutral variants of the same alleles, should be instead considered different specificities. The
28 promoter sequences have also been analyzed, highlighting regions of homology conserved among all the
29 alleles.

30 **Keywords:** self-incompatibility, *S*-locus, *Malus*, *S*-genotyping, apple genome

31 **Acknowledgements:** This work has been partly funded under the EU Seventh Framework Programme by
32 FruitBreedomics project no. 265582: Integrated Approach for Increasing Breeding Efficiency in Fruit Tree
33 Crop, and the Italian Ministry of Education, Universities and Research project PRIN 2015BPM9H3:
34 Investigating Self Incompatibility DEterminants in fruit trees (ISIDE). We would like to thank Dr. Massimo
35 Pindo and Dr. Erika Stefani (Research and Innovation Centre, Fondazione Edmund Mach, San Michele
36 all'Adige, Trento, Italy) for technical support on sequencing reactions.

37 Introduction

38 The domesticated apple (*Malus × domestica* Borkh.) is one of the most important crops in temperate
39 regions; it belongs to the family Rosaceae, tribe Pyreae and subtribe Pyrinae (Potter et al. 2007), together
40 with pears (gen. *Pyrus*) and other cultivated tree species such as quince, loquat and medlar. All the Pyrinae
41 have a distinctive base chromosome number $x=17$, which is considered to have arisen from an ancestor
42 with $x=9$ through an event of autopolyploidization and subsequent aneuploidization (Evans and Campbell
43 2002). The center of origin of domesticated apple is thought to be in Kazakhstan and Central Asia (Vavilov
44 1930), but thanks to its great genetic variability it was adapted to cultivation in different environments
45 ranging from subtropical to subarctic climates. A wide genetic base is maintained also due to self-
46 incompatibility, which forces allogamy by preventing self-fertilization (de Nettancourt 2001).

47 Apple and several other Rosaceae possess the *S*-RNase-based gametophytic self-incompatibility system,
48 which is considered the most widespread and ancient self-incompatibility system in eudicots (Steinbachs
49 and Holsinger 2002). In this system, self-pollen tubes are recognized and rejected by the *S*-locus encoded,
50 pistil-specific ribonuclease (*S*-RNase); compatible pollen tubes, on the other hand, have the ability to
51 inactivate the *S*-RNases through a pool of *S*-locus *F-box Brother* genes (*SFBBs*, Sassa et al. 2007) which are
52 thought to bind it in an allele-specific way and mark it for proteolytic degradation by the ubiquitine-
53 proteasome pathway (reviewed by De Franceschi et al. 2012). The recognition of all non-self *S*-RNase alleles
54 is provided by a pool of different *SFBBs* acting collaboratively (Kubo et al. 2010) and being all encoded
55 within the *S*-locus. The suppression of recombination allows the multi-genic *S*-haplotypes to be inherited as
56 single segregating units, maintaining their integrity and functionality across generations (Wang et al. 2012).
57 On a practical point of view, knowledge of *S*-genotypes of cultivars is fundamental both to establish
58 productive cultivar combinations in orchards and to plan compatible crosses in breeding; therefore, PCR-
59 based *S*-genotyping assays have been developed since the early characterization of the *S*-RNase gene
60 (Janssens et al. 1995).

61 The genome of domesticated apple was sequenced for the first time from ‘Golden Delicious’ (Velasco et al.
62 2010) and subsequently from a doubled haploid derived from the same cultivar (Daccord et al. 2017). The
63 initial polyploidization, the hybridization with wild relatives and the forced allogamy of apple resulted in a
64 complex genome, in which syntenic portions of chromosomes reflect the rearrangements of the two
65 ancestral genomes and a huge variability is maintained (Velasco et al. 2010). The availability of a genome
66 sequence and modern high-throughput sequencing technologies make it possible to characterize the
67 genetic variation of apples by resequencing cultivars and aligning them to the reference genome, also
68 enabling the development of dense SNP-genotyping arrays which efficiently cover the whole genome
69 (Chagné et al. 2012; Bianco et al. 2014; Bianco et al. 2016).

70 While this approach proved to be extremely effective in other genomic regions, the *S*-locus shows a series
71 of features that greatly hamper this process. First of all, different *S*-haplotypes proved to be very variable in
72 terms of number, position and orientation of genes (Minamikawa et al. 2010; Okada et al. 2011; Wang et
73 al. 2012; Okada et al. 2013); this particular structure is the product of a complex evolutionary dynamics of
74 *SFBBs*, in which gene duplication and horizontal transfer are expected to play a major role (De Franceschi et
75 al. 2012; Kubo et al. 2015). The lack of colinearity between haplotypes, together with the massive presence
76 of transposable elements, are the major hurdles when aligning short reads to a reference genome, as the
77 haplotypes borne by the reference cultivar are not likely to show a common structure with accessions

78 having different *S*-alleles. The *S-RNase* gene, moreover, shows a degree of allelic polymorphism much
79 higher than other genes; evidence of positive selection, favoring amino acid change and the rise of new
80 alleles, has been reported by several studies (Ishimizu et al. 1998a; Vieira et al. 2007; Vieira et al. 2010; De
81 Franceschi et al. 2011). Frequency-dependent balancing selection, on the other hand, preserved *S-RNase*
82 alleles from being lost due to frequency fluctuations and genetic drift, maintaining allelic polymorphism
83 generated before speciation in the Pyrinae (Sassa et al. 1996; Ishimizu et al. 1998a; Raspé and Kohn 2002;
84 Dreesen et al. 2010; De Franceschi et al. 2011). As a result, many alleles are present in extant species:
85 around 50 different *S-RNases* have been reported in *Malus* (Kim et al. 2016) and the highest sequence
86 homology can be found among alleles of different species or genera rather than within the same species.
87 The single intron of the *S-RNase* has a length varying from around 100bp to more than 3kb (Takasaki-
88 Yasuda et al. 2013). For all these reasons, a reference sequence does not always provide a good template
89 for the alignment of resequencing reads, either at haplotype or single-gene level.

90 In the present study, we take advantage of the availability of genetic data from 63 previously resequenced
91 apple cultivars (Bianco et al. 2016) to characterize the *S-RNase* gene and its flanking regions, using a
92 dedicated procedure specifically developed to identify and assemble *S-RNase* allele sequences. The full-
93 length sequence of 25 alleles is described, enabling efficient comparisons and resolving previous doubts
94 about *S* specificities of grown apple cultivars.

95

96 **Materials and methods**

97 **Plant material and resequencing**

98 In a previous study, 65 apple accessions were resequenced using an Illumina HiSeq2000 platform with the
99 aim of developing an Axiom®Apple480K SNP genotyping array (Bianco et al. 2016). The resequencing panel
100 included 63 diploid apple cultivars chosen to cover the highest possible genetic diversity of apple
101 germplasm and two doubled haploids; in the present study only the sequence reads obtained from the 63
102 diploid cultivars were used. Freeze-dried leaves used for DNA extraction were provided by various
103 institutions as described by Bianco et al. (2016); the cultivar ‘Mela Rozza’ used for sequencing the allele S6
104 (see later in the text) is housed at the experimental station of Cadriano of the Department of Agricultural
105 Sciences of the University of Bologna (Italy).

106 **Synthetic Sequences**

107 We developed a two-step approach based on the construction of synthetic sequences to determine the *S*-
108 genotype of each cultivar. A first synthetic sequence was built using the information available in literature
109 and on Genbank (<https://www.ncbi.nlm.nih.gov/genbank/>) and used to gain insights on the genotypes
110 available through alignment of the resequencing data. The second synthetic sequence was built to extend
111 the information by means of a local assembly phase.

112 **Alignment of reads to known apple *S-RNases* and determination of *S*-genotypes**

113 The first synthetic sequence contains 32 apple *S-RNase* alleles (see Table 1 for accession numbers)
114 separated from each other by a spacer of 1000 ‘Ns’. Full-length genomic sequences were used when
115 available, but for 13 out of 32 alleles only a portion of the coding sequence has been characterized (Table

116 1). For the S_8 allele, moreover, besides an incomplete coding sequence the intron was also lacking, given
117 that the available sequence was obtained from mRNA. In this case an additional 1000 Ns spacer was
118 manually inserted to separate the two exons. Alleles with multiple neutral variants, such as S_{16} (Matsumoto
119 and Furusawa 2005), were included selecting a single variant to avoid redundancy in the alignment of short
120 reads.

121 S -genotypes were determined by aligning resequencing data against the synthetic sequence. Reads were
122 aligned with BFAST (Homer et al. 2009) keeping all best scoring matches and filtering results allowing 2
123 mismatches; this parameter was set after testing the strategy on the cultivar 'Cox's Orange Pippin', with
124 known S -genotype, allowing 2 to 7 mismatches per read in different runs to determine the best setting to
125 ensure allele-specificity.

126 After aligning and filtering, a coverage analysis was performed to determine the hypothetical alleles. The
127 mean coverage was calculated separately for each exon using SAMtools (Li et al. 2009). S -genotypes were
128 attributed assuming the presence of the two alleles with the highest coverage. We imposed, for calling the
129 second allele, a coverage threshold of 30% with respect to the first one; when this threshold was not
130 reached the presence of an unknown allele was assumed for that cultivar, and it was further analyzed with
131 a different strategy (see "Identification of alleles not included in the synthetic sequence").

132 The strategy was repeated on all the 63 cultivars. Alignments were individually visualized and inspected
133 using Tablet (Milne et al. 2013) to confirm the presence of each allele and to check for possible mismatches
134 with the sequences used as input, as it may occur for neutral variants of the same allele.

135 ***De novo* assembly of flanking regions**

136 To obtain sequence information about flanking regions, sequencing reads had to be assembled *de novo*
137 starting from the ends of known sequences. Targeted *de novo* assemblies were performed using GapFiller
138 version 1.10 (available at <https://www.baseclear.com/services/bioinformatics/basetools/gapfiller/>, Boetzer
139 and Pirovano 2012), making the target regions figure as a gap to fill. The second synthetic sequence was
140 therefore built, in which each S -*RNase* allele identified through the alignment described above was flanked
141 by two artificial 1000 N gaps, as follows:

142 A50 – N1000 – Allele 1 – N1000 – Allele 2 – N1000 – (...) – Allele n – N1000 – C50

143 The initial (A50) and final (C50) strings were inserted to provide an end for the terminal gaps flanking the
144 first and last alleles. This synthetic sequence was used as template for GapFiller, so that each region
145 flanking an allele in the template could be considered by the program as a gap to fill with Illumina reads,
146 under the assumption that only the gaps flanking the two alleles borne by an accession could be filled with
147 reads coming from that accession. Two changes were made with respect to the sequences used in the
148 alignment step: the allele S_{16a} was substituted with S_{16b} (Acc. Number: AF327222 mRNA, full length cds, Van
149 Nerum et al. 2001; AB428430 genomic, partial cds, Kim et al. 2009), as this is the variant actually detected
150 in all the six cultivars in which S_{16} was detected; and S_{6a} was substituted by S_{6b} (Acc. Number AB094493
151 genomic, partial cds, Matsumoto et al. 2003) in the template to be used for the two cultivars which proved
152 to possess this variant. Again, an artificial gap was also placed to separate the two S_8 exons as the intron
153 sequence was unknown. Initially, the S_9 -*RNase* allele was chosen to test the strategy as its flanking regions
154 were sequenced from a BAC library from the cultivar 'Florina' and available at NCBI GeneBank database

155 under the accession number AB270792 (Sassa et al. 2007). GapFiller was run separately on each accession
156 possessing the *S₉-RNase* allele, setting the minimum overlap length to 20 and 3 iterations per individual.
157 The *S₉* region was selected from the GapFiller output for each genotype individually, and the obtained
158 sequences were aligned and compared to the known *S₉* genomic region of ‘Florina’ using ClustalOmega
159 (Sievers et al. 2011). The same approach was then applied to all the 63 cultivars, increasing the number of
160 iterations per individual when necessary to increase the length of the output sequence. The output
161 sequence for each allele carried by each single cultivar was extracted and all the obtained sequences of the
162 same allele were aligned together using ClustalOmega.

163 **Identification of alleles not included in the synthetic sequence**

164 For those cultivars showing only one clearly identifiable *S-RNase* allele, the alignment of reads to the initial
165 synthetic sequence was filtered increasing the number of allowed mismatches per read to 4; the output
166 was then visually inspected to find reads aligning to alleles other than the first one. The reads were
167 isolated, aligned to known alleles using BlastN (Altschul et al. 1990) and manually assembled in short
168 contigs, which were then used to build a separate template for GapFiller as described above, running 10
169 iterations for each individual plus 10 additional iterations when the output sequence was not including the
170 full-length coding sequence of the *S-RNase* allele. The sequences obtained were aligned to known allele
171 using BlastN to find homologies and to determine the putative coding sequence.

172 **Experimental validation of *in silico* assemblies**

173 Genomic DNA was extracted from freeze-dried leaves of selected cultivars using a standard CTAB
174 (cetyltrimethylammonium bromide) protocol (Maguire et al. 1994) and diluted to a final concentration of
175 50ng/μl.

176 RNA extraction was performed from the styles of 10 flowers collected from the cultivar ‘Mela Rozza’ at full
177 bloom, using the protocol described by Zamboni et al. (2008). cDNA was synthesized from 100 ng of total
178 RNA using SuperScript II Reverse Transcriptase (Invitrogen, Carlsbad, CA, USA) following the manufacturer’s
179 instructions. To enable 3’ RACE (rapid amplification of cDNA ends) a tailed oligo(dT) primer was used (5’-
180 gactcgagtcgacatcga-t17-3’) according to the protocol described in Sambrook et al. (1989). The 3’ end of
181 cDNA was then amplified using a gene-specific forward primer and a reverse one designed to anneal on the
182 tail generated by reverse transcription.

183 Primer pairs (Online Resource 1) were designed on the obtained sequences using Primer3 (Untergasser et
184 al. 2012) and PCR was performed in 20μl mixtures containing 1X PCR buffer, 2.5 mM MgCl₂, 0.2 mM of
185 each dNTP, 0.3 mM of each primer, 0.03 U/μl AmpliBiotherm DNA polymerase (Fisher Molecular Biology,
186 Rome, Italy) and 50ng genomic DNA or 1 μl of the reverse transcription mixture diluted 1:10 as template.
187 The thermal profile included an initial denaturation at 95°C for 3’ followed by 30 cycles of denaturation at
188 95° C for 20’’, annealing at 58°C (54°C when the primer MdProm-box1f was used, Online Resource 1) for
189 30’’, extension at 72°C for 1’; and final extension at 72°C for 8’. For cloning and sequencing, primer pairs
190 were designed to amplify the largest possible portion of the sequence including flanking regions; amplicons
191 sized 1810 and 1348bp were obtained for *S₂₃* and *S₅₈*, respectively; for the allele *S8*, however, given the
192 large size of the obtained sequence (4098 bp) we decided to divide it in three partially overlapping chunks
193 (MdS8_5UTR/MdS8_Ex1rev, 1297bp; MdS8_Ex1for/PycomC5r1, 1642bp; and MdS8_Ex2for/MdS8_3UTR,
194 1417bp) covering a total of 3983bp.

195 PCR products were cloned into pGEM-T Easy vectors (Promega, Madison, WI, USA) following
196 manufacturer's instructions and three independent colonies were isolated for each fragment. Plasmids
197 were extracted using GeneJet Plasmid Miniprep Kit (Thermo Scientific, Waltham, MA, USA) and sequenced
198 using universal primers T7 and M13rev, plus consensus primers PycomC1f1 and PycomC5r1 (Sanzol 2009a)
199 as internal oligos when necessary to sequence very large inserts.

200 **Sequence analysis**

201 The 5' and 3' flanking regions obtained for each *S-RNase* allele were aligned using Kalign (Lassmann and
202 Sonnhammer 2005). The promoter regions were analyzed using TSSPlant (Shahmuradov et al. 2017) to
203 determine the putative transcript start site (TSS) and TATA box. Coding regions and deduced proteins were
204 aligned using ClustalW (Larkin et al. 2007) and a phylogenetic tree was built using the Neighbor-Joining (NJ)
205 algorithm as implemented in MEGA6 (Tamura et al. 2013), using complete deletions and the maximum
206 composite likelihood calculated using the three codon positions. Statistical support for the topologies was
207 obtained by bootstrap analysis with 1000 replications.

208

209 **Results and discussion**

210 **Determination of the S-genotype of resequenced cultivars**

211 Illumina reads of 63 apple cultivars (Bianco et al. 2016) were aligned to a synthetic sequence containing 32
212 known *S-RNase* alleles with the purpose of identifying the *S-RNase* alleles specific of each cultivar. This pool
213 of 32 known alleles included all the alleles reported by Matsumoto (2014) with the exclusion of *S₁₈*, for
214 which no sequence is available; the alleles *S₄₁*, *S₄₂*, *S₄₄*, *S₄₅*, *S₄₆*, *S₅₀* and *S₅₃* were added because previously
215 characterized in apple cultivars (De Franceschi et al. 2016) or known to be present in the Italian apple
216 germplasm from previous analyses carried out in our laboratory (unpublished data). This pool is expected
217 to cover most of the allelic variability of domesticated apple at this locus based on recent reports
218 (Matsumoto 2014; Kim et al. 2016; Larsen et al. 2016); however, *S-RNases* from related *Malus* species could
219 also be present in the analyzed apple germplasm (Dreesen et al. 2010; Halász et al. 2011), as well as
220 possible new, still uncharacterized alleles.

221 We assumed that under the appropriate filtering parameters, reads from each cultivar would selectively
222 cover the positions on the synthetic sequence corresponding to the alleles present in its genome; as all
223 cultivars are diploid and forcedly heterozygous at the *S*-locus, we expected to find two matches for each
224 genotype, both with a similar coverage given that the two alleles are present in a 1:1 ratio. We tested our
225 system using a cultivar of known *S*-genotype, 'Cox's Orange Pippin' (*S₅S₉*; Janssens et al. 1995; Bošković and
226 Tobutt 1999) filtering with a maximum number of allowed mismatches ranging from 2 to 7 (i.e. from 2% to
227 7% of the read). Importantly, the filtering of reads had to be strict enough to preserve allele-specificity, but
228 relaxed enough to detect possible variants or point mismatches in the sequence of an allele. A possible
229 source of error in this approach is the presence of repetitive elements in *S-RNase* introns, which could be
230 highly represented between reads and result in a high coverage even when the rest of the allele sequence
231 is not covered; this was the case, in the first runs from 'Cox's Orange Pippin', of alleles *S₃*, *S₁₀*, *S₁₆*, *S₂₅*, *S₄₂*
232 and *S₄₅*, which showed a very high, non-allele specific coverage within introns (data not shown). On the
233 other hand, only exons from *S₅* and *S₉* yielded a good coverage under all the tested conditions. These

234 preliminary results highlighted that non-specific matching of reads on intron sequences could lead to alleles
235 mis-identification; all the subsequent coverage analyses, therefore, were carried out on the two exons
236 separately not taking into consideration reads aligning to introns.

237 The best results from the test cultivar 'Cox's Orange Pippin' were obtained setting the maximum allowed
238 mismatches per read to 2; the hits were distributed only on the two expected alleles and the mean
239 coverage for each of their exons ranged from 10.2 to 12.5, consistently with a >20-fold genome coverage
240 and a 1:1 ratio of the two alleles (Online Resource 2). The same settings were extended to the alignment of
241 reads from all the other cultivars. The coverage obtained for the exons of all the 42 template *S-RNases* is
242 reported on Online Resource 2. Ideally, a 10X coverage on each allele would be expected for a 20X genome
243 coverage; however, the number of reads obtained from each accession and therefore the coverage varied
244 quite significantly, as reported by Bianco et al. (2016).

245 The designed synthetic sequence was suitable for defining 2 alleles for 59 of 63 cultivars (Online Resource
246 2), while only one allele was detected for the remaining 4 cultivars, which indicated that the second one
247 was not included in the synthetic sequence. In total 22 different *S-RNase* alleles were detected in the
248 analyzed accessions. All the six cultivars possessing allele S_{16} featured the S_{16b} variant, while among the
249 three matching allele S_6 , one ('Mela Rozza') possessed S_{6a} and two ('Heta' and 'Maikki') S_{6b} . The lack of
250 significant coverage on more than two alleles for diploid cultivars reflects the accuracy of this alignment
251 strategy for defining the S-genotypes from resequencing data; we could not test this approach with triploid
252 accessions, as none was present in the resequenced pool, but these results suggest that it could be equally
253 effective even with higher ploidy levels, provided that a sufficient genome coverage is obtained.

254 To our knowledge, the S-genotype of 14 of the 63 accessions was analyzed in previous studies; 11 of these
255 are confirmed by our results, namely 'Braeburn' (S_9S_{24}), 'Cox's Orange Pippin' (S_5S_9), 'Delicious' (S_9S_{28}),
256 'Filippa' (S_7S_{24}), 'Fuji' (S_1S_9), 'Jonathan' (S_7S_9), 'Lady Williams' (S_7S_{23}), 'Macoun' (S_3S_{25}), 'McIntosh' ($S_{10}S_{25}$),
257 'Spässerud' (S_1S_7) and 'Åkerö' ($S_1S_?$) for which a single S-allele was reported (Janssens et al. 1995; Sassa et
258 al. 1996; Matsumoto et al. 1999; Van Nerum et al. 2001; Kitahara and Matsumoto 2002a; Kitahara and
259 Matsumoto 2002b; Broothaerts et al. 2004; Garkava-Gustavsson et al. 2008; Nybom et al. 2008;
260 Matsumoto 2014). In three cases, however, the alleles detected did not match the previously reported S-
261 genotype : 'Borowitsky' (syn. 'Oldenburg') was reported as S_3S_7 (Long et al. 2010a) while we detected
262 alleles S_3S_{28} ; 'Priscilla' was reported as S_3S_9 (Broothaerts et al. 2004) or S_9S_{20} (Morita et al. 2009), but our
263 results indicate its genotype to be S_7S_{10} ; and finally 'Worcester Pearmain', previously reported as S_2S_{24}
264 (Kitahara et al. 2000), resulted $S_{24}S_{25}$. In all these cases the coverage on the detected alleles was good (>12
265 for 'Borowitsky'; >11 for 'Priscilla'; >10 for 'Worcester Pearmain', Online Resource 2); the reads aligning to
266 each allele were isolated from each genotype and assembled to a contig, which confirmed to correspond to
267 the detected *S-RNase* allele after the alignment to the corresponding reference sequence (Online Resource
268 3). We therefore propose to correct the S-genotypes of these three cultivars accordingly.

269 ***De novo* assembly of allele sequences and flanking regions**

270 PCR-based S-genotyping assays, in apple as well as in other species, require *a priori* sequence information;
271 most methods currently rely on the use of allele-specific primers, restriction enzymes or a combination of
272 both methods for the identification of *S-RNase* alleles (De Franceschi et al. 2012 and references therein;
273 Larsen et al. 2016). However, not all the known alleles have been sequenced completely. Full-length
274 sequences are available for 19 of the 32 alleles used in our synthetic sequence (Table 1); among the

275 remaining alleles, 2 are incomplete in the first exon (S_5 and S_{44}), 9 in both exons (S_6 , S_{11} , S_{21} , S_{23} , S_{29} , S_{42} , S_{45} ,
276 S_{46} , S_{50} and S_{53}) and one (S_8) lacks the intron and part of the first exon. Extending the sequence information
277 is important especially for those alleles, like S_6 , for which only a small portion of the coding sequence is
278 known. The characterization of flanking regions, the promoter and terminator, is also desirable to improve
279 our knowledge on the *S-RNase* gene and open new possibilities for molecular assays.

280 The high-coverage whole-genome resequencing of 63 cultivars yielded a number reads in which virtually all
281 positions in the genome are expected to be represented; these reads have been aligned against the apple
282 reference genome (Velasco et al. 2010) enabling the study of genetic variation within this species to a level
283 previously unexplored (Bianco et al. 2016). When it comes to the *S-RNase*, however, this approach is not
284 immediately effective due to the extremely high degree of polymorphism of this gene (Ma and Oliveira
285 2002; Vieira et al. 2007; Vieira et al. 2010; De Franceschi et al. 2011), which hampers the process of aligning
286 reads to a reference sequence especially outside the coding region (Ushijima et al. 1998). Therefore, *S-*
287 *RNase* coding and flanking regions are indeed present in the resequencing output, but cannot immediately
288 be obtained by simply aligning to a reference sequence.

289 To exploit this amount of information fragmented into 100bp reads we envisaged a strategy based on a
290 targeted *de novo* assembly. This operation is routinely applied in genomic studies to fill sequence gaps in
291 genome assemblies, using programs such as GapFiller (Boetzer and Pirovano 2012); we reasoned that the
292 same program could help in reconstructing unknown *S-RNase* regions, if they are treated as gaps in an
293 artificial sequence. To test this possibility, we took advantage of the availability of the genomic sequence
294 surrounding the S_9 -*RNase*, which was obtained by a BAC library and deposited on GenBank (Sassa et al.
295 2007). An initial testing template was built by using the S_9 -*RNase* sequence from the start to the stop codon
296 (including intron) and creating two artificial gaps of 1000 nucleotides at both sides. Gapfiller was then run
297 using this template from all the cultivars which proved to possess the allele S_9 ; 10 initial iterations resulted
298 in the elongation of the sequence by 140 to 633bp at the 5' side in the different cultivars (Fig. 1a); this
299 variation is likely due to the different coverage obtained from the accessions. On the opposite side, a much
300 shorter sequence was obtained at the 3' end: in this case, assembly of reads was stopped by the presence
301 of a >70bp microsatellite starting 16 bp downstream the stop codon (Fig. 1b). The newly assembled
302 portions were aligned to the known S_9 region, and proved to be perfectly matching (Fig. 1).

303 This initial test showed that the approach with GapFiller could be effective in assembling reads to extend
304 known allele sequences, but also highlighted that the low coverage obtained from some cultivars and the
305 presence of repetitive DNA close to the coding region are two major problems of this strategy. A template
306 containing all the 22 *S-RNase* alleles detected in the pool of cultivars was then used with GapFiller from all
307 the accessions. After the GapFiller run, all the extended sequences were isolated and the output for the
308 same allele were aligned to check their consistency; importantly, we did not observe any aspecific allele
309 extension, i.e. only the two alleles present in a cultivars were extended by the program (data not shown).
310 The procedure was repeated for those cultivars in which the newly generated portion of sequence was
311 <300 bp upstream the start or downstream the stop codon.

312 The results of this approach are summarized in Table 2. In most cases the program successfully assembled
313 several hundred up to more than one thousand bp on both the 5' and 3' flanking regions. The worst results
314 were observed on the S_6 (*a* and *b*) alleles, for which the full-length coding sequence could not be
315 determined despite the several additional iterations tested. This poor result was likely a consequence of

316 the low coverage obtained from the cultivars 'Mela Rozza', 'Heta' and 'Maikki', also observed on the first
317 alignment of reads (Online Resource 2). Five alleles (S_4 , S_9 , S_{32} , S_{33} , S_{46}) were extended on the 3' side for less
318 than 100bp (Table 1), possibly due to the presence of microsatellites downstream the stop codon, as
319 observed in the initial test with S9 and at least in 3 of these alleles (S_4 , S_9 , S_{46}). Finally, for the allele S50 the
320 assembly was very satisfactory at the 3' end (1028bp downstream the stop codon) but only yielded a 14bp
321 fragment of the 5' flanking region (Table 1); of the three cultivars possessing S_{50} , two had a medium-to-low
322 coverage ('Aport Kuba' and 'Durello di Forli', ranging from 3.4X to 7.1X on S50 exons on the initial
323 alignment) while the third, 'Ag Alma', resulted covered >10X on this allele (Online Resource 2). Therefore,
324 also considering the good results obtained on the 3'side, the inability to assemble the 5' flanking region was
325 unlikely to be due to coverage problems and the reason for this result remains unclear.

326 Identification of new alleles

327 For 4 cultivars, namely for 'Åkerö', 'Antonovka Pamtorutka', 'Ijunscoe Ranee' and 'Young America', no
328 second alleles were clearly found. An attempt to identify them was therefore made using reads that aligned
329 with the relaxed parameters (4 mismatches).

330 After manual inspection, candidate reads aligning on the S_{50} allele were found in 'Åkerö' and chosen for a
331 step of manual assembly that yielded a template for the gapfiller phase. A sequence of 1549 bp was
332 obtained, with 2 exons of 252 and 432 bp separated by an intron of 164 bp. Unfortunately, for 'Antonovka
333 Pamtorutka', the limited number of reads did not allow for the assembly phase; however, the perfect
334 match of Antonovka reads against the contig assembled from 'Åkerö' (data not shown) suggested that the
335 two cultivars share the same allele; this finding was later confirmed by cloning and sequencing the allele
336 from 'Antonovka Pamtorutka' (see next section).

337 The assembled sequence shows a considerable similarity with the S_{50} -*RNase* (93.1 % identity at nucleotide
338 level in the coding region) but nevertheless it can be considered a different allele with a good confidence,
339 as its deduced protein differs from S_{50} for 30 amino acids, showing a 86,8% sequence identity. A blast
340 search of the nucleotide collection (nt) database highlighted a 96% nucleotide identity with S_{37} of *M.*
341 *sylvestris* (NCBI GenBank Acc. Number EU419864, Dreesen et al. 2010); for this allele only a 200bp portion
342 of the coding region is known, encoding a 66 amino acids sequence which shows 5 mismatches (92,4%
343 identity) with the allele from 'Åkerö' and only 2 (97,0% identity) with S_{50} . On the other hand, a 98,7%
344 identity both at the nucleotide and protein level was found with allele S_{121} from *Pyrus communis* (NCBI
345 GenBank Acc. Number EU477839, Sanzol 2009b) suggesting the common origin of this allele before the
346 divergence between the genera *Malus* and *Pyrus*, as frequently observed for *S-RNases* (Ishimizu et al.
347 1998a; Raspé and Kohn 2002). While this manuscript was in preparation, an allele from the crabapple 'Mt
348 Blanc' named S_{58} (Acc. Number MG262529) was released, that shows a high sequence identity with the
349 allele identified in 'Åkerö' and 'Antonovka Pamtorutka'. The two sequences differ for 4 nucleotides, one in
350 the first exon causing a conservative amino acid change (leucine to phenylalanine), one in the intron, and
351 two silent substitutions in the second exon (Online Resource 4). We therefore named the new allele *M. x*
352 *domestica* S_{58} and defined the *S*-genotypes of 'Åkerö' and 'Antonovka Pamtorutka' as S_1S_{58} and S_8S_{58} ,
353 respectively.

354 Using the same approach, hits partially matching S_5 were found for 'Ijunscoe Ranee' and 'Young America',
355 isolated and assembled into sequences to be used as template for GapFiller; the two cultivars showed the
356 presence of the same allele. The obtained sequence consists of two 252 and 438bp exons, a long 1446bp

357 intron, a 360bp 5' flanking region and 267bp downstream the stop codon. The alignment with known
358 sequences using Blast revealed to be almost identical to the sequence of allele S_{39} from *Malus sylvestris*
359 (NCBI GenBank Acc. Number EU419871, Dreesen et al. 2010); even in this case the full-length sequence is
360 not available, but the 204bp portion of the coding sequence showed 3 single nucleotide mismatches with
361 respect to the allele obtained from 'Ijunscoe Ranee' and 'Young America', 2 of which being synonymous
362 substitutions and the third one causing a conservative amino acid change (methionine to isoleucine)
363 (Online Resource 4). The two sequences thus showed a 98,5% identity both at the nucleotide and protein
364 level; although the existence of more relevant differences in other portions of the protein cannot be ruled
365 out, we considered the new allele as the *M. × domestica* variant of *M. sylvestris* S_{39} -*RNase* and defined the
366 S-genotypes of for 'Ijunscoe Ranee' and 'Young America' as $S_{39}S_{46}$ and $S_{28}S_{39}$, respectively.

367 **Validation and sequencing of uncharacterized regions**

368 The initial analysis carried out on the S_9 -*RNase* supported the accuracy of sequence assembly by the
369 program GapFiller; however, we decided to validate the *in silico* sequence assembly by amplifying, cloning
370 and sequencing a set of alleles including S_8 from 'Antonovka', S_{23} from 'Lady Williams' and S_{58} from
371 'Antonovka Pamtorutka'. These alleles were chosen among those for which leaf material was available, to
372 include: an allele with a very large flanking region and intron portions assembled *in silico* (S_8); a previously
373 unknown allele, whose sequence was entirely assembled *in silico* (S_{58}); and an additional control (S_{23}). The
374 obtained sequences were aligned with those generated *in silico* by GapFiller, highlighting a 100% identity
375 (Online Resource 5).

376 Additionally, we further investigated those alleles for which the assembly yielded shorter or incomplete
377 portions, i.e. S_{50} and S_{6a}/S_{6b} . The sequence of S_{50} was extended by GapFiller by 1028bp on the 3' side, but
378 only 14bp upstream the start codon. Interestingly, it resulted quite similar to that of S_{58} , as also highlighted
379 by the partial coverage obtained from 'Åkerö' on S_{50} in the initial alignment (Online Resource 2). Therefore
380 we combined the primer designed on the 5' region of S_{58} (MdS58_5UTR) with MdS50_3UTR designed on
381 the 3' extreme of S_{50} (Online Resource 1), arguing that the high sequence similarity could allow the
382 amplification of both alleles with the same primers. The amplification was tested on 'Durello di Forlì' (S_3S_{50} ,
383 Online Resource 2) and yielded a fragment sized 1745bp, which was cloned and sequenced and proved to
384 include a 343bp fragment of the 5' flanking region. Again, the alignment with the *in silico* obtained
385 sequence highlighted a 100% identity in the overlapping portion (Online Resource 5), confirming once again
386 the accuracy of sequence assembly with GapFiller.

387 For the couple of alleles S_{6a}/S_{6b} , the low coverage obtained for the cultivars 'Mela Rozza', 'Heta' and
388 'Maikki' greatly hampered the assembly process, making it impossible to obtain the complete exon
389 sequences (Table 2). We therefore defined a different strategy to amplify the full-length gene; on the 5'
390 side, we used the consensus (non-allele specific) primer MdProm-box1f designed on a conserved motif
391 found after the alignment of the promoters of all the other alleles (see the next section "Analysis of
392 promoter sequences"), coupled with PycomC5r1 (Online Resource 1). Although none of these primers was
393 specific for S_6 , the amplification from 'Mela Rozza' and 'Heta' yielded two clearly distinct bands from each
394 genotype, as in both cases the S_6 allele was coupled to a much larger one: S_{6a} with S_3 in 'Mela Rozza', S_{6b}
395 with S_{16b} in 'Heta' (S_{6a} and S_{6b} amplified fragments sized about 1000bp, while amplicons from S_3 and S_{16b}
396 were sized >2100bp and >3000bp respectively). After cloning the amplicons, colonies bearing the fragment
397 of the correct allele could be easily distinguished by PCR and selected for sequencing, yielding in both cases

398 the entire 252bp first exon plus a 274bp and 273bp flanking region for S_{6a} and S_{6b} , respectively. The 3' side
399 was amplified from 'Mela Rozza' cDNA through 3' Rapid Amplification of cDNA Ends (RACE); we designed a
400 primer (MdS6a_Ex1for, Online Resource 1) which could selectively amplify S_{6a} excluding S_3 , and coupled it
401 to the primer designed on the oligo(dT) RACE adapter. This allowed us to amplify, clone and sequence a
402 994bp fragment including part of the first exon, the whole intron and second exon (167 and 435bp
403 respectively) plus a 344bp 3' untranslated region (UTR). Interestingly, this latter portion highlighted a very
404 large (230bp) microsatellite region of mixed TA and CA repeats. Three different reverse primers
405 (MdS6a_3UTR1-3, Online Resource 1) were designed on the S_{6a} 3'UTR of 'Mela Rozza', one upstream and
406 two downstream the microsatellite repeats, and tested on 'Heta' coupled with S6a_Ex1for in order to
407 obtain amplification of the corresponding region of S_{6b} . Of these, however, only the upstream one
408 (MdS6a_UTR1), annealing 10bp after the stop codon, proved to work on S_{6b} allowing the characterization of
409 only 9bp of the 3'UTR of this latter allele, but including the complete second exon which made it possible to
410 obtain the full-length coding sequence of the S_{6b} allele from 'Heta'.

411 **Analysis of coding sequences and difference between S_6 (= S_{6a}) and S_{17} (= S_{6b})**

412 Overall, the full length coding sequence was determined for all the 25 *S-RNase* alleles present in the 63
413 resequenced cultivars; these included 10 previously incomplete sequences (S_5 , S_{6a} , S_{6b} , S_8 , S_{11} , S_{23} , S_{39} , S_{46} ,
414 S_{50} and S_{58}). The nucleotide and deduced amino acid sequences were aligned and a Neighbor-Joining tree
415 was built based on the protein alignment (Fig. 2, 3).

416 Surprisingly, the distance (number of amino acid substitutions per site/per sequence) between S_{6a} and S_{6b} -
417 *RNases*, which are proposed to be neutral variants of the same allele, was comparable to that between
418 clearly distinct alleles (Table 3). The complete protein sequences of S_{6a} and S_{6b} differ for 14 amino acids (Fig.
419 2) out of 228 and their calculated distance is 0.063; two other couple of alleles differ for 14 amino acids,
420 S_3/S_{10} and S_1/S_{20} , whose distances were 0.063 and 0.064 respectively. Similar but slightly higher differences
421 were calculated for S_{20}/S_{24} (15 amino acids, 0.069), S_9/S_{46} (16 amino acids, 0.073) and S_1/S_{24} (17 amino acids,
422 0.078).

423 Before the present study, only partial genomic sequences were known for S_{6a} (previously also reported as
424 S_{12}) and S_{6b} (S_{17} or S_{19}), limiting the possibility to compare their protein sequences to a portion of just 54
425 amino acid residues, showing a single difference; nevertheless, the two alleles were thought to encode the
426 same *S*-specificity as they showed an identical hyper-variable region (Fig. 2), which is thought to be mainly
427 responsible for *S*-specific protein interaction (Matsumoto et al. 2003; Morita et al. 2009; Matsumoto 2014).
428 Later, however, it was proved that different *S-RNase* alleles can share identical hyper-variable regions,
429 highlighting that other residues in the protein are indeed also involved in the determination of specificities
430 (Zisovich et al. 2004). As no conclusive evidence supporting their functional identity has been provided so
431 far, some authors considered separately the alleles S_{6a} and S_{6b} , maintaining the designations S_6 and S_{17}
432 (Dreesen et al. 2010; Kim et al. 2016).

433 It is known that neutral variants of the same *S-RNase* allele coexist within a species; it is the case, for
434 example, of apple $S_{16a}/S_{16b}/S_{16c}$ (Matsumoto and Furusawa 2005; Morita et al. 2009) or European pear S_{104-1}/S_{104-2}
435 (Sanzol 2010). In all these cases, however, sequence identity is much higher than within the pair
436 S_{6a}/S_{6b} : S_{16a} , S_{16b} and S_{16c} encode identical proteins, while S_{104-1} and S_{104-2} differ for 2 amino acids; moreover,
437 in both cases, the functional identity of their *S*-specificities was proved after pollination tests, which to our
438 knowledge have not been performed between S_{6a} and S_{6b} ; only the identity between alleles previously

439 reported as S_{17} and S_{19} , which both correspond to S_{6b} , has been functionally proved (Matsumoto et al.
440 2006). On the other hand, alleles S_3 and S_5 from *Pyrus pyrifolia* differ for only 11 amino acids, yet encoding
441 clearly distinct specificities (Ishimizu et al. 1998b).

442 The distribution of the 14 amino acid differences between S_{6a} and S_{6b} further supports a difference in their
443 encoded specificities; although the first three substitutions (residues 3, 9 and 24; Fig. 2) fall in the signal
444 peptide and can therefore be excluded from playing a role in the allele-specific interaction with pollen
445 determinants, at least 5 of the remaining 11 replacements (residues 160 between C4 and C5; 212, 215, 216
446 and 224 after C5, Fig. 2) fall in sites previously reported to be under positive selection (Vieira et al. 2010),
447 and therefore likely involved in determining S -specificity.

448 All these evidences strongly suggest that S_{6a} and S_{6b} encode distinct S -specificities, although a pollination
449 test would be needed to clear any doubt. For this reason, we propose to rename S_{6a} as simply S_6 , and S_{6b} as
450 S_{17} ; accordingly, the S -genotypes of 'Mela Rozza', 'Heta' and 'Maikki' are designated as S_3S_6 , $S_{16b}S_{17}$ and
451 $S_{10}S_{17}$, respectively. Among the sequences available in GenBank, AB094495 from 'Oetwiler ReINETte',
452 AB105061 and EU427461 from 'Citron d'Hiver' correspond to S_6 , while AB094493 from 'Bohnappel' and
453 AB105062 from 'Blenheim Orange' match with S_{17} (sequences from Matsumoto et al. 2003; Dreesen et al.
454 2010).

455 The final S -genotype determined for all 63 cultivars is reported on table 4. Among the 25 different alleles
456 detected, S_3 and S_1 were the most frequent, being detected in 17 and 13 cultivars respectively; S_6 , S_{32} and
457 S_{46} on the contrary appeared just in one individual (Fig. 4). All the sequences characterized in this study
458 have been deposited in GenBank and are available under the accession numbers MG598487 (S_1),
459 MG598488 (S_2), MG598489 (S_3), MG598490 (S_4), MG598491 (S_5), MG598492 (S_6), MG598493 (S_7),
460 MG598494 (S_8), MG598495 (S_9), MG598496 (S_{10}), MG598497 (S_{11}), MG598498 (S_{16b}), MG598499 (S_{17}),
461 MG598500 (S_{20}), MG598501 (S_{23}), MG598502 (S_{24}), MG598503 (S_{25}), MG598504 (S_{26}), MG598505 (S_{28}),
462 MG598506 (S_{32}), MG598507 (S_{33}), MG598508 (S_{39}), MG598509 (S_{46}), MG598510 (S_{50}) and MG598511 (S_{58}).

463 **Analysis of promoter sequences**

464 The S -*RNase* gene is known to be specifically expressed in pistil tissues during flower development (Sassa et
465 al. 1993; Ishimizu et al. 1996). The fine regulation of its expression is thought to depend on the presence of
466 cis-acting elements in the promoter, which despite the strong allelic polymorphism must be present in the
467 promoter of all S -*RNases* (Norioka et al. 2000; Dissanayake et al. 2002). In the present study, 5' flanking
468 regions containing at least part of the S -*RNase* promoter have been obtained for 25 different alleles; the
469 complete alignment is provided as Online Resource 6, while the portion including the 800 positions
470 upstream of the start codon is reported in Fig. 5. The sequences have been individually analyzed with
471 TSSPlant (Shahmuradov et al. 2017), the results are reported in Table 5. For 20 out of 25 alleles the most
472 likely transcript start site was found at a position ranging from -90 to -67, with a conserved TATA box at -
473 123 to -100 (Fig. 5); a more distant TSS and TATA box was predicted for the allele S_{23} , while putative TATA-
474 less promoters were identified for S_6 , S_{17} , S_{20} and S_{24} (Table 5). Based on the sequence alignment (Fig. 5),
475 however, the conserved TATA box identified in most alleles seems to be also present in S_6 , S_{17} , S_{20} and S_{24}
476 while a single substitution occurs in the same region of S_{23} . Despite the different predictions by TSSP for
477 these alleles, all the analyzed S -*RNase* promoters might as well share a conserved structure and a common
478 TSS; experimental evidences are needed to test this hypothesis.

479 Compared to solanaceous *S-RNases*, the Pyrinae promoters show a higher sequence conservation;
480 Solanaceae and Rosaceae however share a single conserved motif named IA-like (Ushijima et al. 1998). This
481 motif falls within an approximately 200bp conserved region designated box1, that was described just
482 upstream the TATA box in the promoters of apple *S-RNase* alleles S_1 and S_9 , Japanese pear S_2 , S_3 , S_4 and S_5 ,
483 and Chinese pear S_{12} , S_{13} and S_{21} , but not identified in apple S_2 (Ushijima et al. 1998; Norioka et al. 2000; Liu
484 et al. 2012); our results show however that the S_2 -*RNase* promoter also contains the conserved box1
485 sequence, but it is more distant from the TATA box due to a 350bp insertion (Fig. 5). Box1 is thought to
486 contain cis-acting regulatory elements driving the *S-RNase* expression; consistently with this hypothesis,
487 transformation experiments performed in *Arabidopsis* with a reporter gene (GUS) under control of the
488 Chinese pear S_{12} -*RNase* promoter carrying sequential deletions showed that a truncated box1 region results
489 in the lack of expression in pistils (Liu et al. 2012). Other regulatory elements, however, must be present
490 upstream box1, as sequential deletions resulted in a gradually decreasing reporter gene expression. A
491 second conserved region named box2 has also been reported, but the degree of conservation is lower than
492 box1; moreover, it was not detected in all promoters (Ushijima et al. 1998). Accordingly, sequence identity
493 seems to fall quickly upstream box1 as visible on the alignment in Online Resource 6.

494 The presence of conserved motifs in the 5' flanking region can be exploited to design consensus PCR primers
495 that would allow the amplification of *S-RNase* alleles including portion of the promoter and, most
496 importantly, the entire first exon. Obtaining full-length coding sequences is fundamental not only to better
497 characterize alleles, but also to make reliable comparisons, as highlighted by the case of S_{6a}/S_6 and S_{6b}/S_{17}
498 which were believed to encode the same specificity because only a small portion of the protein-coding
499 sequence was determined. We designed a primer on a short, highly conserved motif within box1, placed a
500 few bases upstream the IA-like motif (Fig. 5); this primer, MdProm-box1f (5'-arggcabtgccatga-3', Online
501 Resource 1), contains two degenerations in order to include all possible variants within this 16bp region.
502 Importantly, its sequence was found not only in all the apple promoters hereby obtained, but also in the
503 Japanese and Chinese pear alleles so far characterized (Ushijima et al. 1998; Norioka et al. 2000; Liu et al.
504 2012). This primer coupled with a consensus reverse designed on exon 2 (PycomC5r) allowed us to amplify
505 and clone alleles S_6 and S_{17} (see section "Validation and sequencing of uncharacterized regions" above). The
506 same primer pair was also successfully tested on several genotypes of European, Japanese and Chinese
507 pear (Online Resource 7) and wild species of *Malus* and *Pyrus* (data not shown), suggesting that it anneals
508 efficiently on the majority, if not all, the promoters of apple and pear *S-RNases*; the possibility to obtain
509 amplification from more distant genera such as *Sorbus*, *Crataegus* and *Eriobotrya* has to be tested. We
510 suggest the use of MdProm-box1f as a quick and cheap alternative to 5' RACE to characterize the 5' portion
511 of *S-RNases*. Unfortunately, the alignment of 3' flanking regions did not display similar clearly conserved
512 motifs (Online Resource 8).

513 **Conclusions**

514 Using a combined alignment and assembly procedure, specifically devised for the analysis of the *S-RNase*
515 gene from resequencing datasets, we were able to determine the *S*-genotype of 63 apple cultivars and to
516 obtain the full-length sequence, including portions of the flanking regions, for 23 *S-RNase* alleles; traditional
517 cloning and sequencing tests confirmed the accuracy of sequence assemblies. Two additional alleles, S_{6a}
518 and S_{6b} , for which only a short portion of the coding sequence was available, were fully sequenced and
519 proved to encode most likely different *S*-specificities; we therefore propose to rename S_{6a} as S_6 and S_{6b} as
520 S_{17} .

521 The availability of full-length *S-RNase* sequences marks an important step forward in the study of Pyrinae *S-*
522 *RNases* and provides essential information for the development of *S*-genotyping assays. Conserved regions
523 were found in all the promoters, making it possible to amplify and clone 5' flanking regions directly from
524 genomic DNA of apple and pear species using consensus primers.

525

526 **Literature cited**

- 527 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool. *J Mol Biol*
528 215:403–410 . doi: 10.1016/S0022-2836(05)80360-2
- 529 Bianco L, Cestaro A, Linsmith G, Muranty H, Denancé C, Théron A, Poncet C, Micheletti D, Kerschbamer E, Di
530 Pierro EA, Larger S, Pindo M, Van De Weg E, Davassi A, Laurens F, Velasco R, Durel CE, Troggio M
531 (2016) Development and validation of the Axiom®Apple480K SNP genotyping array. *Plant J* 86:62–74 .
532 doi: 10.1111/tpj.13145
- 533 Bianco L, Cestaro A, Sargent DJ, Banchi E, Derdak S, Di Guardo M, Salvi S, Jansen J, Viola R, Gut I, Laurens F,
534 Chagné D, Velasco R, Van De Weg E, Troggio M (2014) Development and validation of a 20K Single
535 Nucleotide Polymorphism (SNP) whole genome genotyping array for apple (*Malus × domestica* Borkh).
536 *PLoS One* 9:e110377 . doi: 10.1371/journal.pone.0110377
- 537 Boetzer M, Pirovano W (2012) Toward almost closed genomes with GapFiller. *Genome Biol* 13:R56 . doi:
538 10.1186/gb-2012-13-6-r56
- 539 Bošković R, Tobutt KR (1999) Correlation of stylar ribonuclease isoenzymes with incompatibility alleles in
540 apple. *Euphytica* 107:29–43 . doi: 10.1023/A:1003516902123
- 541 Broothaerts W, Van Neram I, Keulemans J (2004) Update on and review of the incompatibility (*S-*)
542 genotypes of apple cultivars. *HortScience* 39:943–947
- 543 Chagné D, Crowhurst RN, Troggio M, Davey MW, Gilmore B, Lawley C, Vanderzande S, Hellens RP, Kumar S,
544 Cestaro A, Velasco R, Main D, Rees JD, Iezzoni A, Mockler T, Wilhelm L, van de Weg E, Gardiner SE,
545 Bassil N, Peace C (2012) Genome-wide SNP detection, validation, and development of an 8K SNP array
546 for apple. *PLoS One* 7:e31745 . doi: 10.1371/journal.pone.0031745
- 547 Daccord N, Celton JM, Linsmith G, Becker C, Choisne N, Schijlen E, Van De Geest H, Bianco L, Micheletti D,
548 Velasco R, Di Pierro EA, Gouzy J, Rees DJG, Guérif P, Muranty H, Durel CE, Laurens F, Lespinasse Y,
549 Gaillard S, Aubourg S, Quesneville H, Weigel D, Van De Weg E, Troggio M, Bucher E (2017) High-quality
550 de novo assembly of the apple genome and methylome dynamics of early fruit development. *Nat*
551 *Genet* 49:1099–1106 . doi: 10.1038/ng.3886
- 552 De Franceschi P, Cova V, Tartarini S, Dondini L (2016) Characterization of a new apple *S-RNase* allele and its
553 linkage with the *Rvi5* gene for scab resistance. *Mol Breed* 36:1–11 . doi: 10.1007/s11032-015-0427-x
- 554 De Franceschi P, Dondini L, Sanzol J (2012) Molecular bases and evolutionary dynamics of self-
555 incompatibility in the Pyrinae (*Rosaceae*). *J Exp Bot* 63:4015–4032 . doi: 10.1093/jxb/ers108
- 556 De Franceschi P, Pierantoni L, Dondini L, Grandi M, Sansavini S, Sanzol J (2011) Evaluation of candidate *F-*
557 *box* genes for the pollen *S* of gametophytic self-incompatibility in the Pyrinae (*Rosaceae*) on the basis
558 of their phylogenomic context. *Tree Genet Genomes* 7:663–683 . doi: 10.1007/s11295-011-0365-7
- 559 de Nettancourt D (2001) *Incompatibility and Incongruity in Wild and Cultivated Plants*. Springer-Verlag,

- 560 Berlin, Germany
- 561 Dissanayake DMRKK, Norioka S, Norioka N, Takasaki T, Nakanishi T (2002) Cis-regulatory elements for pistil
562 specific expression in S-RNase promoter region of Japanese pear (*Pyrus pyrifolia* Nakai). *Acta Horti*
563 587:459–465 . doi: 10.17660/ActaHortic.2002.587.60
- 564 Dreesen RSG, Vanholme BTM, Luyten K, van Wynsberghe L, Fazio G, Roldán-Ruiz I, Keulemans J (2010)
565 Analysis of *Malus* S-RNase gene diversity based on a comparative study of old and modern apple
566 cultivars and European wild apple. *Mol Breed* 26:693–709 . doi: 10.1007/s11032-010-9405-5
- 567 Evans RC, Campbell CS (2002) The origin of the apple subfamily (Maloideae; Rosaceae) is clarified by DNA
568 sequence data from duplicated GBSSI genes. *Am J Bot* 89:1478–1484 . doi: 10.3732/ajb.89.9.1478
- 569 Garkava-Gustavsson L, Kolodinska Brantestam A, Sehic J, Nybom H (2008) Molecular characterisation of
570 indigenous Swedish apple cultivars based on SSR and S-allele analysis. *Hereditas* 145:99–112 . doi:
571 10.1111/j.0018-0661.2008.02042.x
- 572 Halász J, Hegedus A, György Z, Pállinger É, Tóth M (2011) S-genotyping of old apple cultivars from the
573 Carpathian basin: Methodological, breeding and evolutionary aspects. *Tree Genet Genomes* 7:1135–
574 1145 . doi: 10.1007/s11295-011-0401-7
- 575 Homer N, Merriman B, Nelson SF (2009) BFAST: An alignment tool for large scale genome resequencing.
576 *PLoS One* 4:e7767 . doi: 10.1371/journal.pone.0007767
- 577 Ishimizu T, Endo T, Yamaguchi-Kabata Y, Nakamura KT, Sakiyama F, Norioka S (1998a) Identification of
578 regions in which positive selection may operate in S-RNase of Rosaceae: Implication for S-allele-
579 specific recognition sites in S-RNase. *FEBS Lett* 440:337–342 . doi: 10.1016/S0014-5793(98)01470-7
- 580 Ishimizu T, Sato Y, Saito T, Yoshimura Y, Norioka S, Nakanishi T, Sakiyama F (1996) Identification and partial
581 amino acid sequences of seven S-RNases associated with self-incompatibility of Japanese pear, *Pyrus*
582 *pyrifolia* Nakai. *J Biochem* 120:326–334 . doi: 10.1093/oxfordjournals.jbchem.a021417
- 583 Ishimizu T, Shinkawa T, Sakiyama F, Norioka S (1998b) Primary structural features of rosaceous S-RNases
584 associated with gametophytic self-incompatibility. *Plant Mol Biol* 37:931–941 . doi:
585 10.1023/A:1006078500664
- 586 Janssens GA, Goderis IJ, Broekaert WF, Broothaerts W (1995) A molecular method for S-allele identification
587 in apple based on allele-specific PCR. *Theor Appl Genet* 91:691–698 . doi: 10.1007/BF00223298
- 588 Kim H, Kakui H, Kotoda N, Hirata Y, Koba T, Sassa H (2009) Determination of partial genomic sequences and
589 development of a CAPS system of the S-RNase gene for the identification of 22 S haplotypes of apple
590 (*Malus × domestica* Borkh.). *Mol Breed* 23:463–472 . doi: 10.1007/s11032-008-9249-4
- 591 Kim HT, Moriya S, Okada K, Abe K, Park JI, Yamamoto T, Nou S (2016) Identification and characterization of
592 S-RNase genes in apple rootstock and the diversity of S-RNases in *Malus* species. *J Plant Biotechnol*
593 43:49–57 . doi: 10.5010/JPB.2016.43.1.49
- 594 Kitahara K, Matsumoto S (2002a) Cloning of the S25cDNA from “McIntosh” apple and an S25-allele
595 identification method. *J Hortic Sci Biotechnol* 77:724–728 . doi: 10.1080/14620316.2002.11511563
- 596 Kitahara K, Matsumoto S (2002b) Sequence of the S10 cDNA from “McIntosh” apple and a PCR-digestion
597 identification method. *HortScience* 37:187–190
- 598 Kitahara K, Soejima J, Komatsu H, Fukui H, Matsumoto S (2000) Complete sequences of the S-genes, Sd-

- 599 and Sh-RNase cDNA in apple. *HortScience* 35:712–715
- 600 Kubo KI, Entani T, Takara A, Wang N, Fields AM, Hua Z, Toyoda M, Kawashima SI, Ando T, Isogai A, Kao TH,
601 Takayama S (2010) Collaborative non-self recognition system in S-RNase-based self-incompatibility.
602 *Science* (80-) 330:796–799 . doi: 10.1126/science.1195243
- 603 Kubo KI, Paape T, Hatakeyama M, Entani T, Takara A, Kajihara K, Tsukahara M, Shimizu-Inatsugi R, Shimizu
604 KK, Takayama S (2015) Gene duplication and genetic exchange drive the evolution of S-RNase-based
605 self-incompatibility in *Petunia*. *Nat Plants* 1:140005 . doi: 10.1038/nplants.2014.5
- 606 Larkin MA, Blackshields G, Brown NP, Chenna R, Mcgettigan PA, McWilliam H, Valentin F, Wallace IM, Wilm
607 A, Lopez R, Thompson JD, Gibson TJ, Higgins DG (2007) Clustal W and Clustal X version 2.0.
608 *Bioinformatics* 23:2947–2948 . doi: 10.1093/bioinformatics/btm404
- 609 Larsen B, Ørgaard M, Toldam-Andersen TB, Pedersen C (2016) A high-throughput method for genotyping S-
610 RNase alleles in apple. *Mol Breed* 36:1–10 . doi: 10.1007/s11032-016-0448-0
- 611 Lassmann T, Sonnhammer ELL (2005) Kalign - An accurate and fast multiple sequence alignment algorithm.
612 *BMC Bioinformatics* 6:298 . doi: 10.1186/1471-2105-6-298
- 613 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R (2009) The
614 Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25:2078–2079 . doi:
615 10.1093/bioinformatics/btp352
- 616 Li T, Long S, Li M, Bai S, Zhang W (2012) Determination S-Genotypes and Identification of Five Novel S-
617 RNase Alleles in Wild *Malus* Species. *Plant Mol Biol Report* 30:453–461 . doi: 10.1007/s11105-011-
618 0345-y
- 619 Liu X ying, Wuyun T na, Zeng H yan (2012) Cloning, characterization and promoter analysis of S-RNase gene
620 promoter from Chinese pear (*Pyrus pyrifolia*). *Gene* 505:246–253 . doi: 10.1016/j.gene.2012.06.017
- 621 Long S, Li M, Han Z, Wang K, Li T (2010a) Characterization of three new S-alleles and development of an S-
622 allele-specific PCR system for rapidly identifying the S-genotype in apple cultivars. *Tree Genet*
623 *Genomes* 6:161–168 . doi: 10.1007/s11295-009-0237-6
- 624 Long S, Li M, Han Z, Zhang B, Wang K, Li T (2010b) Characterization of two novel S-RNase genes and PCR
625 analyzing of S-genotypes of 46 cultivars in *Malus domestica* Borkh. *J Agric Biotechnol* 18:265–271
- 626 Ma RC, Oliveira M (2002) Evolutionary analysis of S-RNase genes from Rosaceae species. *Mol Genet*
627 *Genomics* 267:71–78 . doi: 10.1007/s00438-002-0637-x
- 628 Maguire TL, Collins GG, Sedgley M (1994) A modified CTAB DNA extraction procedure for plants belonging
629 to the family proteaceae. *Plant Mol Biol Report* 12:106–109 . doi: 10.1007/BF02668371
- 630 Matsumoto S (2014) Apple pollination biology for stable and novel fruit production: Search system for
631 apple cultivar combination showing incompatibility, semicompatibility, and full-compatibility based on
632 the S-RNase allele database. *Int J Agron* 2014:1–9 . doi: 10.1155/2014/138271
- 633 Matsumoto S, Furusawa Y (2005) Genomic DNA sequence of S16c (= 16)-RNase in apple: re-numbering of
634 S16 (= 27a)-and S22 (= 27b)-allele to S16a and S16b. *Sci Rep Rac Educ Gifu Univ (Nat Sci)* 29:7–12
- 635 Matsumoto S, Furusawa Y, Kitahara K, Soejima J (2003) Partial Genomic Sequences of S6-, S12-, S13-, S14-,
636 S17-, S19-, and S21-RNases of Apple and Their Allele Designations. *Plant Biotechnol* 20:323–329 . doi:
637 10.5511/plantbiotechnology.20.323

- 638 Matsumoto S, Kitahara K, Komatsu H, Abe K (2006) Cross-compatibility of apple cultivars possessing S-
639 RNase alleles of similar sequence. *J Hortic Sci Biotechnol* 81:934–936 . doi:
640 10.1080/14620316.2006.11512178
- 641 Matsumoto S, Komori S, Kitahara K, Imazu S, Soejima J (1999) S genotypes of 15 apple cultivars and self-
642 compatibility of “Megumi.” *J Japan Soc Hort Sci* 68:236–241
- 643 Milne I, Stephen G, Bayer M, Cock PJA, Pritchard L, Cardle L, Shawand PD, Marshall D (2013) Using tablet for
644 visual exploration of second-generation sequencing data. *Brief Bioinform* 14:193–202 . doi:
645 10.1093/bib/bbs012
- 646 Minamikawa M, Kakui H, Wang S, Kotoda N, Kikuchi S, Koba T, Sassa H (2010) Apple S locus region
647 represents a large cluster of related, polymorphic and pollen-specific F-box genes. *Plant Mol Biol*
648 74:143–154 . doi: 10.1007/s11103-010-9662-z
- 649 Morita J, Abe K, Matsumoto S (2009) S-RNase genotypes of apple cultivars grown in Japan and
650 development of A PCR-RFLP method to identify the S 6-and S 21-RNase alleles. *J Hortic Sci Biotechnol*
651 84:29–34 . doi: 10.1080/14620316.2009.11512475
- 652 Norioka N, Katayama H, Matsuki T, Ishimizu T, Takasaki T, Nakanishi T, Norioka S (2000) Sequence
653 comparison of the 5' flanking regions of Japanese pear (*Pyrus pyrifolia*) S-RNases associated with
654 gametophytic self-incompatibility. *Sex Plant Reprod* 13:289–291 . doi: 10.1007/s004970100067
- 655 Nybom H, Sehic J, Garkava-Gustavsson L (2008) Self-incompatibility alleles of 104 apple cultivars grown in
656 northern Europe. *J Hortic Sci Biotechnol* 83:339–344 . doi: 10.1080/14620316.2008.11512389
- 657 Okada K, Moriya S, Haji T, Abe K (2013) Isolation and characterization of multiple F-box genes linked to the
658 S9- and S10-RNase in apple (*Malus × domestica* Borkh.). *Plant Reprod* 26:101–111 . doi:
659 10.1007/s00497-013-0212-0
- 660 Okada K, Tonaka N, Taguchi T, Ichikawa T, Sawamura Y, Nakanishi T, Takasaki-Yasuda T (2011) Related
661 polymorphic F-box protein genes between haplotypes clustering in the BAC contig sequences around
662 the S-RNase of Japanese pear. *J Exp Bot* 62:1887–1902 . doi: 10.1093/jxb/erq381
- 663 Potter D, Eriksson T, Evans RC, Oh S, Smedmark JEE, Morgan DR, Kerr M, Robertson KR, Arsenault M,
664 Dickinson TA, Campbell CS (2007) Phylogeny and classification of Rosaceae. *Plant Syst Evol* 266:5–43 .
665 doi: 10.1007/s00606-007-0539-9
- 666 Raspé O, Kohn JR (2002) S-allele diversity in *Sorbus aucuparia* and *Crataegus monogyna* (Rosaceae:
667 Maloideae). *Heredity (Edinb)* 88:458–65 . doi: 10.1038/sj/hdy/6800079
- 668 Sambrook J, Fristch EF, Maniatis T (1989) *Molecular cloning: a laboratory manual*, 2nd Ed. Cold Spring Harb
669 Lab Press Cold Spring Harb New York
- 670 Sanzol J (2010) Two neutral variants segregating at the gametophytic self-incompatibility locus of European
671 pear (*Pyrus communis* L.) (Rosaceae, Pyrinae). *Plant Biol* 12:800–805 . doi: 10.1111/j.1438-
672 8677.2009.00277.x
- 673 Sanzol J (2009a) Genomic characterization of self-incompatibility ribonucleases (S-RNases) in European
674 pear cultivars and development of PCR detection for 20 alleles. *Tree Genet Genomes* 5:393–405 . doi:
675 10.1007/s11295-008-0194-5
- 676 Sanzol J (2009b) Pistil-function breakdown in a new S-allele of European pear, S 21 ??, confers self-
677 compatibility. *Plant Cell Rep* 28:457–467 . doi: 10.1007/s00299-008-0645-3

- 678 Sassa H, Hirano H, Ikehashi H (1993) Identification and characterization of stylar glycoproteins associated
679 with self-incompatibility genes of Japanese pear, *Pyrus serotina* Rehd. *MGG Mol Gen Genet* 241:17–
680 25 . doi: 10.1007/BF00280196
- 681 Sassa H, Kakui H, Miyamoto M, Suzuki Y, Hanada T, Ushijima K, Kusaba M, Hirano H, Koba T (2007) S locus
682 F-box brothers: Multiple and pollen-specific F-box genes with S haplotype-specific polymorphisms in
683 apple and Japanese pear. *Genetics* 175:1869–1881 . doi: 10.1534/genetics.106.068858
- 684 Sassa H, Nishio T, Kowyama Y, Hirano H, Koba T, Ikehashi H (1996) Self-incompatibility (S) alleles of the
685 Rosaceae encode members of a distinct class of the T2/S ribonuclease superfamily. *Mol Gen Genet*
686 250:547–557 . doi: 10.1007/s004380050108
- 687 Shahmuradov IA, Umarov RK, Solovyev V V. (2017) TSSPlant: A new tool for prediction of plant Pol II
688 promoters. *Nucleic Acids Res* 45: . doi: 10.1093/nar/gkw1353
- 689 Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H, Remmert M, Söding J,
690 Thompson JD, Higgins DG (2011) Fast, scalable generation of high-quality protein multiple sequence
691 alignments using Clustal Omega. *Mol Syst Biol* 7:539 . doi: 10.1038/msb.2011.75
- 692 Steinbachs JE, Holsinger KE (2002) S-RNase-mediated gametophytic self-incompatibility is ancestral in
693 eudicots. *Mol Biol Evol* 19:825–829 . doi: 10.1093/oxfordjournals.molbev.a004139
- 694 Takasaki-Yasuda T, Nomura N, Moriya-Tanaka Y, Okada K, Iwanami H, Bessho H (2013) Cloning an S-RNase
695 allele, including the longest intron, from cultivars of European pear (*Pyrus communis* L.). *J Hortic Sci*
696 *Biotechnol* 88:427–432 . doi: 10.1080/14620316.2013.11512987
- 697 Tamura K, Stecher G, Peterson D, Filipinski A, Kumar S (2013) MEGA6: Molecular evolutionary genetics
698 analysis version 6.0. *Mol Biol Evol* 30:2725–2729 . doi: 10.1093/molbev/mst197
- 699 Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, Rozen SG (2012) Primer3-new
700 capabilities and interfaces. *Nucleic Acids Res* 40:e115 . doi: 10.1093/nar/gks596
- 701 Ushijima K, Sassa H, Hirano H (1998) Characterization of the flanking regions of the S-RNase genes of
702 Japanese pear (*Pyrus serotina*) and apple (*Malus x domestica*). *Gene* 211:159–167 . doi:
703 10.1016/S0378-1119(98)00105-X
- 704 Van Nerum I, Geerts M, Van Haute a., Keulemans J, Broothaerts W, Nerum I van, Geerts M, Haute A van,
705 Keulemans J, Broothaerts W (2001) Re-examination of the self-incompatibility genotype of apple
706 cultivars containing putative “new” S-alleles. *Theor Appl Genet* 103:584–591 . doi:
707 10.1007/PL00002913
- 708 Vavilov NI (1930) Wild Progenitors of the fruit trees of Turkistan and Cacausus and the problem of the
709 origin of fruit trees. *Report Proc 9th Int Hort Congr* 271–286
- 710 Velasco R, Zharkikh A, Affourtit J, Dhingra A, Cestaro A, Kalyanaraman A, Fontana P, Bhatnagar SK, Troggio
711 M, Pruss D, Salvi S, Pindo M, Baldi P, Castelletti S, Cavaiuolo M, Coppola G, Costa F, Cova V, Dal Ri A,
712 Goremykin V, Komjanc M, Longhi S, Magnago P, Malacarne G, Malnoy M, Micheletti D, Moretto M,
713 Perazzolli M, Si-Ammour A, Vezzulli S, Zini E, Eldredge G, Fitzgerald LM, Gutin N, Lanchbury J, Macalma
714 T, Mitchell JT, Reid J, Wardell B, Kodira C, Chen Z, Desany B, Niazi F, Palmer M, Koepke T, Jiwan D,
715 Schaeffer S, Krishnan V, Wu C, Chu VT, King ST, Vick J, Tao Q, Mraz A, Stormo A, Stormo K, Bogden R,
716 Ederle D, Stella A, Vecchiatti A, Kater MM, Masiero S, Lasserre P, Lespinasse Y, Allan AC, Bus V, Chagné
717 D, Crowhurst RN, Gleave AP, Lavezzo E, Fawcett J a, Proost S, Rouzé P, Sterck L, Toppo S, Lazzari B,
718 Hellens RP, Durel C-E, Gutin A, Bumgarner RE, Gardiner SE, Skolnick M, Egholm M, Van de Peer Y,
719 Salamini F, Viola R (2010) The genome of the domesticated apple (*Malus x domestica* Borkh.). *Nat*

- 720 Genet 42:833–839 . doi: 10.1038/ng.654
- 721 Vieira J, Ferreira PG, Aguiar B, Fonseca NA, Vieira CP (2010) Evolutionary patterns at the RNase based
722 gametophytic self - incompatibility system in two divergent Rosaceae groups (Maloideae and Prunus).
723 BMC Evol Biol 10:200 . doi: 10.1186/1471-2148-10-200
- 724 Vieira J, Morales-Hojas R, Santos RAMM, Vieira CP (2007) Different positively selected sites at the
725 gametophytic self-incompatibility pistil S-RNase gene in the Solanaceae and Rosaceae (Prunus, Pyrus,
726 and Malus). J Mol Evol 65:175–185 . doi: 10.1007/s00239-006-0285-6
- 727 Wang S, Kakui H, Kikuchi S, Koba T, Sassa H (2012) Interhaplotypic heterogeneity and heterochromatic
728 features may contribute to recombination suppression at the S locus in apple (*Malus x domestica*). J
729 Exp Bot 63:4983–4990 . doi: 10.1093/jxb/ers176
- 730 Zamboni A, Pierantoni L, De Franceschi P (2008) Total RNA extraction from strawberry tree (*Arbutus unedo*)
731 and several other woody-plants. iForest - Biogeosciences For 1:122–125 . doi: 10.3832/ifor0465-
732 0010122
- 733 Zisovich AH, Stern RA, Sapir G, Shafir S, Goldway M (2004) The RHV region of S-RNase in the European pear
734 (*Pyrus communis*) is not required for the determination of specific pollen rejection. Sex Plant Reprod
735 17:151–156 . doi: 10.1007/s00497-004-0225-9
- 736

737 **Tables**

738 **Table 1.** List and accession numbers of the *S-RNase* sequences used for building the synthetic alignment
 739 template; alleles from 1 to 34 follow the numeration reported by Matsumoto (2014), while other alleles
 740 were described by Long et al. 2010b; Long et al. 2010a; Dreesen et al. 2010; Li et al. 2012; and De
 741 Franceschi et al. 2016

Allele	Acc. #	Allele	Acc. #	Allele	Acc. #	Allele	Acc. #
1	D50837 ^b EU427454 ^c	9	AB270792	25	AB062100 ^b AB428431 ^c	34	AB540122
2	U12199 ^b HQ693077 ^c	10	AF327221 ^b AB428428 ^c	26	AF016918 ^b AB428432 ^c	41	KT724706
3	U12200 ^b EU427455 ^c	11 ^a	FJ008669 ^c	28	AB035273 ^b AF201748 ^c	42 ^a	EU427453 ^c
4	AF327223 ^b EU427456 ^c	16	AF016919 ^b AB428429 ^c	29 ^a	AY039702 ^c	44 ^a	EU443101 ^c FJ008673 ^c
5 ^a	U19791 ^b AB428427 ^c	20	AB019184 ^b HQ689397 ^c	30	AB035928 ^b AB052268 ^c	45 ^a	FJ008671 ^c
6 ^a	EU427461 ^c	21 ^a	FJ008670 ^c	31	DQ135990	46 ^a	FJ008672 ^c
7	AB032246 ^b AB050634 ^c	23 ^a	AF239809 ^c	32	DQ135991	50 ^a	FJ535241 ^c
8 ^a	AY744080 ^b	24	AF016920 ^b HQ693064 ^c	33	AB540121	53 ^a	FJ602074 ^c

^a full-length coding sequence is not available

^b mRNA

^c partial genomic sequence

742

743

744 **Table 2.** Summary of the sequences obtained by the assembly of reads; the size (bp) of flanking regions,
745 exons and the intron are reported. The number between parentheses indicates the portion of sequence
746 known before the assembly.

Allele	5' flanking	Exon 1	Intron	Exon 2	3' flanking
<i>S</i> ₁	1167	243	344	438	1346
<i>S</i> ₂	877	252	147	435	1053
<i>S</i> ₃	747	252	1290	435	1108
<i>S</i> ₄	912	244	140	440	44
<i>S</i> ₅	536	252 (159)	1158	438	754
<i>S</i> _{6a}	0	246 (138)	167	392 (26)	0
<i>S</i> _{6b}	0	242 (138)	167	34 (26)	0
<i>S</i> ₇	571	252	118	435	866
<i>S</i> ₈	1077	252 (171)	1232 (0)	435	1102
<i>S</i> ₉	1084	252	144	435	85
<i>S</i> ₁₀	766	252	1717	435	1152
<i>S</i> ₁₁	363	246 (150)	175	438 (384)	506
<i>S</i> _{16b}	1275	246	2130	438	1559
<i>S</i> ₂₀	750	243	318	438	1437
<i>S</i> ₂₃	882	252 (132)	147	435 (324)	198
<i>S</i> ₂₄	1013	243	339	438	1361
<i>S</i> ₂₅	1323	252	2359	435	251
<i>S</i> ₂₆	1165	252	159	432	1057
<i>S</i> ₂₈	813	252	169	432	1011
<i>S</i> ₃₂	901	252	151	435	12
<i>S</i> ₃₃	990	261	726	438	68
<i>S</i> ₄₆	541	252 (156)	143	435 (381)	42
<i>S</i> ₅₀	14	252 (156)	162	429 (381)	1028

747

748

Table 3. Estimates of evolutionary divergence between *S-RNase* sequences. The distances calculated as number of amino acid substitutions per site using the Poisson correction model are reported in the lower-left part, while the upper-right shows the number of amino acid differences per sequence.

	S_1	S_2	S_3	S_4	S_5	$S_{6a} (S_6)$	$S_{6b} (S_{17})$	S_7	S_8	S_9	S_{10}	S_{11}	S_{16b}	S_{20}	S_{23}	S_{24}	S_{25}	S_{26}	S_{28}	S_{32}	S_{33}	S_{39}	S_{46}	S_{50}	S_{58}
S_1	-	82	68	72	81	68	69	78	85	70	67	79	77	14	83	17	77	71	67	79	77	89	73	71	73
S_2	0.455	-	76	77	87	76	78	66	87	67	76	76	82	83	86	83	86	70	77	81	78	92	66	71	78
S_3	0.361	0.405	-	76	74	68	68	70	82	66	14	77	76	66	74	69	71	78	71	76	80	79	69	71	72
S_4	0.389	0.414	0.407	-	84	75	74	80	81	75	76	62	36	75	81	74	74	66	70	76	78	83	76	66	69
S_5	0.451	0.483	0.394	0.462	-	83	84	73	89	82	73	85	83	82	89	84	84	90	83	84	50	45	85	85	85
$S_{6a} (S_6)$	0.359	0.407	0.356	0.403	0.457	-	14	74	81	67	68	83	80	70	81	72	73	71	30	74	74	80	68	62	69
$S_{6b} (S_{17})$	0.366	0.421	0.356	0.396	0.464	0.063	-	74	78	69	68	85	80	71	84	75	74	71	28	72	76	81	71	66	71
S_7	0.430	0.343	0.368	0.434	0.385	0.396	0.396	-	86	58	70	81	82	76	79	79	74	69	71	77	66	79	61	70	75
S_8	0.477	0.480	0.445	0.441	0.497	0.441	0.421	0.476	-	76	78	87	85	83	79	85	32	88	77	66	84	89	78	78	87
S_9	0.374	0.347	0.341	0.401	0.448	0.349	0.362	0.295	0.405	-	69	79	70	67	78	70	67	61	65	64	74	86	16	62	65
S_{10}	0.355	0.405	0.063	0.407	0.387	0.356	0.356	0.368	0.418	0.360	-	79	75	65	73	67	66	76	67	76	75	75	70	67	70
S_{11}	0.437	0.407	0.414	0.319	0.469	0.457	0.471	0.441	0.483	0.427	0.427	-	57	81	85	82	80	75	83	75	80	88	79	73	75
S_{16b}	0.423	0.448	0.407	0.172	0.455	0.436	0.436	0.448	0.469	0.368	0.401	0.289	-	79	84	77	76	70	74	77	82	84	73	67	71
S_{20}	0.063	0.462	0.349	0.409	0.458	0.372	0.379	0.416	0.462	0.355	0.342	0.451	0.437	-	80	15	76	72	67	78	77	86	71	73	74
S_{23}	0.462	0.473	0.392	0.441	0.497	0.441	0.462	0.427	0.425	0.418	0.385	0.469	0.462	0.441	-	84	75	79	78	84	78	89	83	77	85
S_{24}	0.078	0.462	0.368	0.403	0.472	0.385	0.405	0.437	0.477	0.374	0.355	0.458	0.423	0.068	0.470	-	78	73	69	80	78	86	73	74	74
S_{25}	0.421	0.473	0.373	0.394	0.462	0.387	0.394	0.394	0.151	0.347	0.341	0.434	0.407	0.414	0.398	0.428	-	77	70	55	77	82	70	69	76
S_{26}	0.381	0.368	0.421	0.345	0.507	0.375	0.375	0.364	0.490	0.312	0.407	0.403	0.370	0.387	0.427	0.394	0.414	-	65	65	79	90	62	60	69
S_{28}	0.355	0.416	0.377	0.372	0.460	0.141	0.131	0.379	0.416	0.339	0.351	0.460	0.398	0.355	0.423	0.368	0.370	0.339	-	70	73	83	69	60	68
S_{32}	0.434	0.438	0.405	0.407	0.462	0.394	0.381	0.414	0.341	0.329	0.405	0.401	0.414	0.428	0.459	0.441	0.276	0.337	0.370	-	75	82	64	72	81
S_{33}	0.421	0.418	0.432	0.421	0.246	0.394	0.407	0.341	0.459	0.392	0.398	0.434	0.448	0.421	0.418	0.428	0.412	0.427	0.390	0.398	-	47	74	79	85
S_{39}	0.509	0.519	0.427	0.455	0.218	0.436	0.443	0.425	0.497	0.476	0.401	0.490	0.462	0.487	0.497	0.487	0.448	0.507	0.460	0.448	0.229	-	89	90	93
S_{46}	0.394	0.341	0.360	0.407	0.469	0.356	0.375	0.312	0.418	0.072	0.366	0.427	0.387	0.381	0.452	0.394	0.366	0.319	0.364	0.329	0.392	0.497	-	65	68
S_{50}	0.383	0.377	0.377	0.347	0.474	0.320	0.345	0.372	0.423	0.320	0.351	0.392	0.353	0.396	0.416	0.403	0.364	0.308	0.310	0.383	0.430	0.510	0.339	-	30

S_{58}	0.394	0.421	0.381	0.364	0.471	0.362	0.375	0.403	0.483	0.337	0.368	0.403	0.377	0.401	0.469	0.401	0.407	0.362	0.357	0.441	0.469	0.530	0.356	0.142	-
----------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	---

Table 4. Summary of the S-genotypes determined for the 63 apple accessions analyzed in this study.

Cultivar	S genotype	Cultivar	S genotype	Cultivar	S genotype
Abbondanza	S ₃ S ₅	Delicious ^a	S ₉ S ₂₈	Malinové Holovouské	S ₃ S ₂₀
Ag Alma	S ₂₃ S ₅₀	Dr. Oldenburg	S ₃ S _{16b}	Mcintosh ^a	S ₁₀ S ₂₅
Aivaniya	S ₁ S ₂₅	Durello Di Forlì	S ₃ S ₅₀	Mela Rosa (Pd)	S ₃ S ₂₅
Ajmi	S ₈ S ₂₃	F22682922	S ₅ S ₂₄	Mela Rozza	S ₃ S ₆
Åkerö ^a	S ₁ S ₅₈	Filippa ^a	S ₇ S ₂₄	Ovčí Hubička	S ₃ S ₂₈
Alfred Jolibois	S ₁ S ₂₈	Fuji ^a	S ₁ S ₉	Panenskéé České	S ₇ S ₁₀
Amadou	S ₂₀ S ₃₃	Fyriki	S _{16b} S ₂₆	Papirovka	S ₁ S ₅
Annurca	S ₇ S ₂₆	Gelata	S ₇ S _{16b}	Patte De Loup	S ₁ S ₂
Antonovka	S ₈ S ₃₂	Godelieve Hegmans	S ₂ S ₇	Pepino Jaune	S ₁ S ₃
Antonovka Pamtorutka	S ₈ S ₅₈	Heta	S _{16b} S ₁₇	Precoce De Karage	S ₄ S ₂₆
Aport Kuba	S ₁ S ₅₀	Hetlina	S ₁ S _{16b}	President Roulin	S ₅ S ₂₄
Belle Et Bonne	S ₃ S ₃₃	Ijunscoe Ranee	S ₃₉ S ₄₆	Priscilla ^b	S ₇ S ₁₀
Borowitzky ^b	S ₃ S ₂₈	Jantarnoe	S ₈ S ₂₆	Reinette Clochard	S ₃ S ₄
Braeburn ^a	S ₉ S ₂₄	Jonathan ^a	S ₇ S ₉	Reinette Dubois	S ₃ S ₁₁
Budimka	S ₁ S _{16b}	Keswick Codlin	S ₄ S ₂₀	Renetta Grigia Torriana	S ₇ S ₂₈
Busiard	S ₁ S ₅	Kmenotvorná	S ₃ S ₇	Rosa (Fi)	S ₂₀ S ₂₈
Cabarette	S ₃ S ₂₄	Košíkové	S ₃ S ₈	Skry	S ₂₄ S ₂₈
Chodské	S ₂ S ₁₀	Kronprins	S ₁₁ S ₂₈	Sonderskow	S ₁₀ S ₃₃
Court-Pendu Henry	S ₃ S ₅	Lady Williams ^a	S ₇ S ₂₃	Spässerud ^a	S ₁ S ₇
Cox's Orange Pippin ^a	S ₅ S ₉	Macoun ^a	S ₃ S ₂₅	Worcester Pearmain ^b	S ₂₄ S ₂₅
De L'Estre	S ₁ S ₂	Maikki	S ₁₀ S ₁₇	Young America	S ₂₈ S ₃₉

^a S-genotypes in agreement with previous reports (Janssens et al. 1995; Sassa et al. 1996; Matsumoto et al. 1999; Van Nerum et al. 2001; Kitahara and Matsumoto 2002a; Kitahara and Matsumoto 2002b; Broothaerts et al. 2004; Garkava-Gustavsson et al. 2008; Nybom et al. 2008; Matsumoto 2014)

^b S-genotypes in disagreement with previous reports (Kitahara et al. 2000; Broothaerts et al. 2004; Morita et al. 2009; Long et al. 2010a)

Table 5. Prediction of transcript start site (TSS) and TATA box position in the promoter region of *S-RNases*.

Allele	TSS position	TSS score	TATA position	TATA score
<i>S</i> ₁	-78	1.988	-111	8.193
<i>S</i> ₂	-83	1.988	-116	8.193
<i>S</i> ₃	-78	1.988	-111	5.506
<i>S</i> ₄	-78	1.988	-111	7.781
<i>S</i> ₅	-77	1.989	-111	8.193
<i>S</i> ₆	-50	1.888	TATA-	
<i>S</i> ₇	-68	1.988	-102	8.193
<i>S</i> ₈	-67	1.986	-100	6.901
<i>S</i> ₉	-73	1.988	-105	8.193
<i>S</i> ₁₀	-78	1.988	-111	5.506
<i>S</i> ₁₁	-72	1.989	-105	8.193
<i>S</i> _{16b}	-78	1.989	-111	8.193
<i>S</i> ₁₇	-50	1.887	TATA-	
<i>S</i> ₂₀	-322	1.991	TATA-	
<i>S</i> ₂₃	-199	1.968	-230	4.240
<i>S</i> ₂₄	-333	1.993	TATA-	
<i>S</i> ₂₅	-77	1.988	-111	8.193
<i>S</i> ₂₆	-74	1.989	-108	8.193
<i>S</i> ₂₈	-90	1.989	-123	8.193
<i>S</i> ₃₂	-78	1.986	-111	6.972
<i>S</i> ₃₃	-71	1.989	-104	8.193
<i>S</i> ₃₉	-77	1.989	-111	8.193
<i>S</i> ₄₆	-79	1.988	-113	8.193
<i>S</i> ₅₀	-78	1.988	-111	8.193
<i>S</i> ₅₈	-76	1.988	-109	8.193

Figure captions

Fig. 1 Alignment of the newly assembled flanking regions of the *S₉-RNase* to the genomic sequence. The varying length obtained at the 5' side (a) after 10 iterations depended essentially on the different coverage obtained from the tested cultivars, while at the 3' side (b) assembly of reads was hindered by the presence of a simple sequence repeat

Fig. 2 Alignment of deduced protein sequences of the 25 *S-RNase* alleles detected in the pool of resequenced cultivars. Black and grey backgrounds mark conserved sites and conservative substitutions, respectively; the five conserved (C1-C5) and one hyper-variable (RHV) regions of Rosaceous *S-RNases* are underlined

Fig. 3 Phylogenetic tree built on the alignment of *S-RNase* protein sequences using the Neighbor-Joining method; the percentage of replicate trees in which the associated sequences clustered together in the bootstrap test (1000 replicates) are shown next to the branches. The evolutionary distances were computed using the Poisson correction method and are expressed as the number of amino acid substitutions per site

Fig. 4 Bar plot of allele frequencies among the 63 cultivars

Fig. 5 Alignment of *S-RNase* promoters. Transcript start site (TSS) and TATA box were predicted using TSSPlant (see text); the conserved region box1, the motif IA-like and the annealing site of primer MdProm-box1f are indicated

Electronic Supplementary Material

Online Resource 1. List and sequence of primers used in this study.

Online Resource 2. Coverage obtained from the alignment of Illumina reads on the synthetic sequence.

Online Resource 3. Verification of *S-RNase* alleles detected in 'Borowitsky', 'Priscilla' and 'Worcester Pearmain'.

Online Resource 4. Alignment of *Malus* hybrid 'Mt Blanc' *S₅₈* with *M. × domestica S₅₈-RNase* and *M. sylvestris S₃₉* with *M. × domestica S₃₉-RNase*.

Online Resource 5. Alignment of sequences derived from assembly of reads and traditional cloning and sequencing.

Online Resource 6. Alignment of *S-RNase* 5' flanking regions.

Online Resource 7. Test amplification with primer MdProm-box1f from a set of cultivars of different pear species.

Online Resource 8. Alignment of *S-RNase* 3' flanking regions.