Title: Assessment of population differentiation and linkage disequilibrium in

*Solanum pimpinellifolium* using genome-wide high-density SNP markers

Authors: Ya-Ping Lin, Chu-Yin Liu, Kai-Yi Chen

Institutional affiliations:

Department of Agronomy, National Taiwan University, Taipei, Taiwan, 10617

Short running title: RADseq in *S. pimpinellifolium*

Keywords: *Solanum pimpinellifolium*, population differentiation, linkage disequilibrium pattern, restriction site associated DNA sequencing.

Corresponding author: Kai-Yi Chen

Office mailing address: No.1, Sec. 4, Roosevelt Rd., Da'an Dist., Department of Agronomy, National Taiwan University, Taipei City, Taiwan, 10617

Phone number: +886 2 3366 4766

Email address: kaiychen@ntu.edu.tw

1 **ABSTRACT**

2      To mine new favorable alleles for tomato breeding, we investigated the

3 feasibility of utilizing *Solanum pimpinellifolium* as a diverse panel of

4 genome-wide association study through the restriction site-associated DNA

5 sequencing technique. Previous attempts to conduct genome-wide association

6 study using *S. pimpinellifolium* were impeded by an inability to correct for

7 population stratification and by lack of high-density markers to address the

8 issue of rapid linkage disequilibrium decay. In the current study, a set of

9 24,330 SNPs was identified using 99 *S. pimpinellifolium* accessions from the

10 Tomato Genetic Resource Center. Approximately 84% *Pst*I site-associated

11 DNA sequencing regions were located in the euchromatic regions, resulting in

12 the tagging of most SNPs on or near genes. Our genotypic data suggested

13 that the optimum number of *S. pimpinellifolium* ancestral subpopulations was

14 three, and accessions were classified into seven groups. In contrast to the

15 SolCAP SNP genotypic data of previous studies, our SNP genotypic data

16 consistently confirmed the population differentiation, achieving a relatively

17 uniform correction of population stratification. Moreover, as expected, rapid

18 linkage disequilibrium decay was observed in *S. pimpinellifolium*, especially in

19 euchromatic regions. Approximately two-thirds of the flanking SNP markers

20 did not display linkage disequilibrium. Our result suggests that higher density

21 of molecular markers and more accessions are required to conduct the

22 genome-wide association study utilizing the *Solanum pimpinellifolium*

23 collection.

24          **INTRODUCTION**

25          The wild tomato species *Solanum pimpinellifolium* is a native perennial

26   shrub in Ecuador and Peru and is believed to have originated in northern Peru

27   and then diversified into several subpopulations after it migrated to Ecuador

28   and southern Peru (Rick *et al.* 1977; Zuriaga *et al.* 2009; Blanca *et al.* 2012,

29   2015). The accessions in northern Peru display higher genetic variation and a

30   higher outcrossing rate than those in southern Peru (Rick *et al.* 1977; Caicedo

31   and Schaal 2004). Recent studies suggested that *S. pimpinellifolium* can be

32   divided into at least three subpopulations: one in Peru, one in northern

33   Ecuador and one in the mountains of Ecuador (Zuriaga *et al.* 2009; Blanca *et*

34   *al.* 2012, 2015). In addition, the major climatic parameters, such as

35   temperature and precipitation, show unidirectional gradient changes from

36   southern Peru towards Ecuador (Zuriaga *et al.* 2009). Because geographic

37   distributions of distinct *S. pimpinellifolium* subpopulations also aligned from

38   south to north, it was proposed that the genetic distances between

39   subpopulations were correlated with climatic differences (Zuriaga *et al.* 2009;

40   Blanca *et al.* 2012, 2015).

41          *S. pimpinellifolium* is an attractive resource for tomato breeding because

42   it can freely cross with cultivated tomatoes and introduces novel alleles into the

43   limited gene pool of cultivated tomatoes (Tanksley and Mccouch 1997;

44   Spooner *et al.* 2005; Moyle 2008). *S. pimpinellifolium* has been used as a

45   genetic resource for disease resistance and fruit quality traits in tomato

46   breeding (Grandillo *et al.* 2011; Víquez-zamora *et al.* 2014; Capel *et al.* 2015).

47   A core collection of *S. pimpinellifolium* was developed by AVRDC for the

48    purpose of preservation and utilization (Rao *et al.* 2012). This core collection

49    has been used to mine novel alleles of salt tolerance via genome-wide

50    association study (GWAS) (Rao *et al.* 2015).

51        GWAS utilizes linkage disequilibrium (LD), the non-random association

52    between marker alleles and alleles conferring targeted phenotypes in a given

53    population of germplasm, to map quantitative trait loci (QTLs) (Soto-Cerda and

54    Cloutier 2012). The average ranges of LD decay in different collections of

55    cultivated tomatoes varied from 6 to 13 cM (Sim *et al.* 2012a; Pascual *et al.*

56    2015; Bauchet *et al.* 2017). It is expected that the range of LD decay is smaller

57    in *S. pimpinellifolium* populations because *S. pimpinellifolium* presents larger

58    genetic variation than cultivated tomatoes (Blanca *et al.* 2012, 2015; Ranc *et al.*

59    2012; Bauchet *et al.* 2017). Indeed, the SolCAP array with 7,720 SNPs did not

60    achieve full LD coverage across all chromosomal regions for GWAS using

61    *S. pimpinellifolium* accessions (Sim *et al.* 2012a, 2012b; Bauchet *et al.* 2017).

62    The restriction site-associated DNA sequencing (RADseq) technique might

63    provide an inexpensive solution to address this challenge (Davey and Blaxter

64    2010). The RADseq technique limits sequencing resources at the vicinity of

65    restriction enzyme cutting sites and therefore provides flexibility of

66    experimental design in terms of the trade-off between cost-effectiveness and

67    marker densities (Chen *et al.* 2014; Bhakta *et al.* 2015).

68        The objective of the current study was to develop genome-wide

69    high-density SNP markers for a subset of *S. pimpinellifolium* collections from

70    the Tomato Genetic Resource Center (TGRC) through the RADseq approach.

71    The population differentiation and the range of LD decay were then assessed.

72    **MATERIALS AND METHODS**

73    **Plant materials**

74    All plant materials and their information were obtained from TGRC (Table

75    S1; http://tgrc.ucdavis.edu/). In this study, 12 accessions from Ecuador and 87

76    accessions from Peru were utilized. According to their mating types, 43

77    accessions were facultative self-compatible (FSC) and 56 accessions were

78    autogamous self-compatible (ASC). Seeds were propagated by self-pollination

79    for two generations using the method of single-seed descent in a greenhouse.

80    Young leaves collected from plants of these single-seed descendent seeds

81    were used for DNA extraction.

82    **RAD sequencing**

83    Total genomic DNA was extracted from young leaves using a modified

84    CTAB method (Fulton *et al.* 1995) and purified with a DNeasy Blood & Tissue

85    Kit (QIAGEN, Venlo, Netherland) following the manufacturer's instructions.

86    *Pst*I-digested DNA libraries were prepared following the protocol of Etter *et al.*

87    (Etter *et al.* 2011). Four RADseq libraries were constructed, and each was

88    sequenced in one lane of an Illumina HiSeq2000 flow cell (100 bp single-end

89    reads) (Illumina Inc., San Diego, CA, USA).

90    **SNP calling**

91    Reads were analyzed with Stacks version 1.37 (Catchen *et al.* 2013) and

92    with CLC Genomics Workbench software version 6.5.1 (QIAGEN, Venlo,

93    Netherlands) ("CLC Genomics Workbench 6.5.1"). First, the *process_radtags*

94    command in Stacks filtered out low-quality reads with Q scores less than 20.

95    The remaining reads were mapped to the tomato reference genome SL2.50

96    (Fernandez-Pozo *et al.* 2015) using the "Map Reads to Reference" tool in the

97    CLC Genomics Workbench software. Considering that genetic variation

98    between the tomato reference genome *S. lycopersicum* and

99    *S. pimpinellifolium* is larger than genetic variation within *S. lycopersicum*,

100    mapping parameters were set as 0.5 for the length fraction and 0.9 for the

101    similarity fraction. The reads of the same individual in different lanes were

102    merged together. In the subsequent analyses using Stacks, the *ref_map.pl*

103    command set the parameter –m (minimum read depth to create a stack) as 10,

104    and the *populations* command set the parameter –p (minimum number of

105    populations a locus must be present) as 75. SNPs with a minor allele

106    frequency of less than 0.05 were further excluded, and a set of 24,330 SNP

107    markers was obtained. ITAG2.4 gene model from SGN was used as the

108    reference gene annotation.

109    **Identification of insertion/deletion (InDel) and simple sequence repeat**

110    **(SSR) markers**

111    InDels were identified from the Sequence Alignment Map (SAM) files of

112    all *S. pimpinellifolium* accessions with a read depth of no less than two using

113    the "InDels and Structural Variants" tool and the "Compare variants" tool

114    provided in the CLC Genomics Workbench software ("CLC Genomics

115    Workbench 6.5.1"). InDel markers in the form of tandem repeated sequences

116    were classified as SSR markers.

**Population differentiation**

117

118     To avoid redundant SNP markers used in the subsequent analyses, only

119     one SNP that showed complete LD ($r^2 = 1$) in the same sequencing block

120     around a *Pst*I site was kept whenever more than one SNP existed in the

121     sequencing block. This process resulted in a total of 19,993 SNP markers

122     extracted from the set of 24,330 SNPs noted above. Principle component

123     analysis (PCA) was performed using TASSEL5.0 (Bradbury *et al.* 2007).

124     Population differentiation was investigated via ADMIXTURE (Alexander *et al.*

125     2009). Calculation of pairwise $F_{st}$ (Weir and Cockerham 1984) and analysis of

126     molecular variance (AMOVA) (Excoffier *et al.* 1992) were conducted in the R

127     packages hierfstat (Goudet and Jombart 2015) and StAMPP (Pembleton *et al.*

128     2013), respectively.

**Estimate of genetic variation and LD**

129

130     The genotypes of the 24,330 SNP markers were used to estimate genetic

131     variation and LD in this *S. pimpinellifolium* population. Genetic variation within

132     overall accessions and each of seven groups was assessed based on

133     observed heterozygosity and the within-population gene diversity (expected

134     heterozygosity) using the R package hierfstat (Goudet and Jombart 2015).

135     Pairwise $r^2$ values between SNP markers were calculated to assess overall

136     extent of LD via plink1.9 within a 1-Mb window (Gaunt *et al.* 2007) and fit by

137     non-linear regression (Remington *et al.* 2001). The baseline of the $r^2$ value was

138     set at 0.1 (Bauchet *et al.* 2017). To assess the local LD along each

139     chromosome, we defined the basic unit for local LD as the sequencing region

140     surrounding a *Pst*I site, usually 186 bp long, which has at least one SNP with a

8

141  minor allele frequency greater than 0.05 in the *S. pimpinellifolium* population.

142  For each pair of consecutive basic units, the average $r^2$ was calculated

143  between two SNPs in different basic units and plotted along the left *Pst*I cutting

144  site based on the physical position. The heterochromatin regions were marked

145  according to the genetic map of EXPIM 2012 and the physical map of the

146  tomato reference genome (Sim *et al.* 2012b).

147  **Analysis of SolCAP array data of *S. pimpinellifolium***

148  The SolCAP data of 214 samples of *S. pimpinellifolium* were downloaded

149  from previous studies (Blanca *et al.* 2012, 2015; Sim *et al.* 2012a). A set of

150  4,326 bi-allelic SNPs was first extracted and filtered with the criteria that minor

151  allele frequency is greater than 0.05 and the proportion of missing genotypes

152  is less than 25%. These SNP-filtering criteria are the same as the criteria

153  applied to the SNP dataset generated in this study. This procedure resulted in

154  2,817 SNPs. Subsequently, population differentiation was investigated by

155  ADMIXTURE (Alexander *et al.* 2009). Because some accessions appeared in

156  different SolCAP genotyping studies and their genotypes were not completely

157  matched, different suffixes—"_2012S," "_2012B," and "_2015B,"—were added

158  to the accession name to indicate their origins from references Sim *et al.*

159  2012a, Blanca *et al.* 2012, and Blanca *et al.* 2015, respectively. In addition, the

160  percentage of identical SNP genotypes between accessions with the same

161  name was calculated based on the 4,326 SNP genotypes and excluding

162  missing values.

9

163    **Data avalibility**

164     All the sequences of RADseq are available at the NCBI SRA database,

165    and the BioProject Number is PRJNA358110. Supplemental files available at

166    FigShare: https://doi.org/10.6084/m9.figshare.7010495.v1 for supplemental

167    figures; https://doi.org/10.6084/m9.figshare.7010492.v1 for supplemental

168    tables.


169                                    **RESULTS**

170    **Identification of 24,330 SNPs from *Pst*I-digested DNA libraries**

171     A total of 655,973,270 short DNA reads were obtained from four lanes of

172    the Illumina HiSeq2000 flow cell and were divided into 99 parts according to

173    barcode sequences. Each part was derived from the DNA of a

174    *S. pimpinellifolium* accession and contained at least 3.7 million DNA reads,

175    except for LA2647 (Table S1). To ensure the accuracy of SNP calling and

176    genotype calling, two criteria were set: one was that the read depth aligning to

177    the reference sequences was equal to or greater than 10, and the other was

178    that at least 75% of the accessions showed genotypes associated with a

179    defined SNP marker.

180     Among the 82,814 *Pst*I sites in the tomato reference sequence SL2.50,

181    only 23,988 *Pst*I sites were around the sequenced DNA reads (Table S2). The

182    sequenced regions included 0.54% of the SL2.50 reference sequences and

183    12,790 annotated genes (Table 1). Interestingly, approximately 84% of the

184    sequenced *Pst*I sites were located in the euchromatic regions (Table S2).

185    Nevertheless, no significant difference was observed in the proportion of

10

186    sequencing regions for SNP discovery between the euchromatic regions

187    (68.85%) and the heterochromatic regions (60.59%) (Table S2). A total of

188    67,804 SNPs were identified in the sequenced regions of 99

189    *S. pimpinellifolium* accessions, and 24,330 of them had a minor allele

190    frequency greater than 0.05.

191         In the genotypic dataset of 24,330 SNP markers (Table S3), the missing

192    proportion of each accession ranged from 0.72% to 15.92%, except for

193    LA2647, for which the value was 65.68% due to a low number of sequencing

194    reads (Table S1). Regarding the features of these 24,330 SNPs, 16,365 SNPs

195    were located in 7,383 annotated genes (Table 1) and the remaining SNPs

196    were located in the intergenic regions. In addition, 3,068 InDels (Table S4) and

197    107 SSR markers (Table S5) were obtained. In the subsequent analyses, only

198    SNP markers were utilized, and the genotypic data of the LA0411 accession

199    was dropped because the observed heterozygosity of LA0411 was

200    inconsistent with its mating type (Table S1).

201

202    **Genetic differentiation of *S. pimpinellifolium* corresponded to the**

203    **geographic area**

204         The collection of 98 *S. pimpinellifolium* accessions in this study was

205    divided into seven groups corresponding to three ancestral populations using

206    the ADMIXTURE software (Figure 1A and Figure S1). The seven groups

207    included three groups with pure ancestry, three groups with an admixture of

208    two different ancestries, and one group with an admixture of three ancestries.

209    As expected, accessions in each group were clustered together in the PCA

11

210    plot, in which principal component 1 (PC1), PC2, PC3, PC4 and PC5

211    explained 16.04%, 8.00%, 3.94%, 3.12% and 2.54% of the variation,

212    respectively (Figure 1B). Interestingly, most accessions in the same group

213    were in the same vicinity in terms of their collection sites (Figure 1C). In

214    addition, different ancestral groups were spread in somewhat distinct

215    geographic areas along the coastline from Ecuador to southern Peru (Figure

216    1C). The geographic distribution of these groups appeared in the following

217    order from north to south: the pure red ancestral group, the admixture group

218    with red-blue ancestries, the pure blue ancestral group, the admixture group

219    with blue-green ancestries, and the pure green ancestral group (Figure 1C).

220    This geographic distribution showed a trend in which the admixture groups

221    were located between their corresponding pure ancestral groups.

222        To compare genetic variation within pure ancestral groups or within

223    admixture groups, the within-population gene diversity of each group was

224    calculated. The blue group and the red-blue group showed the highest genetic

225    variation among the pure ancestral groups and the admixture groups,

226    respectively (Table 2). Both groups were in northern Peru, which indicated that

227    northern Peru is the origin of *S. pimpinellifolium*. Pairwise $F_{st}$ confirmed the

228    population differentiation (Table S6), and AMOVA revealed that the variation

229    between groups was 41.96% (p-value < 0.001).

230        The differentiation of two mating types, FSC and ASC, was expected

231    because non-random mating would disrupt the Hardy-Weinberg equilibrium,

232    leading to population structure (Weir and Cockerham 1984; Holsinger and

233    Weir 2009). In this collection, most FSC accessions were clustered in northern

234    Peru, while ASC accessions were scattered in Ecuador and central and

12

235  southern Peru, along the western side of the Andes Mountains to the coast

236  (Figure S2). The pairwise $F_{st}$ (0.0029) of FSC and ASC was significant

237  (p-value < 0.001). However, PCA presented unclear clusters between FSC

238  and ASC (Figure S3). In addition, the variation between FSC and ASC was

239  only 5.91% despite the significance of AMOVA (p-value < 0.001).

240  **Rapid LD decay**

241  Overall LD decay was estimated for the mapping resolution in GWAS. In

242  this population, the non-linear regression curve dropped very quickly (Figure

243  S4). Following the non-linear regression curve, the overall LD decay was

244  within 18 Kb when the baseline of the $r^2$ value was set at 0.1 (Table 3 and

245  Figure 2A). The fastest LD decay was within 10 Kb on chromosome 9 while the

246  slowest decay was within 30 Kb on chromosome 4 (Table 3 and Figure S5).

247

248  **Heterogeneity of genetic recombination within each chromosome**

249  LD decay has often been estimated for each chromosome (Sim *et al.*

250  2012a; Bauchet *et al.* 2017). However, the LD decay per chromosome was

251  insufficient to capture the local variations of historically accumulated

252  recombination events because the tomato genome comprises more than 75%

253  heterochromatin, which usually suppresses recombination events (Sim *et al.*

254  2012a). The local LD profile of individual chromosomes was assessed based

255  on the average $r^2$ value of flanking sequencing units that contained at least

256  one SNP marker. Two major trends were observed (Figure 2B and Figure S6).

257  Marker density in the heterochromatic regions was lower than that in the

258  euchromatic regions, and approximately two-thirds of the $r^2$ values were less

13

259    than 0.1 (Table 3). The latter observation indicated that these flanking SNP

260    markers were not in a state of linkage disequilibrium.

261

262    **DISCUSSION**

263    **A similar distribution between genes and SNPs was identified in the**

264    **vicinity of *Pst*I cutting site throughout the genome**

265        The observation that 67.26% (16,365 to 24,330) of the SNPs were

266    located in the annotated gene regions (Table 1) implied a correlation between

267    the distribution of the identified SNPs in the current study and the distribution

268    of the annotated genes. Additional observations in the current study indicated

269    a preference for genomic DNA digestion by the *Pst*I restriction enzyme in the

270    euchromatic regions: only 28.97% (23,988 to 82,814) of *Pst*I sites were found

271    in the deep sequencing regions, and 83.55% (20,043 to 23,988) of the deep

272    sequencing regions were located in the euchromatic region (Table S2). It is

273    worth noting that the current RADseq protocol did produce low coverage of

274    sequencing reads in certain *Pst*I sites (with a read depth less than 10), and

275    these *Pst*I sites were filtered by the criteria of SNP and genotype calling;

276    therefore, the deep sequencing regions indicated that their read depths were

277    no less than 10. Incidentally, because SNPs can be identified only in the

278    sequenced regions, it is a reasonable deduction that most SNPs found in the

279    current study are located in the euchromatic regions. Plotting the annotated

280    genes, the expected *Pst*I sites, the *Pst*I sites in the deep sequencing regions,

281    and the 24,330 SNPs identified in the current study (Figure 3A, 3B, 3C, and 3D,

14

282    respectively), shows clearly that the annotated tomato genes, the *Pst*I sites in

283    the deep sequencing regions, and identified SNPs are mainly located in the

284    euchromatic regions.

285        *Pst*I is a methylation-sensitive restriction enzyme and recognizes the

286    sequences "CTGCAG" (Dobritsa and Dobritsa 1980). The study of the

287    genome-wide methylation pattern in tomato leaves and immature fruits

288    revealed that the gene-rich euchromatic regions at the distal ends of

289    chromosomes were characterized as the regions with low levels of cytosine

290    methylation at the "CG", "CTG", and "CAG" sequences, and the

291    pericentromeric heterochromatin regions were the regions with high levels of

292    cytosine methylation (Zhong *et al.* 2013). Because the young tomato leaves

293    were used as the DNA source to construct the RADseq libraries, it is

294    reasonable to infer that the *Pst*I-digested RADseq-targeted chromosomal

295    regions were concentrated in the gene-rich euchromatic regions. Therefore,

296    one can emphasize the sequencing resources on euchromatic regions via *Pst*I

297    RADseq when preparing candidate gene research for tomatoes.

298

299    **The discrepancy in inferences of population differentiation of *S.***

300    ***pimpinellifolium***

301        The estimation of the best subpopulation number (K) is very important in

302    GWAS because population structure is integrated as a correction to eliminate

303    the inflated significance due to confounding effects (Korte and Farlow 2013). If

304    the best K of a certain population could not be confirmed, the results of GWAS

305    would be unreliable. However, several previous studies did not achieve the

15

306    same best K of *S. pimpinellifolium*: 10 SSR markers for 248 individuals

307    obtained an unclear K (Zuriaga *et al.* 2009); 48 SSR markers for 190

308    individuals revealed a best K = 2 with admixtures following a K = 5 (Rao *et al.*

309    2012); finally, the SolCAP array for two collections of 63 and 112 individuals

310    obtained the same best K = 3 with admixtures (Blanca *et al.* 2012, 2015). Our

311    study obtained the best K = 3, but our ancestral and admixture groups were

312    different from the latter studies of the SolCAP array (Blanca *et al.* 2012, 2015).

313    These previous studies suggested that the *S. pimpinellifolium* population was

314    differentiated into three ancestral groups: one in the northern Ecuador; another

315    in the mountainous area from southern Ecuador extending to northern Peru,

316    and the third in the low-altitude areas of Peru, along with certain admixtures

317    (Blanca *et al.* 2012, 2015). In contrast, our study showed that the

318    *S. pimpinellifolium* accessions were clustered into three pure ancestral groups,

319    with one in Ecuador (the red group), another in northern Peru (the blue group),

320    and the third in southern Peru (the green group), as well as three clearly

321    identified admixture groups (Figure 1C).

322        To investigate the potential reasons for the inconsistent conclusions

323    between the current study and the previous studies based on SNP markers,

324    genotypic data of the *S. pimpinellifolium* accessions made from the SolCAP

325    array in three previous studies were obtained from internet (Blanca *et al.* 2012,

326    2015; Sim *et al.* 2012a) and a meta-analysis was conducted using our

327    workflow (please see details in the "Materials and Methods" section) (Table

328    S7). A total of 214 samples representing 126 accessions were divided into 11

329    groups via ADMIXTURE using filtered genotypes of 2,817 SNP markers.

330    However, the results of this meta-analysis pose two problems. First, the

16

331   cross-validation error did not confirm that K = 11 was the optimal grouping

332   method (Figure S7 and Figure S8). This condition can be explained as a low

333   population structure in a population with high genetic diversity, which may

334   have resulted from frequent gene flow (Gevaert *et al.* 2013).

335   Overrepresentation of common SNPs within *S. pimpinellifolium* accessions on

336   the SolCAP genotyping array could be the other reason given that the SolCAP

337   array was originally created to explore the genetic variation within cultivated

338   tomatoes and to map genes (Hamilton *et al.* 2012; Sim *et al.* 2012a, 2012b).

339   The second problem was that certain samples belonging to the same

340   accession were not clustered in the same group. For example, two samples of

341   the BGV007104 accession, which shared 93.23% genotypic identity in this

342   SNP set (Table S8) and were labeled as BGV007104_2012B and

343   BGV007104_2015B, were assigned to different groups (Figure S8). The wrong

344   grouping for the samples of the same accession may result from an insufficient

345   number of SNP markers that were unable to capture similarity within the same

346   accession when the sample size increases. In an empirical study of

347   *Arabidopsis halleri*, a few thousand SNP markers were required to estimate

348   the genetic diversity among populations with different genetic variation

349   (Fischer *et al.* 2017). The latter problem prevented meaningful comparisons

350   between the inference of population differentiation in the current study and that

351   of the meta-analysis. Regardless of the inconsistent best K among these

352   studies, the results of pairwise $F_{st}$ and AMOVA statistically supported the

353   subpopulations in this study, suggesting that this set of high-density SNPs

354   could stably estimate the best K.

355

17

**A group of individuals would be a better representative of an accession**

356

357    An accession should be represented by a group of samples rather than

358 only a single individual since an accession in its natural habitat is composed of

359 a group of individuals, especially when gathering accessions with high

360 diversity. However, under circumstances with limited resources, we instead

361 prepared a collection representing a population rather than only several

362 accessions because our final goal was to apply *S. pimpinellifolium* in GWAS.

363 In addition, the mating system and the propagation method of *S.*

364 *pimpinellifolium* made the variations between accessions greater than that

365 within accessions. Therefore, our only option was to involve as much diversity

366 as possible to enhance the efficiency of GWAS.

367

368 **High genetic variation leads to rapid LD decay**

369    The observed and expected heterozygosity of this population were

370 0.0761 and 0.2786, respectively, slightly higher than in previous research

371 (Blanca *et al.* 2012, 2015). Since *S. pimpinellifolium* was detected with up to a

372 40% outcrossing rate (Rick *et al.* 1977) and demonstrated high genetic

373 variation, it is expected to cause rapid LD decay. In this study, LD decay was

374 within 18 Kb throughout the genome, which was much shorter than in

375 cultivated tomatoes (Sim *et al.* 2012a; Bauchet *et al.* 2017). However, such

376 high genetic variation requires much more markers to enable the

377 comprehensive detection in GWAS. The 900-Mb tomato genome requires at

378 least 50,000 markers to cover the entire genome evenly. Therefore, acquiring

379 more SNPs using different methods is essential to conduct a GWAS in the

18

380    *S. pimpinellifolium* population. One possible approach is to increase the

381    sample size evenly for each subpopulation (Brachi *et al.* 2011). Since

382    approximately 64% of alleles were rare in this population, the augmentation of

383    the subpopulation size may adjust rare alleles to common alleles, potentially

384    increasing the SNPs without extending coverage. Another possible strategy is

385    exome sequencing, a selective genome sequencing technology that selects

386    desired sequencing regions by the hybridization of designed probes (Kaur and

387    Gaikwad 2017). Based on tomato genome sequence information, such as the

388    gene model or EST database, one could design different sets of probes to limit

389    sequencing regions (Ruggieri *et al.* 2017). Given the approximately 110 Mb

390    total gene length in the ITAG2.4 gene model, the potential coverage could

391    reach 12% and all target the gene region. This exome sequencing strategy

392    may be able to increase SNPs without increasing the population size.

393

394    **A reproductive strategy would reduce the genetic diversity of *S.***

395    ***pimpinellifolium***

396    These accessions were propagated using single-seed descent for two

397    generations. Therefore, the heterozygosity would be reduced compared to the

398    original specimens, especially for FSC accessions. Here, we revealed that an

399    ASC accession, LA0411, presented 40.25% heterozygosity, which highlighted

400    the contradiction of self-fertilization consequence. Lacking the same accession

401    as a reference in previous studies and considering the 0 to 22% heterozygosity

402    of other accessions in the original published research (Rick *et al.* 1977), we

403    could remove only LA0411 from our analyses based on the fact that its

19

404     heterozygosity was too high for an ASC accession.

## ACKNOWLEDGMENTS

## LITERATURE CITED

416     Alexander, D. H., J. Novembre, and K. Lange, 2009 Fast Model-Based Estimation of Ancestry in Unrelated Individuals. Genome Res. 1655–1664.

418     Bauchet, G., S. Grenier, N. Samson, J. Bonnet, L. Grivet *et al.*, 2017 Use of modern tomato breeding germplasm for deciphering the genetic control of agronomical traits by Genome Wide Association study. Theor. Appl. Genet. 130: 875–889.

422     Bhakta, M. S., V. A. Jones, and C. E. Vallejos, 2015 Punctuated distribution of recombination hotspots and demarcation of pericentromeric regions in *Phaseolus vulgaris* L. PLoS One 10(1): e0116822.

425     Blanca, J., J. Cañizares, L. Cordero, L. Pascual, M. J. Diez *et al.*, 2012

426    Variation Revealed by SNP Genotyping and Morphology Provides Insight

427    into the Origin of the Tomato. PLoS One 7(10): e48198.

428  Blanca, J., J. Montero-Pau, C. Sauvage, G. Bauchet, E. Illa *et al.*, 2015

429    Genomic variation in tomato, from wild ancestors to contemporary

430    breeding accessions. BMC Genomics 16: 257.

431  Brachi, B., G. P. Morris, and J. O. Borevitz, 2011 Genome-wide association

432    studies in plants: the missing heritability is in the field. Genome Biol. 12:

433    232.

434  Bradbury, P. J., Z. Zhang, D. E. Kroon, T. M. Casstevens, Y. Ramdoss *et al.*,

435    2007 TASSEL: Software for association mapping of complex traits in

436    diverse samples. Bioinformatics 23: 2633–2635.

437  Caicedo, A. L., and B. A. Schaal, 2004 Population structure and

438    phylogeography of *Solanum pimpinellifolium* inferred from a nuclear gene.

439    Mol. Ecol. 13: 1871–1882.

440  Capel, C., A. Fernández, J. Manuel, V. Lima, S. Francesc *et al.*, 2015

441    Wide-genome QTL mapping of fruit quality traits in a tomato RIL population

442    derived from the wild-relative species *Solanum pimpinellifolium* L. Theor.

443    Appl. Genet. 128: 2019–2035.

444  Catchen, J., P. A. Hohenlohe, S. Bassham, A. Amores, and W. A. Cresko, 2013

445    Stacks: An analysis tool set for population genomics. Mol. Ecol. 22: 3124–

446    3140.

447  Chen, A. L., C. Y. Liu, C. H. Chen, J. F. Wang, Y. C. Liao *et al.*, 2014

448    Reassessment of QTLs for late blight resistance in the tomato accession

449    L3708 using a restriction site associated DNA (RAD) linkage map and

450    highly aggressive isolates of *Phytophthora infestans.* PLoS One 9(5):

451    e96417.

452    CLC Genomics Workbench 6.5.1.

453    Davey, J. L., and M. W. Blaxter, 2010 RADseq: Next-generation population

454        genetics. Brief. Funct. Genomics 9: 416–423.

455    Dobritsa, A. P., and S. V. Dobritsa, 1980 DNA protection with the DNA

456        methylase M · *Bbv*I from *Bacillus brevis* var. GB against cleavage by the

457        restriction endonucleases *Pst*I and *Pvu*II. Gene 10: 105–112.

458    Etter, P. D., S. Bassham, P. A. Hohenlohe, E. A. Johnson, and W. A. Cresko,

459        2011 SNP Discovery and Genotyping for Evolutionary Genetics Using

460        RAD Sequencing. Mol. Methods Evol. Genet. 772: 1–19.

461    Excoffier, L., P. E. Smouse, and J. M. Quattro, 1992 Analysis of molecular

462        variance inferred from metric distances among DNA haplotypes:

463        Application to human mitochondrial DNA restriction data. Genetics 131:

464        479–491.

465    Fernandez-Pozo, N., N. Menda, J. D. Edwards, S. Saha, I. Y. Tecle *et al.*, 2015

466        The Sol Genomics Network (SGN)-from genotype to phenotype to

467        breeding. Nucleic Acids Res. 43: D1036–D1041.

468    Fischer, M. C., C. Rellstab, M. Leuzinger, M. Roumet, F. Gugerli *et al.*, 2017

469        Estimating genomic diversity and population differentiation- an empirical

470        comparison of microsatellite and SNP variation in *Arabidopsis halleri*. BMC

471        Genomics 18: 1–15.

472    Fulton, T. M., J. Chunwongse, and S. D. Tanksley, 1995 Microprep protocol for

473        extraction of DNA from tomato and other herbaceous plants. Plant Mol.

474        Biol. Report. 13: 207–209.

475    Gaunt, T. R., S. Rodríguez, and I. N. Day, 2007 Cubic exact solutions for the

476  estimation of pairwise haplotype frequencies: implications for linkage

477  disequilibrium analyses and a web tool "CubeX". BMC Bioinformatics 8:

478  428.

479 Gevaert, S. D., J. R. Mandel, J. M. Burke, and L. A. Donovan, 2013 High

480  genetic diversity and low population structure in Porter's sunflower

481  (*Helianthus porteri*). J. Hered. 104: 407–415.

482 Goudet, J., and T. Jombart, 2015 hierfstat: Estimation and Tests of Hierarchical

483  F-Statistics.

484 Grandillo, S., R. Chetelat, S. Knapp, D. Spooner, I. Peralta *et al.*, 2011 *Wild*

485  *Crop Relatives: Genomic and Breeding Resources Vegetables* (C. Kole,

486  Ed.). Springer Heidelberg Dordrecht London New York.

487 Hamilton, J. P., S. Sim, K. Stoffel, A. Van Deynze, C. R. Buell *et al.*, 2012 Single

488  Nucleotide Polymorphism Discovery in Cultivated Tomato via Sequencing

489  by Synthesis. Plant Genome 5: 17–29.

490 Holsinger, K. E., and B. S. Weir, 2009 Genetics in geographically structured

491  populations: defining, estimating and interpreting $F_{ST}$. Nat. Rev. Genet. 10:

492  639–650.

493 Kaur, P., and K. Gaikwad, 2017 From Genomes to GENE-omes: Exome

494  Sequencing Concept and Applications in Crop Improvement. Front. Plant

495  Sci. 8: 1–7.

496 Korte, A., and A. Farlow, 2013 The advantages and limitations of trait analysis

497  with GWAS: a review. Plant Methods 9: 29.

498 Moyle, L. C., 2008 Ecological and evolutionary genomics in the wild tomatoes

499  (Solanum Sect. Lycopersicon). Evolution. 62: 2995–3013.

500 Pascual, L., N. Desplat, B. E. Huang, A. Desgroux, L. Bruguier *et al.*, 2015

501     Potential of a tomato MAGIC population to decipher the genetic control of

502     quantitative traits and detect causal variants in the resequencing era. Plant

503     Biotechnol. J. 13: 565–577.

504   Pembleton, L. W., N. O. I. Cogan, and J. W. Forster, 2013 StAMPP: An R

505     package for calculation of genetic differentiation and structure of

506     mixed-ploidy level populations. Mol. Ecol. Resour. 13: 946–952.

507   Ranc, N., S. Munos, J. Xu, M. C. Le Paslier, A. Chauveau *et al.*, 2012

508     Genome-wide association mapping in tomato (*Solanum lycopersicum*) is

509     possible using genome admixture of *Solanum lycopersicum* var.

510     *cerasiforme*. G3 2: 853–864.

511   Rao, E. S., P. Kadirvel, R. C. Symonds, S. Geethanjali, and A. W. Ebert, 2012

512     Using SSR markers to map genetic diversity and population structure of

513     *Solanum pimpinellifolium* for development of a core collection. Plant Genet.

514     Resour. 10: 38–48.

515   Rao, E. S., P. Kadirvel, R. C. Symonds, S. Geethanjali, R. N. Thontadarya *et al.*,

516     2015 Variations in *DREB1A* and *VP1.1* genes show association with salt

517     tolerance traits in wild tomato (*Solanum pimpinellifolium*). PLoS One 10:

518     1–19.

519   Remington, D. L., J. M. Thornsberry, Y. Matsuoka, L. M. Wilson, S. R. Whitt *et*

520     *al.*, 2001 Structure of linkage disequilibrium and phenotypic associations in

521     the maize genome. Proc. Natl. Acad. Sci. U. S. A. 98: 11479–11484.

522   Rick, C., J. Fobes, and M. Holle, 1977 Genetic variation in *Lycopersicon*

523     *pimpinellifolium*: Evidence of evolutionary change in mating systems. Plant

524     Syst. Evol. 127: 139–170.

525   Ruggieri, V., I. Anzar, A. Paytuvi, R. Calafiore, R. A. Cigliano *et al.*, 2017

526  Exploiting the great potential of Sequence Capture data by a new tool,

527  SUPER-CAP. DNA Res. 24: 81–91.

528  Sim, S. C., A. van Deynze, K. Stoffel, D. S. Douches, D. Zarka *et al.*, 2012a

529  High-density SNP genotyping of tomato (*Solanum lycopersicum* L.)

530  reveals patterns of genetic variation due to breeding. PLoS One 7(9):

531  e45520.

532  Sim, S. C., G. Durstewitz, J. Plieske, R. Wieseke, M. W. Ganal *et al.*, 2012b

533  Development of a large SNP genotyping array and generation of

534  high-density genetic maps in tomato. PLoS One 7: e40563.

535  Soto-Cerda, B. J., and S. Cloutier, 2012 Association mapping in plant genomes,

536  in *Genetic Diversity in Plants*, edited by C. Mahmut. InTech, Rijeka.

537  Spooner, D. M., I. E. Peralta, and S. Knapp, 2005 Comparison of AFLPs with

538  other markers for phylogenetic inference in wild tomatoes [Solanum L.

539  section Lycopersicon (Mill.) Wettst.]. Taxon 54: 43–61.

540  Tanksley, S. D., and S. R. Mccouch, 1997 Seed Banks and Molecular Maps :

541  Unlocking Genetic Potential from the Wild. Science. 277: 1063–1066.

542  Víquez-zamora, M., M. Caro, R. Finkers, Y. Tikunov, A. Bovy *et al.*, 2014

543  Mapping in the era of sequencing : high density genotyping and its

544  application for mapping TYLCV resistance in *Solanum pimpinellifolium*.

545  BMC Genomics 15: 1–10.

546  Weir, B. S., and C. C. Cockerham, 1984 Estimating F-Statistics for the Analysis

547  of Population Structure. Evolution. 38: 1358–1370.

548  Zhong, S., Z. Fei, Y. R. Chen, Y. Zheng, M. Huang *et al.*, 2013 Single-base

549  resolution methylomes of tomato fruit development reveal epigenome

550  modifications associated with ripening. Nat. Biotechnol. 31: 154–159.

25

551    Zuriaga, E., J. M. Blanca, L. Cordero, A. Sifres, W. G. Blas-Cerdán *et al.*, 2009

552        Genetic and bioclimatic variation in *Solanum pimpinellifolium*. Genet.

553        Resour. Crop Evol. 56: 39–51.

554

555    **Figure legend**

556    **Figure 1**. Ancestry and geographic distribution of 98 *Solanum pimpinellifolium*

557    accessions from the Tomato Genetics Resource Center. A) Model-based

558    ancestry for each accession. B) Principle component analysis of the

559    *S. pimpinellifolium* population. C) Geographical distribution of the 98

560    *S. pimpinellifolium* accessions. Symbol and color codes are as follows: square

561    symbols with red, blue and green colors were used to indicate three pure

562    ancestry groups corresponding to the same colors in the ancestry plot; triangle

563    symbols with goldenrod, purple and aquamarine colors were used to present

564    the three admixture groups with red-green, red-blue and blue-green mixing

565    ancestries, respectively; black circle symbols were used for the group with
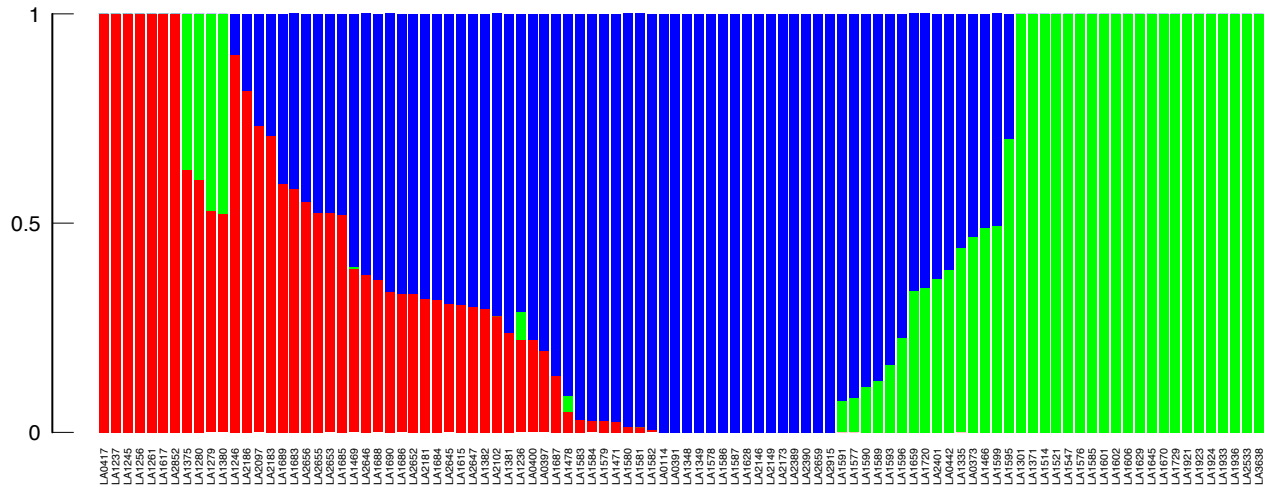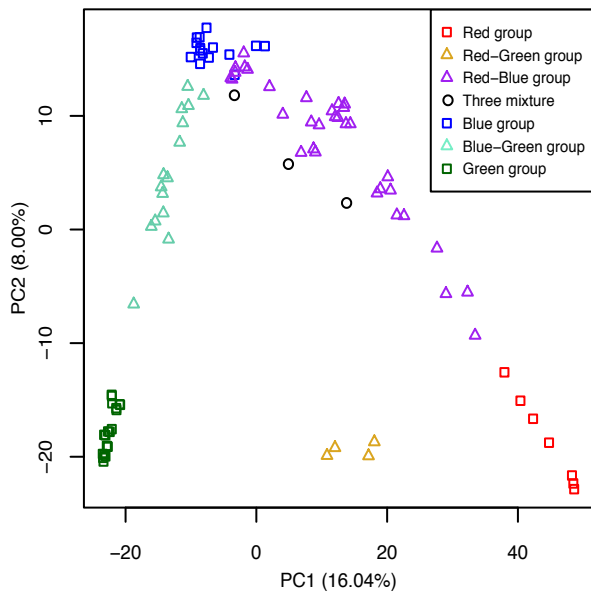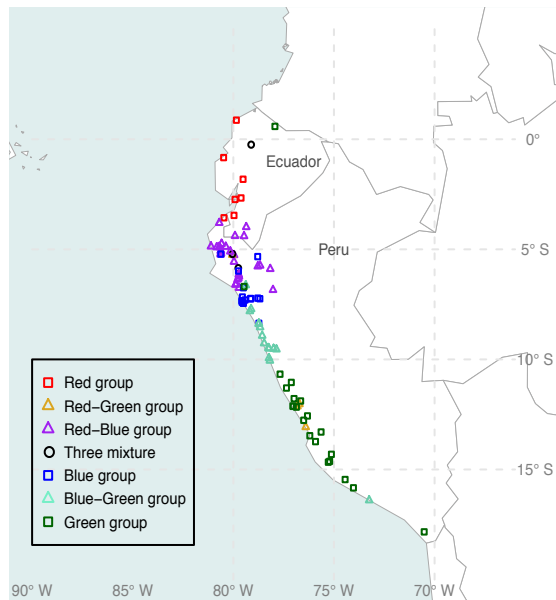
566    admixture of three ancestries.

567

568    **Figure 2**. Visualization for LD. A) The 50 Kb interval of overall LD decay. The

569    red curve indicates non-linear regression, and the black dotted line indicates

570    the baseline of $r^2$ at 0.1. B) The local LD of chromosome 1. The red dotted line

571    indicates the baseline of $r^2$, and the orange line indicated the heterochromatic
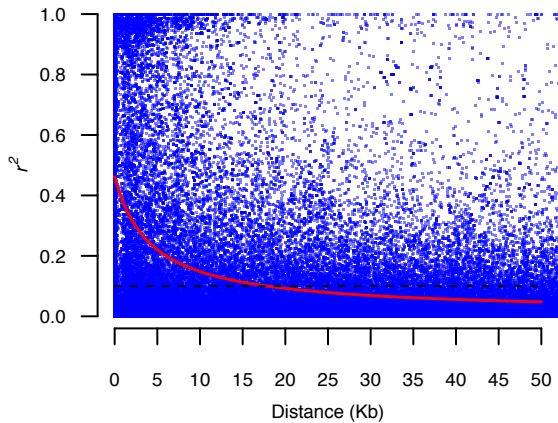
572    region.

573

574    **Figure 3.** The distributions of ITAG2.4 gene model, *Pst*I cutting sites and

575    SNPs throughout the genome. Each section indicates one chromosome, with

26

576     labeling on the circumference. Circles A, B, C, and D indicate the distribution of

577     ITAG2.4 genes, expected *Pst*I cutting sites, *Pst*I cutting sites in the deep

578     sequencing regions and RADseq SNPs, respectively. The black lines in the

579     inner D layer indicate the heterochromatic regions.

580

**A** K=3

**B** PCA

**C**

- □ Red group
- △ Red–Green group
- △ Red–Blue group
- ○ Three mixture
- □ Blue group
- △ Blue–Green group
- □ Green group

**Zoom−in Overall LD Decay**

**Chr. 1**