1

# European Vintage tomatoes galore: a result of farmers combinatorial assorting/swapping of a few diversity rich loci

Running title: Vintage European tomatoes diversification

Jose Blanca[1], Clara Pons[1,2], Javier Montero-Pau[1], David Sanchez-Matarredona[1], Peio Ziarsolo[1], Lilian Fontanet[3], Josef Fisher[4], Mariola Plazas[1], Joan Casals[5], Jose Luis Rambla[2], Alessandro Riccini[6], Samuela Pombarelli[7], Alessandra Ruggiero[7] Maria Sulli[8], Stephania Grillo[7], Angelos Kanellis[9], Giovanni Giuliano[8], Richard Finkers[10], Maria Cammareri[7], Silvana Grandillo[7], Andrea Mazzucato[6], Mathilde Causse[3], Maria José Díez[1], Jaime Prohens[1], Dani Zamir[4], Joaquin Cañizares[1],. Antonio Jose Monforte[2 (*)], Antonio Granell[2(*)]

[1] Instituto de Conservación y Mejora de la Agrodiversidad Valenciana (COMAV-UPV), Universitat Politècnica de València, València, Spain
J. Blanca: jblanca@btc.upv.es
Clara Pons cpons@upvnet.upv.es
J. Montero-Pau: javier.montero@uv.es.
D Sanchez-Matarredona: david.sanchez.matarredona@gmail.com
P. Ziarsolo: pziarsolo@gmail.com
M. Plazas: maplaav@btc.upv.es
M.J. Díez: mdiezni@btc.upv.es
J. Prohens: jprohens@btc.upv.es
J. Cañizares: jcanizares@upv.es


[2] Instituto de Biología Molecular y Celular de Plantas (IBMCP). Consejo Superior de Investigaciones Científicas (CSIC), Universitat Politècnica de València, València, Spain
Clara Pons cpons@upvnet.upv.es
Jose Luis Rambla: jlrambla@ibmcp.upv.es
Antonio Jose Monforte: amonforte@ibmcp.upv.es
Antonio Granell: agranel@ibmcp.upv.es

33

34  [3] INRAE, UR1052, Génétique et Amélioration des Fruits et Légumes, 67 Allée des Chênes,

35  Centre de Recherche PACA, Domaine Saint Maurice, CS60094, Montfavet, 84143, France

36  Mathilde Causse: mathilde.causse@inrae.fr

37  Lilian Fontanet: Lilian.fontanet@hmclause.com

38

39

40

41  [4] Hebrew Univ Jerusalem, Robert H Smith Inst Plant Sci & Genet Agr, Rehovot, Israel

42  Joseph Fisher. Josef.fisher@mail.huji.ac.il

43  Dani Zamir. dani.zamir@mail.huji.ac.il

44

45  [5] Department of Agri-Food Engineering and Biotechnology/Miquel Agustí Foundation, UPC-

46  BarcelonaTech, Campus Baix Llobregat, Esteve Terrades, 8, 08860 Castelldefels, Spain

47  Joan Casals: joan.casals-missio@upc.edu

48

49  [6] Department of Agriculture and Forest Sciences (DAFNE), Università degli Studi della Tuscia,

50  Viterbo, Italy

51  Andrea Mazzucato: mazz@unitus.it

52  Alessandro Riccini: alessandroriccini@gmail.com

53

54  [7] Institute of Biosciences and BioResources (IBBR), National Research Council of Italy (CNR),

55  Via Università 133, 80055 Portici, Italy

56  Silvana Grandillo:  silvana.grandillo@cnr.it

57  Samuela Pombarelli: palombieri@unitus.it

58  Stephania Grillo: grillo@unina.it

59  Maria Cammareri: maria.cammareri@cnr.it

60  Alessandra Ruggiero: alessandra.ruggiero@ibbr.cnr.it

61

62  [8] Italian National Agency for New Technologies, Energy and Sustainable Economic Development

63  (ENEA), Casaccia Research Centre, Rome, Italy

64  Maria Sulli: maria.sulli@enea.it

65  Giovanni Giuliano: giovanni.giuliano@enea.it

66

2

67    [9] Group of Biotechnology of Pharmaceutical Plants, Laboratory of Pharmacognosy, Department

68    of Pharmaceutical Sciences, Aristotle University of Thessaloniki, 54124 Thessaloniki, Greece

69    Angelos K. Kanellis: kanellis@pharm.auth.gr

70

71    [10] Wageningen Univ & Res, Plant Breeding, POB 386, NL-6700 AJ Wageningen, Netherlands

72    Richard Finkers: richard.finkers@wur.nl

73

74    J. M.-G. Current address: Cavanilles Institute of Biodiversity and Evolutionary Biology (ICBiBE),

75    Universitat de València, 46022, Valencia, Spain

76

77    L. M. Current address: HM Clause, Portes-lès-Valence, France,

78

79    (*) Corresponding authors

80    Antonio Jose Monforte, amonforte@ibmcp.upv.es

81    Antonio Granell, agranell@ibmcp.upv.es

82

83

84

85    Date of submission:  July 8th 2021

86    Five figures, eight Supplementary figures, three Supplementary tables.

87    Word count: 6,896

88

# Highlight

The high phenotypic diversity observed among European vintage varieties was created by traditional farmers by combining very few polymorphic loci subjected to balancing selection.


# Abstract

A comprehensive collection of 1,254 tomato accessions corresponding to European heirlooms and landraces, together with modern varieties, early domesticates and wild relatives, were analyzed by genotyping by sequencing. A continuous genetic gradient between the vintage and modern varieties was observed. European vintage tomatoes displayed very low genetic diversity, with only 298 loci out of 64,943 variants being polymorphic at the 95% threshold. European vintage tomatoes could be classified in several genetic groups. Two main clusters consisting of Spanish and Italian accessions showed a higher genetic diversity than the rest varieties, suggesting that these regions might be independent secondary centers of diversity and with a different history. Other varieties seem to be the result of a more recent complex pattern of migrations and hybridizations among the European regions. Several polymorphic loci were associated in a GWAS with fruit morphological traits in the European vintage collection, and the corresponding alleles were found to contribute to the distinctive phenotypic characteristic of the genetic varietal groups. The few highly polymorphic loci associated with morphological traits in an otherwise diversity-poor genome suggests a history of balancing selection, in which tomato farmers maintained the morphological variation by applying a high selective pressure within different varietal types.

Keywords: Crop evolution, diversification, selection, genotyping by sequencing, GWAS, SNP, fruit morphology


Abbreviations


GBS: Genotyping by Sequencing

GWAS: Genome-Wide Association Analysis

LD: Linkage Disequilibrium

LSL: Long Shelf-Life

120    MAF: Minimum Allele Frequency

121    PcoA: Principal Coordinate Analyses

122    QTL: Quantitative Trait Locus

123    SLL: *Solanum lycopersicum* L. var. *lycopersicum*

124    SLC: *S. lycopersicum* var. *cerasiforme*

125    SNP: Single Nucleotide Polymorphism

126    SP: *S. pimpinellifolium*


# Introduction

127

128    The widespread tomato crop (*Solanum lycopersicum* L. var. *lycopersicum*; SLL) originated in

129    Mesoamerica in a region corresponding to today's Mexico as a result of the *S. lycopersicum* L.

130    *var. cerasiforme* (SLC) (Blanca *et al.,* 2012; Blanca *et al.,* 2015; Razifard *et al.,* 2020). Tomato

131    was later brought to Europe, and the Italian botanist Mattioli in 1544 already described varieties

132    with flat, round and segmented fruit types (McCue 1952). This indicated that tomato had probably

133    arrived to Europe in different shapes from America  (Luckwill, 1943; Sanfuentes-Echevarria 2006;

134    Sahagún 1577). Tomato was not immediately adopted for consumption by Europeans, as it was

135    considered at different times and regions as: poisonous, aphrodisiac, ornamental, valuable for

136    sauces and soups, miracle cure and, finally, a fresh salad ingredient (Harvey 2004). It was only

137    as late as the mid-19th century that the tomato became a regular component of the diet in Britain

138    and North America (Harvey 2004). On the contrary, the tomato was better received, extensively

139    cultivated, and consumed as food by the 18th century in Southern Europe, which therefore could

140    have become a secondary center of diversity (Boswell 1937; Bauchet and Causse 2012). As a

141    result of this long tradition of use a large number of traditional varieties are currently available

142    along the Mediterranean basin showing an impressive phenotypic diversity in terms of fruit

143    appearance, adaptation to local conditions and culinary use. Despite the interest for unveiling the

144    population history and the processes that gave rise to the domestication of tomato (Blanca *et al.,*

145    2015; Razifard *et al.,* 2020), there are yet no detailed genetic analyses of the diversification history

146    of the European traditional tomato varieties.

147

148    The extent and type of the molecular variation in the tomato clade has been extensively analyzed

149    in previous studies. The first molecular studies, carried out with isoenzymes, determined that the

150    worldwide cultivated SLL was less variable than the wild *S. pimpinellifolium* (SP) and that the wild,

151    feral and semi-domesticated *S. lycopersicum* var. *cerasiforme* (SLC) was genetically closer to

152   SLL than to SP (Rick *et al.,* 1974; Rick and Fobes 1975). A clear trend of diversity reduction was
153   already observed at the species/subspecies level, probably due to bottlenecks associated with
154   migrations and to the selection pressure imposed by humans during the early domestication
155   stages and development of cultivars from SP to SLC, and lastly, to SLL, (Blanca *et al.,* 2012,
156   2015, Razifard *et al.,* 2020).

158   Despite this limited SLL diversity, several molecular studies have unveiled the worldwide genetic
159   structure within SLL, dividing it into four major groups: processing and fresh market, cherry and
160   vintage tomatoes (Williams and St. Clair 1993; Robbins *et al.,* 2011; Sim *et al.,* 2011; Casals *et*
161   *al.,* 2019). The first three groups correspond to modern tomato varieties created by breeders in
162   the 20th century, characterized by their different culinary use and the introgression of wild species
163   genes, mainly to increase disease resistance and also to develop new type of cultivars. Vintage
164   cultivars are defined as those developed by traditional farmers by intuitive breeding and were
165   cultivated (and some of them are still nowadays locally) before the advent of professional
166   breeding. In this study, landraces, traditional and heirlooms are considered as synonymous of
167   vintage. Park *et al.,* (2004) found genetic differentiation between vintage and modern cultivars. A
168   more comprehensive analysis using 7,720 SolCAP single nucleotide polymorphisms (SNP) from
169   over 426 accessions confirmed the previously described fresh, processing, and vintage groups,
170   at the same time finding two extra clusters located between SLL and SP that corresponded to
171   cultivated and wild cherry tomatoes (Sim *et al.,* 2012). Blanca *et al.,* (2012; 2015) also obtained
172   the fresh, processing, and vintage split and clarified the status of the cherry tomatoes: some of
173   them were SLC from South America, Mesoamerica, and the subtropical regions, while others
174   were modern cherry tomatoes obtained by hybridizing cultivated SLL with wild SP. Blanca *et al.,*
175   (2015), compared with a rarefaction analysis the genetic diversities of the different groups and
176   found that vintage SLL and SLC from outside Peru and Ecuador had the lowest diversity, whereas
177   Peruvian and Ecuadorian SP and SLC had much higher diversities.

179   The studies mentioned above differentiated the modern varieties from the vintage ones, but none
180   of them found any structure within the vintage tomato group. García-Martínez *et al.,* (2006)
181   studied a collection of vintage Spanish cultivars belonging to the varietal groups "Muchamiel",
182   "Pera", and "Moruno" with 19 microsatellite and amplified fragment length polymorphism markers
183   and managed to differentiate the "Pera" type from the other two groups. Mazzucato *et al.,* (2008)
184   dissected a collection of 36 Italian vintage accessions by using 29 microsatellites, and Sacco *et*
185   *al.,* (2015) found differences between 61 Italian vintage varieties and 26 American ones. Current

6

186 genomic sequencing technologies allow finding variable molecular markers even in very narrow
187 genetic contexts. Thus, recently, Esposito *et al.,* (2020), using double digest restriction-site
188 associated DNA sequencing (ddRAD-seq), was able to obtain a sufficient number of SNPs to
189 study the differentiation of a special type of vintage tomatoes cultivated in Spain and Italy, called
190 "de penjar" or "da serbo", characterized by their long shelf-life (LSL). Overall, "de penjar/da serbo"
191 varieties tended to cluster together, showing certain genetic differentiations when compared with
192 other vintage and modern cultivars, but some level of admixture was also found. These former
193 studies were focused on a limited number of accessions from a narrow local diversity and
194 therefore a broader view is clearly needed to better understand the history and relationships of
195 the European vintage varieties.
196
197 In the present study, the genomes of 1,254 European tomato accessions collected from Southern
198 European seed banks were partially sequenced by Genotyping by Sequencing (GBS, Elshire *et*
199 *al.,* 2011; Baird *et al.,* 2008) to obtain genotypes for unbiased markers. Based on these, the
200 genetic structure, diversity, and the association between the polymorphic loci with the
201 morphological variation in that collection were analyzed to shed light on the history of the making
202 up of the diverse vintage European tomatoes.

# Material and methods

## Materials

205 A total of 1,254 tomato accessions were analyzed in this study. One thousand forty four of these
206 accessions are part of the collection of the EU TRADITOM project (www.traditom.eu). Seeds
207 composing the TRADITOM collection were obtained from the genebanks of the Institute for the
208 Conservation and Improvement of Valencian Agrodiversity at the Polytechnic University of
209 Valencia (COMAV-UPV, Valencia, Spain), of the Balearic Island University (UIB, Mallorca, Spain),
210 the Station d`Amelioration des Plantes Maraicheres of the French National Institute for
211 Agricultural Research, (INRA, Montfavet, France), of the Department of Agriculture and Forest
212 Sciences of the University of Tuscia (UNITUS, Viterbo, Italy), of Institute of Biosciences and
213 Bioresources of the Italian National Council of Research (CNR-IBBR, Portici, Italy), of of the
214 Agricultural Research Center of Macedonia and Thrace of the National Agricultural Research
215 Foundation (GGB-NAGREF, Thessaloniki, Greece) and the seed collections of the Miquel Agustí

216  Foundation of the Polytechnic University of Catalunya (FMA-UPC, Casteldefels, Spain),  of

217  BioEconomy of the Italian National Council of Research (CNR-IBE, Catania, Italy), of ARCA 2010

218  S.C.ar.l. (ARCA, Acerra, Italy), of the University of Reggio Calabria (UNIRC, Reggio Calabria,

219  Italy), of the Robert H. Smith Faculty of Agriculture, Food and Environment of the Hebrew

220  University of Jerusalem (HUJI-ARO, Rehovot, Israel). An additional set of 110 accessions were

221  obtained from the COMAV genebank (http://www.upv.es/contenidos/BGCOMAV/indexi.html) that

222  contained 10 wild accessions from the Galapagos Islands, one accession of each wild species *S.*

223  *habrochaites*, *S. chmielewskii* and *S. peruvianum*, 36 *S. pimpinellifolium* accessions from Peru

224  (SP) and North Ecuador (SP_NECu) and 52 *S. lycopersicum* var. *cerasiforme* (SLC) accessions,

225  three modern and 20 SPxSL (*S. pimpinellifolium* x *S. lycopersicum* hybrids, corresponding to

226  cherry cultivars and other crosses between the two species). Passport data can be found in

227  Supplementary Table S1. The germplasm collection was extensively phenotyped in the

228  TRADITOM project (Pons *et al.,* 2017, and in preparation). The dataset corresponding to fruit

229  morphology and color traits obtained at the HUJI-ARO trial was used and analyzed for this article

230  (Supplementary Table S2).

## 231  DNA extraction, library preparation and sequencing

232  Genomic DNA was isolated from young leaves of 5-10 seedlings per accession, using the DNeasy

233  96 Plant Mini Kit (Qiagen, Germany). Genotype-By-Sequencing (GBS) was performed following

234  the procedure reported by Elshire (2011). Briefly, DNA was digested with the restriction enzyme

235  *Ape*K I, barcoded libraries were prepared to track each accession and the DNA sequence

236  corresponding to the region flanking the *Ape*K I site was obtained on an Illumina HiSeq 2000

237  platform by LGC Genomics GmbH (Berlin, Germany). Following the Variant Call Format standard,

238  we used the term sample to refer to one genotyping experiment from one accession.

239

## 240  Read mapping, SNP calling and SNP filtering

241  FastQC was used to evaluate the quality of the sequenced reads, and these were mapped against

242  the *S. lycopersicum* genome build 2.5 (Sato *et al.,* 2012) using BWA mem (Li 2013). After

243  mapping, the PHRED quality of 3 aligned nucleotides from each read end was set to 0 in order to

244  avoid possible false positives caused by misalignments (Li 2011). Mapping statistics were

245  calculated with the samtools stats command (Li *et al.,* 2009).

246

247  SNP calling was carried out by freebayes (Garrison and Marth 2012) with the following
248  parameters: a minimum mapping quality of 57, 5 best alleles, 20 minimum base quality, 0.05
249  maximum mismatch read alignment rate, 10 minimum coverage, 2 minimum alternate allele
250  count, and 0.2 minimum alternate fraction. To avoid regions in the reference genome with
251  potential assembly problems, the Heinz 1706 reads used to build the reference genome were
252  mapped against the reference assembly version SL2.50, a 50X mean coverage was obtained,
253  and when any region had a coverage higher than 200X was removed from the SNP calling.

254

255  SNP and genotype processing were carried out by using the variation Python library located at
256  https://github.com/JoseBlanca/variation. To create the tier1 SNP set to be used in the rest of the
257  analyses, the genotypes with a quality lower than 5 were set to missing, and the variants with a
258  SNP quality lower than 50, an observed heterozygosity higher than 0.1, and a call rate lower than
259  0.6 were filtered out. Moreover, in order to avoid false positives, only variants in which the minor
260  allele was found in more than 2 samples were kept. This filtering was carried out with the
261  "create_tier1.py" script. For some analyses, the pericentromeric regions, that seldom recombine,
262  were removed as part of the heterochromatin. To locate the pericentromeric regions a piecewise
263  regression analysis was applied to the relationship between the genetic distance and the physical
264  positions of the markers of the EXPIM map (Sim *et al.,* 2012). Regression analyses were done
265  using the segmented R library (Muggeo 2003). The calculated pericentromeric regions were:
266  chromosome 1, from 5488553 to 74024603, chromosome 2, from 0 to 30493730, chromosome
267  3, from 16493431 to 50407653, chromosome 4, from 7406888 to 50551374, chromosome 5, from
268  9881466 to 58473554, chromosome 6, from 3861081 to 33077717, chromosome 7, from 4056987
269  to 58629226, chromosome 8, from 4670213 to 54625578, chromosome 9, from 6225214 to
270  63773642, chromosome 10, from 3775719 to 55840828, chromosome 11, from 10947270 to
271  48379978, and chromosome 12, from 5879033 to 61255621.

## PCoA and genetic structure, Diversities and Linkage disequilibrium

273  The genetic structure and the division in subpopulations were determined by conducting a series
274  of hierarchical Principal Coordinate Analyses (PCoA). The PCoAs were carried out with a subset
275  of the variants after filtering. The variant filtering process was comprised of several steps. First,
276  only the euchromatic variants were considered, and from those only the ones with a call rate lower
277  than 0.95, also the ones in which the minor allele was present in less than 3 samples were
278  removed. From the remaining variants, 2000 evenly distributed across the genome were selected.

279  Furthermore, in order to avoid overrepresentation of large regions with complete linkage
280  disequilibrium, when several consecutive variants had a correlation higher than 0.95, only one of
281  them was kept. Finally, pairwise distances between samples were calculated (Kosman and
282  Leonard 2005), and from the distance matrix, a PCoA (Krzanowski and Krzanowski 2000) was
283  generated following the pycogent implementation. These methods were implemented in the
284  do_pca.py script. Additionally, the genetic structure was also estimated with fastSTRUCTURE
285  (Raj *et al.,* 2014).

286

287  The observed and expected heterozygosity and the number of variants per genetic group were
288  calculated considering only the variants variable in the samples involved in the analysis. The script
289  that implemented these analyses is calc_diversities2.py. The allele spectrum figure was plotted
290  by the script calc_maf_trends.py and the rarefaction curves by rarefaction_analysis.py.

291

292  The linkage disequilibrium (LD) was calculated between euchromatic markers with a major allele
293  frequency lower than 0.98 following the Rogers and Huff method for loci with unknown phase
294  (Rogers and Huff 2009).

## GWAS and allele frequencies

296  A heatmap plot that represents the major allele frequency in each group was generated according
297  to a dendrogram by the method implemented in the Python seaborn library
298  (https://seaborn.pydata.org/) and was plotted by the get_most_diverse_snps.py script.
299  A Genome-Wide Association Analysis (GWAS) was carried out with the Genesys R package
300  (Gogarten *et al.,* 2019) on the set of polymorphic variants (95% threshold). The quantitative
301  characters were normalized by using the Box and Cox transformation implemented by the Python
302  scipy library (https://www.scipy.org/). The character normality was checked with a qqplot plotted
303  by the Python statsmodels library (Seabold and Perktold 2010). The correction for genetic
304  structure was calculated with a Principal Component Analysis on the filtered variants implemented
305  by the SNPRelate R library with a 0.3 linkage disequilibrium threshold (Zheng *et al.,* 2012). The
306  quantitative trait associations were tested with the Wald method, and the binomial ones by the
307  Score one. To account for the multiple tests, a Bonferroni threshold was applied. The step-by-
308  step implementation of the GWAS analysis can be analyzed in the gwas.py script.

## Genetic group distances

Two genetic distances among groups were calculated and compared: Nei and Dest (Peakall and Smouse 2006; 2012). They were implemented by the Python variation library and the cacl_pop_dists.py script. From those distances both a neighbor joining tree and a split network were calculated using SplitsTree (Huson and Bryant 2006).

# Results

## High through-put genotyping of a European vintage tomato collection

To genetically characterize vintage European tomatoes, a total of 1,254 tomato accessions were used (Supplementary Table S1). That set included an extensive representation of the extant European vintage tomato variability constituted by 506 accessions from Spain, 305 from Italy, 203 from Greece, 96 from France, and 58 from other origins, with 25 modern commercial cultivars, 39 SP and 22 SLC accessions (the two last ones of American origin) used as references. A total of 3,700 million reads with a mean phred quality of 33.5 were obtained after genotyping-by-sequencing, providing an average of 2.9 million reads per sample. Out of those, 99.0% were successfully mapped to the tomato reference genome (v2.50), but only 55.9% were kept after applying the MAPQ filter with a 57 threshold. These reads were mostly properly paired (96.1%). Of all of the genomic positions that comprise the reference genome, 0.79% had a per sample average sequencing coverage higher than 5X, 0.46% higher than 10X and 0.21% higher than 20X. The complete sequencing and mapping statistics for all samples are available in Supplementary Table S3 and the number of positions per megabase with more than 5 reads in at least 90% of the samples is represented in Supplementary fig 1. Finally, 448,121 variants were called by freebayes, and after filtering them, a working dataset of 64,943 variants was created.

## Genetically defining true European vintage tomatoes and its relationship with American relatives

To genetically position the European tomato collection relative to South and Mesoamerican germplasms, which represent early domestication and improvement steps (Blanca *et al.,* 2015),

11

338  the observed variability of European tomato was analyzed together with SP, SLC, SLxSP hybrids

339  and a sample of modern cultivars including modern fresh and processing cultivars. A series of

340  PCoAs (Fig. 1 and 2) was performed comparing vintage and modern vintage collections. The

341  genetic classification based on these PCoAs can be found in Supplementary Table S1 under the

342  header rank1 classification.

343

344  The PCoA performed with this expanded collection (Fig. 1A and 1B), showed that the green fruited

345  and Galapagos wild species, Peruvian SP (SP), Northern Ecuadorian SP (SP_NEcu), Ecuadorian

346  SLC (SLC_Ecu), Peruvian and Mesoamerican SLC (SLC_Peru_MA), as well as several SP x SL

347  hybrids and admixtures (SPxSL), formed a series of clusters that were clearly separated from the

348  modern and European vintage tomatoes (Fig. 1A and 1B), with the Peruvian and Mesoamerican

349  SLC (SLC_Peru_MA) being the closest American group to the European vintage tomatoes. To

350  obtain a further insight into the genetic architecture of the European tomato, the genetic data was

351  analyzed by using fastSTRUCTURE (Raj *et al.,* 2014). The model marginal likelihoods reached a

352  plateau by four populations (Supplementary Fig. 2). When this result was compared with the

353  PCoA classification, the four fastSTRUCTURE populations were found to correspond to: SP,

354  modern tomatoes, and two distinct vintage populations (Supplementary Fig. 2). It is remarkable

355  that, according to fastSTRUCTURE, the modern tomato, that has been obtained after crossing

356  varieties from different sources, was identified as an original population whereas all the wild and

357  semi-domesticated SLCs, including the Ecuadorian, the Peruvian, and the Mesoamerican ones,

358  appeared as admixtures.

359

360

361  A continuous gradient from vintage to modern rather than clearly split groups was observed in the

362  PCoA plots (Fig. 1A, 1C and 1D). To define the limits between modern and vintage in the PCoA,

363  we chose Heinz 1706 as the reference (in pink, Fig. 1 and 2), since it was one of the first tomato

364  varieties reported to include introgressions from wild *Solanum* species on chromosomes 4, 9, 11

365  and 13 (Sato *et al.,* 2012; Causse *et al.,* 2013; Menda *et al.,* 2014), typical of modern cultivars

366  carrying mainly disease resistance genes.

367

368  PCoA-based classifications indicate that a total of 24.9% of the accessions labelled as vintage

369  according to their passport data mapped outside the vintage genetic cluster in the PCoA space

370  and were localized within the modern and SPxSL genetic groups (Fig. 1). This indicates that either

371  they have been misclassified or correspond to a mixture between vintage and modern varieties.

12

372 To find introgressions in European tomatoes, a haplotype analysis was performed to reveal
373 haplotypes not typically found in the vintage materials. For this, the genome was divided into
374 windows and, in every one of them, the Kosman distances were calculated from the non-vintage
375 samples to the haplotypes found in the vintage samples. When the analyzed non-vintage sample
376 haplotype had a non-zero distance to any of the vintage ones, it was marked as distant from the
377 vintage collection. Several accessions mapping close to the modern varieties in the PCoA space
378 were consistently found to include haplotypes not present in the vintage group (Supplementary
379 Fig. 3) and, despite these being initially catalogued as vintage, it was clear that they actually came
380 from modern breeding programs or were the result of a cross with modern cultivars, and thus
381 were reclassified as modern genetic material (see Supplementary Table S1).

382

383 The modern materials (including both modern references and the vintage reclassified as modern)
384 were spread across the PCoAs according to their use: fresh or processing, and also to their
385 degree of introgression (Fig. 1C, 1D, and 2, Supplementary Fig. 3). PCoAs, when applied only to
386 the modern accessions resulted in four groups (Fig. 2): modern processing, modern and
387 processing long-shelf-life (LSL), modern fresh 1 and modern fresh 2 (Fig. 1 and 2). Modern
388 processing tomatoes, the most distant group to Heinz 1706, were characterized by introgressions
389 that included almost the entire chromosome 5 and the beginning of chromosome 11, and small
390 introgressions in chromosomes 2, 3, 4 and 11 (Supplementary Fig. 3). Modern fresh 1 tomatoes,
391 distributed across the PCoAs between Modern processing and Heinz 1706, were characterized
392 by having a large introgression at the beginning of chromosome 11, a small one at the end of the
393 same chromosome, and another introgression at the beginning of chromosome 6. Modern fresh
394 2 group, which is closer to Heinz 1706, was characterized by having an introgression at the
395 beginning of chromosome 11 (Supplementary Fig. 3).  The modern LSL and processing group
396 was genetically very close to Heinz 1706 (in blue, Fig. 1C and 1D, sharing a large part of
397 chromosome 9, including an introgression considered to be the result of the introduction of the
398 *Tm-2* gene, conferring resistance to Tomato Mosaic Virus, in modern breeding programs. All of
399 these haplotypes could be used for the identification of non-true European vintage tomatoes.

## Diversity among European vintage tomatoes

401 European vintage tomatoes are usually considered to have low genetic diversity (Blanca *et al.,*
402 2015). Therefore, it was important to calculate the number of polymorphic variants present in our
403 collection of European vintage tomatoes, the largest collection analyzed by sequencing thus far,

13

404 and to compare it with the variability present in the wild SP, the wild and semi-domesticated SLC,

405 and the modern cultivars. The number of variants within the European vintage collection was quite

406 large (26,129), it was even larger than the number found in SP (19,164), in SLC (7,690), or in the

407 materials classified as modern (17,328). However, this comparison could be biased in favor of

408 the vintage collection because of the larger number of samples in vintage 890, compared to SP

409 24, SLC 42, and modern 243.

410

411 To correct for this factor, diversity indexes were calculated with the same number of samples (20)

412 (fig 3A) and the analysis was repeated 100 times, with a different set of 20 samples chosen at

413 random each time. Both the Nei diversity and the percentage of polymorphic variants (with a 95%

414 threshold) was much higher in the wild SP that in any other group, and, even more relevant, both

415 indexes were the lowest, by far, in vintage. The analysis indicated that the many of the variants

416 found in the vintage collection could not be considered polymorphic. At 95 % threshold, the

417 vintage collection contained only 298 polymorphic variants. This scarcity in polymorphic variants

418 in the European vintage group can also be observed in the allele frequency spectrum

419 (Supplementary Fig. 4) in which it is clear that most variants were almost fixed in the vintage

420 collection.

421

422 To better compare the amount of genetic variability in each major cultivated group

423 (SLC_Peru_MA, vintage, and modern) a rarefaction analysis was carried out. In this analysis, the

424 samples were added one at a time, to check if the number of variants, including the ones at very

425 low frequencies, reached a maximum when more samples were considered (fig 3B). The number

426 of variants found in the vintage group was always lower than in the modern and SLC_Peru_MA

427 groups. However, the total number of variants within the vintage collection kept increasing as

428 more samples were added. However, the number of polymorphic variants did stabilize with a few

429 samples. Finally, the Nei diversity decreased (Supplementary Fig. 6) when more samples were

430 added. This decrease was due to the high number of variants found within vintage that were close

431 to fixation.

## Linkage disequilibrium

433 The linkage disequilibrium (LD) was calculated for the genetic groups with enough polymorphic

434 markers (Minimum Allele Frequency (MAF)> 0.02 threshold) (Supplementary Fig. 6), which were

435 SP, SLC from Peru and Mesoamerica, modern, and European vintage varieties. Wild SP showed

436   the lowest LD ($r^2$=0.42) at 5 kb and it was also the group in which LD decreased the fastest, being

437   only $r^2$=0.2 at 25 Kb. In SLC, $r^2$ was 0.8 at 5 kb and at 1000 Kb it was still 0.4. Vintage had the

438   highest LD at 25 Kb ($r^2$=0.97); however, it decreased to the lowest value ($r^2$=0.05) at 1000 Kb.

439   The modern accessions had a high LD both at 25 Kb ($r^2$=0.9) and at 1000 Kb ($r^2$=0.6). The LD

440   found at 1000 Kb is likely due to population substructure. SLC and modern had high long range

441   LDs, perhaps because modern included both fresh and processing accessions, which were

442   clearly separated in the PCoAs, and SLC contained accessions from Peru and Mesoamerica, two

443   geographically distant areas. Additionally, modern cultivars often contain introgressions from wild

444   species, including disease resistance genes, that span large regions for which recombination is

445   usually suppressed. SP is also known to have a clear population structure (Blanca *et al.,* 2012)

446   and also showed some long range LD, which clearly supports the conclusion that LD is not due

447   just to gamete disequilibrium, but to other causes too. The vintage accessions showed the lowest

448   LD at 1000 Kb perhaps because it has a less remarkable population substructure.

## 449   Classification of vintage tomato clusters

450   To further classify true vintage tomatoes, a series of PCoAs (Supplementary Fig. 7) were

451   performed. A genetic group was created when several samples that grouped together in the

452   PCoAs shared their geographic origin or traditional variety name, or some aspect of their

453   phenotype (Supplementary Table S2) e.g. fruit shape and size. Most vintage samples could be

454   classified into 27 different genetic groups, using this PCoA strategy, and were named as *"Balearic*

455   *cherry",* "*Bell pepper", "Cor de bou", "Greek (grc)", "Italian (ita) ellipsoid", "Ita grc", "Ita small",*

456   *"Lemonia", "Liguria", "Long Shelf Life (LSL) da serbo", "LSL heart", "LSL penjar cat", "LSL penjar*

457   *vlc", "LSL piennolo", "LSL ramellet", "Marmande", "Montserrat", "Muchamiel", "Palosanto pometa*

458   *1", "Palosanto pometa 2", "Pera girona", "Pimiento", "San Marzano", "Scatolone di bolsena",*

459   *"Spagnoletta", "Tondo piccolo", "Valenciano".*

460

461   Two connected clusters of genetic groups (for the sake of clarity we will use "group" to refer to a

462   PCoA group of samples and "cluster" to talk about a cluster of groups) were observed in PCoA

463   (Supplementary Fig. 7A and 7B). Within the cluster at the center of PCoA, we found the genetic

464   groups *"LSL ramellet"* , *"LSL penjar vlc"* , *"LSL penjar cat"*, *"LSL ramellet", "Marmande",*

465   *"Montserrat"*, "*Bell pepper", "Lemonia", "Muchamiel", "Palosanto pometa 1", "Palosanto pometa*

466   *2", "Pera girona", "Scatolone di bolsena", "Spagnoletta" and "Valenciano"* (Supplementary Fig. 7C

467   and  7H). These genetic groups belong to Spain, with the exception of *"Marmande",* "*Bell pepper"*

468 and *"Palosanto pometa 1"*, which were represented in all four Mediterranean countries (Spain,

469 Italy, France and Greece), the Italian *"Scatolone di bolsena"* and *"Spagnoletta",* and the

470 Greek *"Lemonia"* (Fig. 4A). Outside the central cluster, but close, we found groups of big

471 tomatoes: "*Liguria*", with accessions mainly collected in Italy, and "*Cor de bou*" and "*Pimiento*",

472 present in all four countries (Supplementary Fig. 7C and 7D, Fig. 4A and 4B). A second cluster

473 included mostly Italian accessions classified into the *"Ita ellipsoid", "Ita small", "LSL da serbo",*

474 *"LSL piennolo", "San Marzano",* and *"Tondo piccolo"* genetic groups, and also some Greek and

475 Spanish accessions included in the *"grc", "Ita grc"*, *and "Balearic cherry"* genetic groups, all

476 characterized by having a small size with no or weak ribbing (Supplementary Fig. 7A, 7B. 7M-R,

477 Fig. 4A and 4B). In summary, the PCA separated vintage accessions mainly by country of origin

478 and fruit size. It is interesting to note that the LSL-type accessions, which were highly represented

479 in the collection, were not grouped together, but rather segregated by country: the Italian LSL

480 varieties were found within the Italian cluster, and the Spanish LSL within the Spanish cluster.

481 Several accessions located between the Spanish and the Italian clusters could not be grouped

482 by passport data or any other characteristic.

483

## Allele frequencies across the genome in Vintage groups and their relationship with phenotypic diversity

486 A clustering of the vintage genetic in groups based on a distance tree was calculated using the

487 polymorphic variants (95% threshold) (Fig. 4A). This analysis showed that the defined genetic

488 groups had quite distinct allele frequencies along the genome. Concomitantly, the genetic groups

489 also showed enrichment in specific phenotypic characteristics related to their horticultural

490 classification. For example, varieties belonging to the genetic groups "*Pera girona*", "*LSL ramelle*t"

491 and "*LSL penjar vlc*" have colourless skin, while "*Balearic cherry*", "*Tondo piccolo*", *LSL piennolo*",

492 "*Lemonia*", "*LSL heart*", "*LSL Penjar vlc*", "*grc*", and "*San Marzano*" showed mostly weak ribbing,

493 and, finally, "*Spanoletta*" was characterized by its fasciation (Fig. 4B). Moreover, the fruit size was

494 also different for different genetic groups and clusters *"LSL heart"*, "*Ita ellipsoid", "Ita small", "LSL*

495 *da serbo", "LSL piennolo", "San Marzano", "Tondo piccolo*", *"grc", "Ita grc" and "Balearic cherry"*

496 were characterized by having a small size, while the rest were medium or large in size.

497 Furthermore, several noticeable clusters of genetic groups with common phenotypic traits could

498 be observed. For instance, there was a cluster formed by small-fruited, slightly-ribbed, long shelf-

499 life and processing Italian genetic groups which included the well-known Italian *"da Serbo"* and

16

500    *"San Marzano"* tomatoes. Another cluster was comprised mainly by long shelf-life colourless-
501    skinned Spanish tomatoes, which included the *"LSL penjar cat"*, *"LSL penjar vlc"*, and the *"LSL*
502    *ramellet"* groups. Interestingly, this cluster also included the Catalonian big fruited *"Monserrat"*
503    group which, in contrast to the others, were fasciated and used for fresh consumption. Close to
504    this cluster were some of the most typical Spanish vintage fresh-market varieties: *"Valenciano",*
505    *"Muchamiel"* and *"Palosanto pometa 2".* In addition, big tomatoes appertaining to *"Liguria"*, "*Cor*
506    *de bou*", and "*Pimiento*" clustered together.

507

508    Some of the genetic differences between the groups could be due to genetic drift not related with
509    the phenotypic variability generated during the history that gave rise to the different vintage
510    varieties, but allele frequencies of genes involved in the phenotypic variation could have been
511    selected either inadvertently or consciously by traditional farmers. In order to elucidate whether
512    the differentiating variants were associated to the phenotypic variation observed in the different
513    genetic groups a GWAS analysis was carried out using selected fruit characters (Fig. 4C).
514    Two of the main phenotypic characteristics differentiating the vintage tomatoes are fruit weight
515    (fw) and ribbing (Fig. 4B). In the GWAS analysis fruit weight was associated to variants on
516    chromosome 1, 3 and 11. The MAF analysis indicated that most of the small fruited tomatoes
517    such *"LS heart"*, "*Ita ellipsoid", "Ita small", "LSL da serbo", "LSL piennolo", "San Marzano", "Tondo*
518    *piccolo*", *"grc", "Ita grc"*, and *"Balearic cherry"* shared the fixation of the same allelic variant in
519    chromosome 1. The pattern found in chromosome 3 was similar, except for the *"LS heart"* and
520    "*LSL piennolo"* groups.
521    For ribbing, GWAS revealed association with variants on chromosomes 1, 7, 10 and 11. The
522    chromosome 1 region was fixed in the weak ribbed groups *"Balearic cherry"*, *"Tondo piccolo*" and
523    *"LSL piennolo"*. In contrast, almost all medium and large tomatoes, with the exception of
524    *"Pimiento"* and "*Spanish LSL*" fruits (both showing no or weak ribbing) had a fixed common variant
525    in chromosome 11 that was associated by GWAS with fruit weight, ribbing at calyx end, and fruit
526    shape index.

527

528    Another trait differentiating vintage tomato cultivars was skin colour, for instance, most Spanish
529    LSL as well as tomatoes included in the *"Cor de bou"*, *"Montserrat"*, and "*Pera girona*" genetic
530    groups had a colourless skin, which resulted in pinkish fruit (Fig. 4B). GWAS found association
531    with this pink color in chromosomes 1 (two regions), 3, 5, and 10. The GWAS and MAF analysis
532    comparison (Fig. 4A and 4C) showed that different pink genetic groups had different allelic
533    composition in the associated genomic regions, what might reflect a complex genetic control.

17

534    Fruit shape was associated with regions in chromosomes 2, 5, 10, and 12. The region in
535    chromosome 2 was fixed in *"Marmande"* and "*Scatolone di bolsena*", two groups that are well
536    known for having flat fruits. In addition, "*Valenciano*", "*Pimiento*", and "*Liguria*" had the minor allele
537    almost fixed in the chromosome 10 region. High frequency minor alleles, almost fixed in the
538    regions associated to fruit shape in the GWAS, were also observed in other genetics groups such
539    as the Italian "*LSL da serbo*", in chromosome 5, and "*Ita ellipsoid*" and "*Tondo Piccolo*", in
540    chromosome 6 as well as in "*Cor de bou*" and *"Pimiento"* groups, in chromosome 12.
541    In the case of use, associated variants were found in chromosomes 10 and 11, but, in this case,
542    no clear relationship was found between allelic frequencies among the tomato genetic groups and
543    GWAS.

## Network analyses supports the differentiation between Spanish and Italian vintage tomatoes and the occurrence of hybridization events in vintage tomatoes across Europe

547    To study the genetic relationships between accessions and groups of accessions, a network
548    based on pairwise Dest group distances was created with Splitstrees. Evolutionary relationships
549    are often represented as an unique tree under the assumption that evolution is a branching or
550    tree-like process (Husson 1998). However, real data does not always clearly support a tree.
551    Phylogenetic split decompositions represented in a network may be evidence for conflicting
552    reticulated phylogenies due to gene flow and/or hybridization (Husson 1998).
553    The splitrees network of European tomato is depicted in Fig. 5. The group organization in the
554    network was structured, like the PCoAs (Supplementary figure 7), in two main country-related
555    clusters. One cluster was comprised mainly of Spanish vintage groups, which included the
556    Spanish LSL, *"Muchamiel"*, and *"Montserrat"* types, and another cluster was mostly comprised by
557    the small fruited Italian LSL and processing groups. Interestingly, the "*Liguria*" group clustered
558    with Spanish clusters, although the branch that linked it with the core Spanish clusters was quite
559    large.
560    The degree of reticulation found (Fig. 5) suggested that hybridizations might have occurred
561    between the ancestors of accessions collected from the same geographical regions. On the other
562    hand, the groups that included accessions from different countries, such as "*Marmande*",
563    "*Pimiento*", "*Cor de bou*" or "*Palosanto pometa 1*", were located between the Spanish and Italian
564    clusters.

18

565 These groups of mixed origin could be more modern and derived from hybridization from old
566 Spanish and Italian varieties or, alternatively, they could be very old varieties found across Europe
567 before the Spanish and the Italian diversification started. To check those possibilities, a
568 rarefaction analysis was performed of the number of polymorphic sites found in these three
569 clusters was calculated (Supplementary Fig. 8). The number of polymorphic sites was clearly
570 higher in the Italian and Spanish clusters and much lower in the mixed origin cluster, an evidence
571 that supports that Spain and Italy were secondary centers of diversity for the European tomato,
572 whereas the varieties included in the mixture cluster would be more recent.
573

# Discussion

574

## Very low, but discriminant, variation in vintage European tomatoes

575

576
577 The genetic diversity of this European vintage collection was very low when compared with the
578 diversity found in SP or even in SLC. While this result is in agreement with previous surveys on
579 worldwide SLL accessions carried out with the SolCAP SNP platform (Sim *et al.,* 2012, Blanca *et*
580 *al.,* 2012; Blanca *et al.,* 2015), the current analysis represents the first estimate obtained using a
581 comprehensive representation of vintage European tomatoes, and it is relevant to study the role
582 of Europe as a secondary center for tomato diversification. The low level of diversity found in
583 these traditional materials was quite striking: after sequencing 0.8% of the genome, only 298
584 polymorphic variants at the 95% level were found. This result is quite remarkable when we
585 consider the high phenotypic diversity of vintage tomatoes. Moreover, the high linkage
586 disequilibrium found in these traditional vintage materials suggests that it is rather unlikely that
587 the total number of polymorphic blocks would grow much even if whole genome sequences were
588 to be obtained.
589

590 Previous studies demonstrated a strong bottleneck during the SLC tomato's travel from Ecuador
591 and Peru to Mesoamerica (Blanca *et al.,* 2015, Lin *et al.,* 2014; Razifard *et al.,* 2020). However,
592 it is remarkable that despite the low genetic diversity found in vintage European tomatoes there
593 are still a few highly polymorphic loci within this tomato gene pool. Some of this variation could
594 be due to the random nature of genetic drift. However, the association study carried out with major
595 phenotypical/morphological traits found that a sizeable fraction of those diverse loci were

596  associated with the vintage fruit phenotypical/morphological variation. Therefore, it is quite likely
597  that many of those polymorphic loci had been under balancing selection (Delph and Kelly, 2014)
598  during the diversification process and were in fact responsible for a sizeable part of the tomato
599  phenotypic variation, or, at least, in linkage disequilibrium with the variants selected. It may seem
600  paradoxical that the high diversity of shapes, colors, sizes, uses, and other agronomic traits in the
601  vintage group could be maintained by such a poor gene pool, but it seems that the selection
602  carried out by the traditional growers in favor of this agronomic diversity resulted in a desert of
603  variation, with just a handful of scattered polymorphic loci. This is consistent with two highly
604  polymorphic SNPs found by Muños *et al.,* (2011) in the *lc* locus. These were highly polymorphic,
605  but were surrounded by loci with "drastically reduced" diversity. Thus, they seemed to be the
606  result of selection for low or high number of locules in different materials.

607  Recently, structural variants (SV) were studied in tomato using new long-read sequencing
608  technologies and new analysis algorithms (Alonge *et al.,* 2020; Domínguez *et al.,*, 2020). A large
609  number of structural variants were identified and were mostly generated by transposons and
610  related repeats. Similar to the variants studied here, most structural variants had a very low
611  frequency, and the majority were singletons.

612  Therefore, the phenotypic diversity present in European vintage tomatoes seems to have been
613  built by remixing/reshuffling/swapping very few polymorphisms with the selection pressure
614  associated with the creation of new varietal types and to the adaptation of these types to different
615  regional environments.

616

617  ## Tomato History: tomato movement in Europe

618

619  The distribution of the genetic variability in the European vintage tomatoes showed mostly a
620  continuous gradient. However, the Spanish and Italian varieties occupied opposite regions of the
621  PCoA space what supports a genetically differentiation among varieties originated in those
622  countries. The lack of clear-cut limits may be due to migrations between different regions and
623  countries and subsequent intercrossing. Despite this difficulty, the genetic vintage groups
624  proposed here were differentiated by characteristics such as: their main geographic origin, use,
625  fruit morphology, and varietal name. The genetic groups sometimes corresponded with the
626  varietal type, such as in "Valenciano", "Muchamiel", "Penjar" or "Piennolo". However, the match
627  between the proposed genetic group and the sample varietal name was seldom complete. For
628  instance, the "*Cor de bou*" group included two "Valenciano" samples, one "Russe", and one

629　"Costoluto". This may be due to the limitations of the genotyping or genetic classification
630　methodology utilized or to erroneous passport data, as it may not be trivial for a standard grower
631　to evaluate the sometimes subtle varietal differences. Other genetic groups, such "Italian small"
632　showed no clear associations to any variety name.. Finally, cultivars previously classified as
633　belonging to some variety, such as "Marmande", were included in many different genetic groups.
634　It is likely that the popularity of some varietal types such as "Marmande", made some growers
635　prone to apply the label to any variety that displayed the typical morphological characteristics of
636　a well-known varietal type. Thus, the "Marmande" tomatoes are characterized by its production
637　of large and multi-locule tomatoes, and any other variety with a similar fruit morphology could
638　have been labeled as "Marmande".

639

640　One clear example of mistaken identities and/or inadvertent out crossing is provided by the
641　vintage samples that were found to include haplotypes not found in the vintage core and to be
642　genetically closer to the modern varieties than to the vintage materials in the PCoA. It is not even
643　trivial to define the borderline between vintage and modern varieties. One could think that until
644　the 1950´s most varieties were heirlooms and landraces maintained by small farmers, but the real
645　history is more complex. When tomato cultivation was popularized in the 19th Century in France,
646　England, and the USA some of the varieties were already provided by seed companies (Boswell
647　1937), and there were seed shipments documented between countries, for instance, from
648　England to the Canary Islands (Amador *et al.,* 2012). Moreover, from 1910 onwards, professional
649　breeding efforts created new varieties adapted to long-distance shipping and for processing
650　(Boswell 1937). These efforts did not yet include wild materials, so their results are not easy to be
651　differentiated in a PCoA analysis. It is only when shortly afterwards, breeders started introgressing
652　wild species alleles for disease resistance, that the varieties created were different enough to be
653　easily differentiated in the PCoA analysis. In any case, the vintage-modern limit has to be
654　somewhat conventional, although a characteristic of modern cultivars compared with vintage
655　varieties is the introgression of genes from wild species. Therefore, true vintage cultivars were
656　defined based on the absence of wild species haplotypes.

657

658　Most of these introgressions seem to be related to disease resistance genes as the *Cladosporium*
659　*fulvum* resistance gene *Cf-2* in chromosome 6, *Tm-2* (resistance Tomato Mosaic virus,) in
660　chromosome 9. It is likely that the modern genetic variability has been combined with the true
661　traditional varieties, so some materials catalogued in the genebanks as traditional are in fact a
662　mixture of traditional and modern. This is to be expected, as the seed collectors/genebanks label

663   as vintage any material considered as such by the farmer from whom the seeds where collected.
664   Although European small farmers often save their own tomato seed, they may occasionally
665   purchase or get plantlets from markets or nurseries or save seeds from modern varieties
666   purchased in the market and introduce them in their fields. This may lead after several years of
667   reproduction and farmer selection to complex hybridizations and mixings. Clearly, there have
668   been many opportunities for introgressing modern haplotypes into the vintage materials, such as
669   unintentional crosses. This phenomenon could be thought of as blurring the boundaries of a
670   supposedly pure vintage population, but one may also think that this leakage had the positive
671   unintended consequence of increasing the very low diversity of the vintage pool, and it is also the
672   case that evolution consists of change and adaptation of local varieties (Casañas *et al.,*, 2017).

674   The allele frequency based tree (Fig. 4) defined three major clusters: Spanish, Italian, and Mixed
675   origin. The mixed origin groups are basal in the Fig.4 tree, have longer branch lengths, and occupy
676   an intermediate position between the Italian and Spanish clusters in the Dest network (Fig. 5).
677   These results are compatible with the hypothesis that Italy and Spain formed two centers of
678   diversity. The differentiation of Italian and Spanish gene pools is exemplified by the long LSL
679   varieties from both countries. Italian and Spanish LSL varieties were clustered apart from each
680   other with only a small number of samples from the other country, so it seems as if the origin of
681   the long shelf life tomatoes in both countries was independent. The transformation from a fresh
682   to a long shelf-life variety is likely due to a limited number of loci, as observed in Fig. 4 in which
683   the Catalonian fresh "Montserrat" type is closely related to the Catalonian long shelf-life "Penjar"
684   type. Esposito et al (2020) also observed geographic differentiation of the Italian and Spanish
685   long shelf-life varieties. Therefore, although there may have been migrations from Italy to Spain
686   and vice versa, these may not have been extensive enough to erase the genetic differences
687   between the Italian and Spanish varieties

689   Regarding the mixture cluster, the groups included in it are basal in the Fig.4 tree, they have
690   longer branch lengths, and occupy an intermediate position between the Italian and Spanish
691   clusters in the Dest network (Fig. 5). Moreover, the rarefraction analysis supports that it included
692   varieties derived from the two secondary centers of diversity. This could be the result of long of
693   tomato cultivation tradition in Southern Europe, being the groups included in this cluster
694   developed from hybridizations between the two centers of diversity. New mutations, other
695   introductions of tomatoes from America or new genes from varieties developed worldwide might
696   also be involved in the history of the groups of mixed origin.

697

698    A complex pattern of migrations can also be inferred in several genetic groups as the "*Cor de*

699    *bou*" group that included varieties from most countries: French "Coeur de boeuf", Italian "Cuor di

700    bue", Catalonian "Cor de bou", Castillian "Corazon de toro", and "Navarran corazón de fitero".

701    Also, the Italian "Spagnoletta" group was closely related with the "*Marmande*" group comprised

702    by French, Spanish, Greek, and Italian accessions. Other genetic groups with mixed geographic

703    origin are "*Liguria*", "*Cor de bou*", "*Pimiento*", "*Palosanto Pometa 1*"and "*Marmande*".

704

## Do a few Polymorphic genes differentiate the true European vintage tomato genetic groups?

705

706

707

708    In order to shed light on the apparent contradiction between the low genetic diversity and the

709    large phenotypic variation of European vintage tomatoes, a GWAS was carried out with the

710    polymorphic variants and some of the most obvious morphological traits (fruit morphology, color,

711    and ripening behavior).

712    Variants located in the genomic regions of previously identified loci involved in fruit weight, and

713    likely involved in domestication and diversification, were associated with this trait in the GWAS

714    performed with the European vintage collection. Most of the small fruited tomatoes shared fixed

715    variation regions in chromosomes 1 and 3 which mapped close to previously-described

716    Quantitative Trait Loci (QTL)and genes associated with fruit size: *fw1.1* (Grandillo *et al.,* 1999)

717    and *fw3.2/KLUH* and *ENO* (Chakrabarti *et al.,* 2013; Yuste-Lisbona *et al.,* 2020) (Fig. 4A and 4B).

718    In contrast, almost all medium and large tomatoes shared a region in chromosome 11 that

719    mapped close to FAS (Xu *et al.,* 2015) and *fw11.3*/CSR (Mu *et al.,* 2017), with both genes playing

720    a known role in controlling fruit size and fasciation. No more associations were observed in other

721    genomic regions for fruit weight, so it seems reasonable to think that these QTLs can be

722    responsible, at least in part, for the fruit size variability among the European vintage tomatoes.

723    Regarding fruit shape, two of the associated regions found included known genes. The

724    chromosome 2 region included previously mapped QTLs as heart shape *hrt2.2* (heart shape),

725    *pblk2.2* (proximal end blockiness), *psh2.2* (shoulder height), *piar2.2* (indentation area) (Brewer *et*

726    *al.,* 2007) and *ovate* (Liu *et al.,* 2002), and the region in chromosome 10 is located close to, where

727    the original copy of the *sun* locus was found (Xiao *et al.,* 2008). In the case of skin color, a different

728    pattern was characteristic of different pink genetic groups. "*LSL Penjar* vlc" and "*LSL ramellet*"

23

729   shared a variant at the end of chromosome 1 that matched a region that was previously

730   associated with skin color, the colorless-peel and locus (Ballester *et al.,* 2010), while "*Pera*

731   *Girona*" had the minor allele for the other chromosome 1 variant, which is located at the beginning

732   of the chromosome and maps close to the SlCMT3 (Gallusci *et al.,* 2016) and PSY3 (Li *et al.,*

733   2008) genes involved in epigenetic ripening regulation and carotene biosynthesis, respectively.

734

735   The current analysis suggests that fruit morphology variability among European vintage tomatoes

736   could be the consequence of the combination of a relatively low number of genes, as suggested

737   by Rodriguez *et al.,* (2011), including *fw3.2/KLUH, ENO, FAS,* SUN, and OVATE. On the other

738   hand, skin color could be a consequence of y-locus and other genes related to phenylpropanoid

739   metabolism. Interestingly, SV mutations have been found in *fw3.2/KLUH*, *FAS* and *SUN* that

740   supports the impact of SV on tomato phenotypic diversity (Alonge *et al.,* 2020, Dominguez *et al.,*

741   2020). Also some cryptic variation hidden in the Mesoamerican tomatoes may have emerged in

742   European tomatoes after generating new combinations and divergent selection by the traditional

743   farmers as found for the jointless trait in tomato (Soyk *et al.,* 2017, Soyk *et al.,* 2019, Alonge *et*

744   *al.,* 2020).

745   ## Impact on genebank and on farm variability management

746

747   Many of the few polymorphic genetic variants, within the very low diversity European vintage

748   tomatoes, appeared to be associated with phenotypic variation. This has implications for the

749   conservation efforts carried out by the genebanks. Thousands of European vintage tomatoes are

750   maintained in many of those genebanks. However, the cost of these conservation efforts could

751   be severely reduced if only these few polymorphic loci were taken into account. Of course, that

752   would ignore most variants, the ones found in very low frequencies, but conserving these low

753   frequency alleles, that in many cases would be neutral, and thus not associated with any

754   phenotypic variation, requires a sizeable investment. An alternative would be to identify the alleles

755   associated with a phenotype, however, that would require an exhaustive phenotypic

756   characterization.

757

758   Most of the European accessions analyzed here were collected from farmers in the 1950's to

759   1980's, and as landraces, they are appreciated, competitive, and cultivated varieties. The genetic

760   diversity of many other crops has also been maintained as landraces that evolved on-farm.

761   However, this diversity is continuously under threat by the introduction of new modern varieties

762 derived from a limited gene pool that have replaced the vintage varieties. It is generally believed
763 that most of the accessions in seed banks do not contribute to modern varieties (Tanksley and
764 McCouch, 1997) and this is also the case for tomato. Our Identification of the morphological and
765 genetic structure present in the European vintage tomato gene pool will be important to guarantee
766 access to that variability as the basis of the development of new varieties or evolved landraces in
767 the future (Casañas et al 2017).

## Conclusion

769 The entrepreneurship of many local European farmers during the last five hundred years has
770 managed to create a very complex and diverse set of tomato varieties adapted to different local
771 tastes and morphological preferences. These localized activities did not restrain those farmers
772 from importing other interesting novelties developed by other farmers elsewhere, thus generating
773 a much larger set of varietal tomato types that are characterized by an exuberant diversity that
774 serves as a variety for fresh, processing, and long shelf-life uses.

775

776 The current report shows that such a plethora of different types has been created from an original
777 material devoid of genetic diversity, by exploiting very few polymorphic loci subjected to balancing
778 selection.

779

## Supplementary Data

781

782 Fig. S1 Number of genomic positions with high coverage and number of variants per megabase
783 along the genome in all accessions.
784 Fig. S2. FastSTRUCTURE analysis.
785 Fig. S3. Introgressed regions along the genome detected in the modern genetic groups
786 Fig. S4. Major Allele Frequency spectrum in vintage, modern, and SCL_Peru_MA
787 Fig. S5. Rarefaction analysis of the expected heterozygosity for each genetic group
788 Fig. S6. Genome-wide linkage disequilibrium (LD) decay in wild, *S. lycopersicum* var. *cerasiforme*
789 (SLC), vintage, and modern accession groups.
790 Fig. S7. Hierarchical Principal Coordinate Analysis of European vintage tomato varieties
791 Fig. S8. Rarefaction analysis of the number of polymorphic variants (95% threshold)

792

793 Table S1. Accessions analyzed in this study.

796

## Acknowledgments

798

809

## Author contributions

811

812

813    JB and JC analysed the data and drafted the manuscript. J M-P, D S-M, PZ, RF analysed the

814    data. CP obtained the DNA, field trial phenotypic data and revised manuscript.

815    LF, JF, MP, JLR, ARiccini, SP, ARuggiero and MS obtained field trial phenotypic data. JC

816    obtained field trial phenotypic data, selected and provided vintage varieties.SG, AK, GG, MC, SG,

817    AM, MC, MJD, JP, selected and provided vintage varieties and revised the manuscript. DZ

818    coordinated the field trial. AJM and AG conceived and coordinated the study and revised the

819    manuscript.

820

## Data Availability

822    The sequence data can be found in NCBI (https://www.ncbi.nlm.nih.gov/sra) under the accession

823    number PRJNA722111.

824

# References

**Alonge M, Wang XG, Benoit M, K *et al.,*.** 2020. Major impacts of widespread structural variation on gene expression and crop improvement in tomato. Cell **182**, 145-161.

**Amador L, Santos B, Ríos D.** 2012. Variedades tradicionales de tomates de Canarias. Tenerife (Spain): Cultivos y Tecnología Agraría de Tenerife.

**Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, Selker EU, Cresko WA, Johnson EA**. 2008. Rapid SNP discovery and genetic mapping using sequenced RAD markers. Plos One **3**, e3376.

**Ballester AR, Molthoff J, de Vos R, Hekkert BTL, Orzaez D, Fernandez-Moreno JP, Tripodi P, Grandillo S, Martin C, Heldens J, Ykema M, Granell A, Bovy A**. 2010. Biochemical and molecular analysis of pink tomatoes: deregulated expression of the gene encoding transcription factor S1MYB12 leads to pink tomato fruit color. Plant Physiology **152**, 71-84.

**Bauchet G, Causse M**. 2012. Genetic diversity in tomato (*Solanum lycopersicum*) and its wild relatives. . In: Caliskan M, ed. Genetic diversity in plants. Rijeka (Croatia): InTech, 138-162

**Blanca J, Canizares J, Cordero L, Pascual L, Diez MJ, Nuez F**. 2012. Variation Revealed by SNP Genotyping and Morphology Provides Insight into the Origin of the Tomato. Plos One **7**, e48198

**Blanca J, Montero-Pau J, Sauvage C, Bauchet G, Illa E, Diez MJ, Francis D, Causse M, van der Knaap E, Canizares J**. 2015. Genomic variation in tomato, from wild ancestors to contemporary breeding accessions. BMC Genomics **16**, 257

**Boswell VR. 1937**. improvement and genetics of tomatoes, peppers and eggplant. In: Yearbook of the United States Department of Agriculture. Washington (MA): Goverment Print Office, 177-206

**Brewer MT, Moyseenko JB, Monforte AJ, van der Knaap E**. 2007. Morphological variation in tomato: a comprehensive study of quantitative trait loci controlling fruit shape and development. Journal of Experimental Botany **58**, 1339-1349.

**Casals J, Rivera A, Sabate J, del Castillo RR, Simo J**. 2019. cherry and fresh market tomatoes: differences in chemical, morphological, and sensory traits and their implications for consumer acceptance. Agronomy-Basel **9**, 18.

**Casanas F, Simo J, Casals J, Prohens J**. 2017. toward an evolved concept of landrace. Frontiers in Plant Science **8**, 7.

**Causse M, Desplat N, Pascual L, Le Paslier MC, Sauvage C, Bauchet G, Berard A, Bounon R, Tchoumakov M, Brunel D, Bouchet JP**. 2013. Whole genome resequencing in tomato reveals variation associated with introgression and breeding events. BMC Genomics **14**, 14.

**Chakrabarti M, Zhang N, Sauvage C, Munos S, Blanca J, Canizares J, Diez MJ, Schneider**

860 **R, Mazourek M, McClead J, Causse M, van der Knaap E**. 2013. A cytochrome P450
861 regulates a domestication trait in cultivated tomato. Proceedings of the National Academy of
862 Sciences of the United States of America **110**, 17125-17130.

863 **Delph LF, Kelly JK**. 2014. On the importance of balancing selection in plants. New Phytologist
864 **201**, 45-56.

865 **Dominguez M, Dugas E, Benchouaia M, Leduque B, Jimenez-Gomez JM, Colot V,**
866 **Quadrana L**. 2020. The impact of transposable elements on tomato diversity. Nature
867 Communications **11**, 4058

868 **Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, Mitchell SE**. 2011. A
869 robust, simple genotyping-by-sequencing (GBS) approach for high diversity species. Plos One
870 **6**, e19379

871 **Esposito S, Cardi T, Campanelli G, Sestili S, Diez MJ, Soler S, Prohens J, Tripodi P**. 2020.
872 ddRAD sequencing-based genotyping for population structure analysis in cultivated tomato
873 provides new insights into the genomic diversity of Mediterranean 'da serbo' type long shelf-life
874 germplasm. Horticulture Research **7**, 134

875 **Gallusci P, Hodgman C, Teyssier E, Seymour GB**. 2016. DNA methylation and chromatin
876 regulation during fleshy fruit development and ripening. Frontiers in Plant Science **7**, 14.

877 **Garcia-Martinez S, Andreani L, Garcia-Gusano M, Geuna F, Ruiz JJ**. 2006. Evaluation of
878 amplified fragment length polymorphism and simple sequence repeats for tomato germplasm
879 fingerprinting: utility for grouping closely related traditional cultivars. Genome **49**, 648-656.

880 **Garrison E, Marth G**. 2012. Haplotype-Based Variant Detection from Short-Read Sequencing.
881 ArXiv, 1207.3907.

882 **Gogarten SM, Sofer T, Chen H, Yu CY, Brody JA, Thornton TA, Rice KM, Conomos MP**.
883 2019. Genetic association testing using the GENESIS R/Bioconductor package. Bioinformatics
884 **35**, 5346-5348.

885 **Grandillo S, Ku HM, Tanksley SD**. 1999. Identifying the loci responsible for natural variation in
886 fruit size and shape in tomato. Theoretical and Applied Genetics **99**, 978-987.

887 **Harvey M**. 2004. Exploring the tomato: transformations of nature, society and economy.
888 Cheltenham(UK): Edward Elgar Publishing.

889 **Huson DH, Bryant D**. 2006. Application of phylogenetic networks in evolutionary studies.
890 Molecular Biology and Evolution **23**, 254-267.

891 **Kosman E, Leonard KJ**. 2005. Similarity coefficients for molecular markers in studies of
892 genetic relationships between individuals for haploid, diploid, and polyploid species. Molecular
893 Ecology **14**, 415-424.

894 **Krzanowski WJ, Krzanowski WJ**. 2000. Principles of multivariate analysis: a user's
895 perspective. 2nd Edition. New York(NY): Oxford University Press.

896  **Li FQ, Vallabhaneni R, Wurtzel ET**. 2008. PSY3, a new member of the phytoene synthase
897  gene family conserved in the poaceae and regulator of abiotic stress-induced root
898  carotenogenesis. Plant Physiology **146**, 1333-1345.

899  **Li H**. 2011. Improving SNP discovery by base alignment quality. Bioinformatics **27**, 1157-1158.

900  **Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin
901  R, Genome Project Data P**. 2009. The Sequence Alignment/Map format and SAMtools.
902  Bioinformatics **25**, 2078-2079.

903  **Li H**. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
904  ArXiv, 1303.3997.

905  **Lin T, Zhu GT, Zhang JH, K** *et al.,*. 2014. Genomic analyses provide insights into the history of
906  tomato breeding. Nature Genetics **46**, 1220-1226.

907  **Liu JP, Van Eck J, Cong B, Tanksley SD**. 2002. A new class of regulatory genes underlying
908  the cause of pear-shaped tomato fruit. Proceedings of the National Academy of Sciences of the
909  United States of America **99**, 13302-13306.

910  **Luckwill LC**. 1943. The genus *Lycopersicon*; an historical, biological, and taxonomic survey of
911  the wild and cultivated tomatoes. Aberdeen University Studies 120. Aberdeen: The University
912  Press.

913  **Mazzucato A, Papa R, Bitocchi E, Mosconi P, Nanni L, Negri V, Picarella ME, Siligato F,
914  Soressi GP, Tiranti B, Veronesi F**. 2008. Genetic diversity, structure and marker-trait
915  associations in a collection of Italian tomato (*Solanum lycopersicum* L.) landraces. Theoretical
916  and Applied Genetics **116**, 657-669.

917  **McCue GA**. 1952. The history of the use of the tomato: an annotated bibliography. Annals
918  Missouri Botanical Garden **39,** 289-348.

919  **Menda N, Strickler SR, Edwards JD, Bombarely A, Dunham DM, Martin GB, Mejia L,
920  Hutton SF, Havey MJ, Maxwell DP, Mueller LA**. 2014. Analysis of wild-species introgressions
921  in tomato inbreds uncovers ancestral origins. BMC Plant Biology **14**, 16.

922  **Mu Q, Huang ZJ, Chakrabarti M, Illa-Berenguer E, Liu XX, Wang YP, Ramos A, van der
923  Knaap E**. 2017. Fruit weight is controlled by Cell Size Regulator encoding a novel protein that is
924  expressed in maturing tomato fruits. Plos Genetics **13**, e1006930.

925  **Muggeo VMR**. 2003. Estimating regression models with unknown break-points. Statistics in
926  Medicine **22**, 3055-3071.

927  **Munos S, Ranc N, Botton E, Berard A, Rolland S, Duffe P, Carretero Y, Le Paslier MC,
928  Delalande C, Bouzayen M, Brunel D, Causse M**. 2011. increase in tomato locule number is
929  controlled by two Single-Nucleotide Polymorphisms located near WUSCHEL. Plant Physiology
930  **156**, 2244-2254.

931  **Park YH, West MAL, St Clair DA**. 2004. Evaluation of AFLPs for germplasm fingerprinting and

932    assessment of genetic diversity in cultivars of tomato (*Lycopersicon esculentum* L.). Genome
933    **47**, 510-518.

934    **Peakall R, Smouse PE**. 2006. Genalex 6: Genetic analysis in Excel. Population genetic
935    software for teaching and research. Molecular Ecology Notes 6, 288–95.

936

937    **Peakall R, Smouse PE**. 2012. GenAlEx 6.5: genetic analysis in Excel. Population genetic
938    software for teaching and research-an update. Bioinformatics **28**, 2537-2539.

939    **Pons C, Blanca B, Cañizares J, Ziasolo P, Finkers R, Rambla JL; da Silva GE, Zacarias L,**
940    **Monforte AJ, Granell A**. 2017. Variability in fruit ripening within the European traditional pool of
941    tomato varieties. SOLCUC2017. Valencia, Spain, September 6-7, 2017.

942    **Raj A, Stephens M, Pritchard JK**. 2014. fastSTRUCTURE: Variational inference of population
943    structure in large SNP data sets. Genetics **197**, 573-589.

944    **Razifard H, Ramos A, Della Valle AL, Bodary C, Goetz E, Manser EJ, Li X, Zhang L, Visa S,**
945    **Tieman D, van der Knaap E, Caicedo AL**. 2020. Genomic evidence for complex domestication
946    history of the cultivated tomato in Latin America. Molecular Biology and Evolution **37**, 1118-
947    1132.

948    **Rick CM, Fobes JF**. 1975.Allozyme variation in the cultivated tomato and closely related
949    species. Bulletin of the Torrey Botanical Club **102**, 376-384.

950    **Rick CM, Zobel RW, Fobes JF**. 1974. Four peroxidase loci in red-fruited tomato species:
951    genetics and geographic distribution. Proceedings of the National Academy of Sciences of the
952    United States of America **71**, 835-839.

953    **Robbins MD, Sim SC, Yang WC, Van Deynze A, van der Knaap E, Joobeur T, Francis DM**.
954    2011. Mapping and linkage disequilibrium analysis with a genome-wide collection of SNPs that
955    detect polymorphism in cultivated tomato. Journal of Experimental Botany **62**, 1831-1845.

956    **Rodriguez GR, Munos S, Anderson C, Sim SC, Michel A, Causse M, Gardener BBM,**
957    **Francis D, van der Knaap E**. 2011. Distribution of SUN, OVATE, LC, and FAS in the Tomato
958    Germplasm and the Relationship to Fruit Shape Diversity. Plant Physiology **156**, 275-285.

959    **Rogers AR, Huff C**. 2009. Linkage Disequilibrium Between Loci With Unknown Phase.
960    Genetics **182**, 839-844.

961    **Sacco A, Ruggieri V, Parisi M, Festa G, Rigano MM, Picarella ME, Mazzucato A, Barone A**.
962    2015. Exploring a tomato landraces collection for fruit-related traits by the aid of a high-
963    throughput genomic platform. Plos One **10**, e0137139.

964    **Sahagún B**. 1577. Historia general de las cosas de Nueva España por el fray Bernardino de
965    Sahagún: el Códice Florentino. Libro X: del pueblo, sus virtudes y vicios, y otras naciones.
966    Florence: Biblioteca Medicea Laurenziana

967  **Sanfuentes-Echeverría O**. 2006. Europa y su percepción del nuevo mundo a través de las
968  especies comestibles y los espacios americanos en el siglo XVI Historia (Santiago) **39**, 531-
969  556.

970  **Sato S, Tabata S, Hirakawa H,** *et al.,* 2012. The tomato genome sequence provides insights
971  into fleshy fruit evolution. Nature **485**, 635-641.

972  **Seabold S, Perktold J**. 2010. Statsmodels: econometric and statistical modeling with Python.
973  In: Walt S, Willman J Eds. 9th Python in Science Conference (SciPy 2010). Austin, TX, June 28-
974  July 3, 2010 Proceedings, **92-96**.

975  **Sim SC, Robbins MD, van Deynze A, Michel AP, Francis DM**. 2011. Population structure and
976  genetic differentiation associated with breeding history and selection in tomato *(Solanum*
977  *lycopersicum* L.). *Heredity* **106**, 927-35.

978  **Sim SC, Durstewitz G, Plieske J, Wieseke R, Ganal MW, Van Deynze A, Hamilton JP, Buell**
979  **CR, Causse M, Wijeratne S, Francis DM**. 2012a. Development of a large SNP genotyping
980  array and generation of high-density genetic maps in tomato. Plos One **7**, e40563

981  **Sim SC, Van Deynze A, Stoffel K, Douches DS, Zarka D, Ganal MW, Chetelat RT, Hutton**
982  **SF, Scott JW, Gardner RG, Panthee DR, Mutschler M, Myers JR, Francis DM**. 2012b. High-
983  density SNP genotyping of tomato (*Solanum lycopersicum* L.) reveals patterns of genetic
984  variation due to breeding. Plos One **7**, e45520.

985  **Soyk S, Lemmon ZH, Oved M, Fisher J, Liberatore KL, Park SJ, Goren A, Jiang K, Ramos**
986  **A, van der Knaap E, Van Eck J, Zamir D, Eshed Y, Lippman ZB**. 2017. Bypassing negative
987  epistasis on yield in tomato imposed by a domestication gene. Cell **169**, 1142-1155.

988  **Soyk S, Lemmon ZH, Sedlazeck FJ, Jimenez-Gomez JM, Alonge M, Hutton SF, Van Eck J,**
989  **Schatz MC, Lippman ZB**. 2019. Duplication of a domestication locus neutralized a cryptic
990  variant that caused a breeding barrier in tomato. Nature Plants **5**, 471-479.

991  **Tanksley SD, McCouch SR**. 1997. Seed banks and molecular maps: Unlocking genetic
992  potential from the wild. Science **277**, 1063-1066.

993  **Williams CE, Stclair DA**. 1993. Phenetic relationships and levels of variability detected by
994  Restriction-Fragment-Length-Polymorphism and Random Amplified Polymorphic DNA analysis
995  of cultivated and wild accessions of *Lycopersicon esculentum*. Genome **36**, 619-630.

996  **Xiao H, Jiang N, Schaffner E, Stockinger EJ, van der Knaap E**. 2008. A retrotransposon-
997  mediated gene duplication underlies morphological variation of tomato fruit. Science **319**, 1527-
998  1530.

999  **Xu C, Liberatore KL, MacAlister CA, Huang ZJ, Chu YH, Jiang K, Brooks C, Ogawa-**
1000 **Ohnishi M, Xiong GY, Pauly M, Van Eck J, Matsubayashi Y, van der Knaap E, Lippman**
1001 **ZB**. 2015. A cascade of arabinosyltransferases controls shoot meristem size in tomato. Nature
1002 Genetics **47**, 784-792.

1003 **Yuste-Lisbona FJ, Fernandez-Lozano A, Pineda B, Bretones S, Ortiz-Atienza A, Garcia-**

1004 **Sogo B, Muller NA, Angosto T, Capel J, Moreno V, Jimenez-Gomez JM, Lozano R**. 2020.
1005 ENO regulates tomato fruit size through the floral meristem development network. Proceedings
1006 of the National Academy of Sciences of the United States of America **117**, 8187-8195.

1007 **Zheng XW, Levine D, Shen J, Gogarten SM, Laurie C, Weir BS**. 2012. A high-performance
1008 computing toolset for relatedness and principal component analysis of SNP data. Bioinformatics
1009 **28**, 3326-3328.

1010

# Figure Legends

1012

1013 Fig. 1. Principal Coordinate Analysis (PCoA) including cultivated tomato (*Solanum lycopersicum*
1014 var. *lycopersicum*, SLL): vintage European tomato, modern cultivars with different culinary use
1015 (fresh, processing and long shelf life, lsl), *S. lycopersicum* var. cerasiforme (SLC) from different
1016 origin [Peru, Mesoamerica (MA) Ecuador (Ecu)], together with several American wild relatives: *S.*
1017 *pimpinellifolium* (SP), *S. cheesmaniae*, *S. galapagense* (Galápagos), *S. peruvianum*, *S.*
1018 *chmielewskii* and *S. habrochaites* (green) and SPxSL hybrids. The modern cultivar Heinz1706
1019 was included as reference. (A) First and second principal components (dim1 and dim2) from the
1020 PCoA using all accessions analyzed in this study. (B) First and third components (dim1 and dim3)
1021 from the same PCoA. C) First and second components (dim1 and dim2) from PCoA using only S.
1022 lycopersicum var. lycopersicum samples. D) First and third components (dim1 and dim3) from the
1023 previous PCoA The percentage of explained variance for each principal component is indicated
1024 on each axis.

1025

1026 Fig. 2. Principal Coordinate Analysis (PCoA) of modern cultivars. (A) and (B) the three first
1027 principal components (dim1, dim2 and dim3) from the PCoA considering all modern cultivars and
1028 cv. Heinz1706 as reference. (C) and (D) PCoA including only modern fresh 2 and Long Shelf Llife
1029 (LSL) and modern processing genetic groups. The variance accounted for each principal
1030 component is depicted on each axis.

1031

1032 Fig. 3. Genetic diversity for the rank1 genetic groups. (A) Genetic diversity estimated by the
1033 expected heterozygosity and the percentage of polymorphic variants (95% threshold). The
1034 indexes were calculated 100 times taking 20 samples at random from each genetic group. The
1035 mean and standard deviation are shown. (B) Rarefaction analysis of the number of variants found
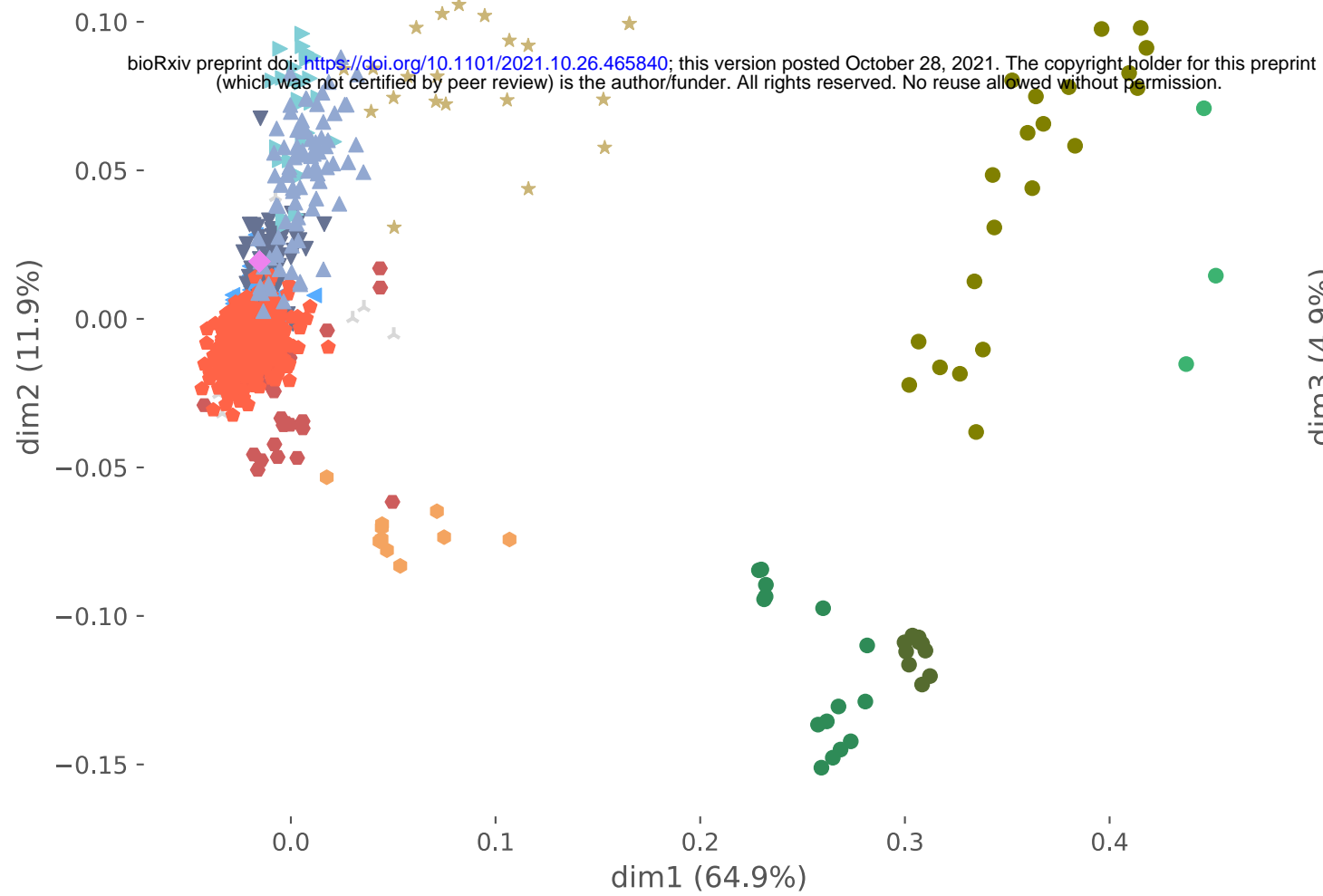1036 in each genetic group. Axis X shows the number of samples, Axis Y shows the number of variants.
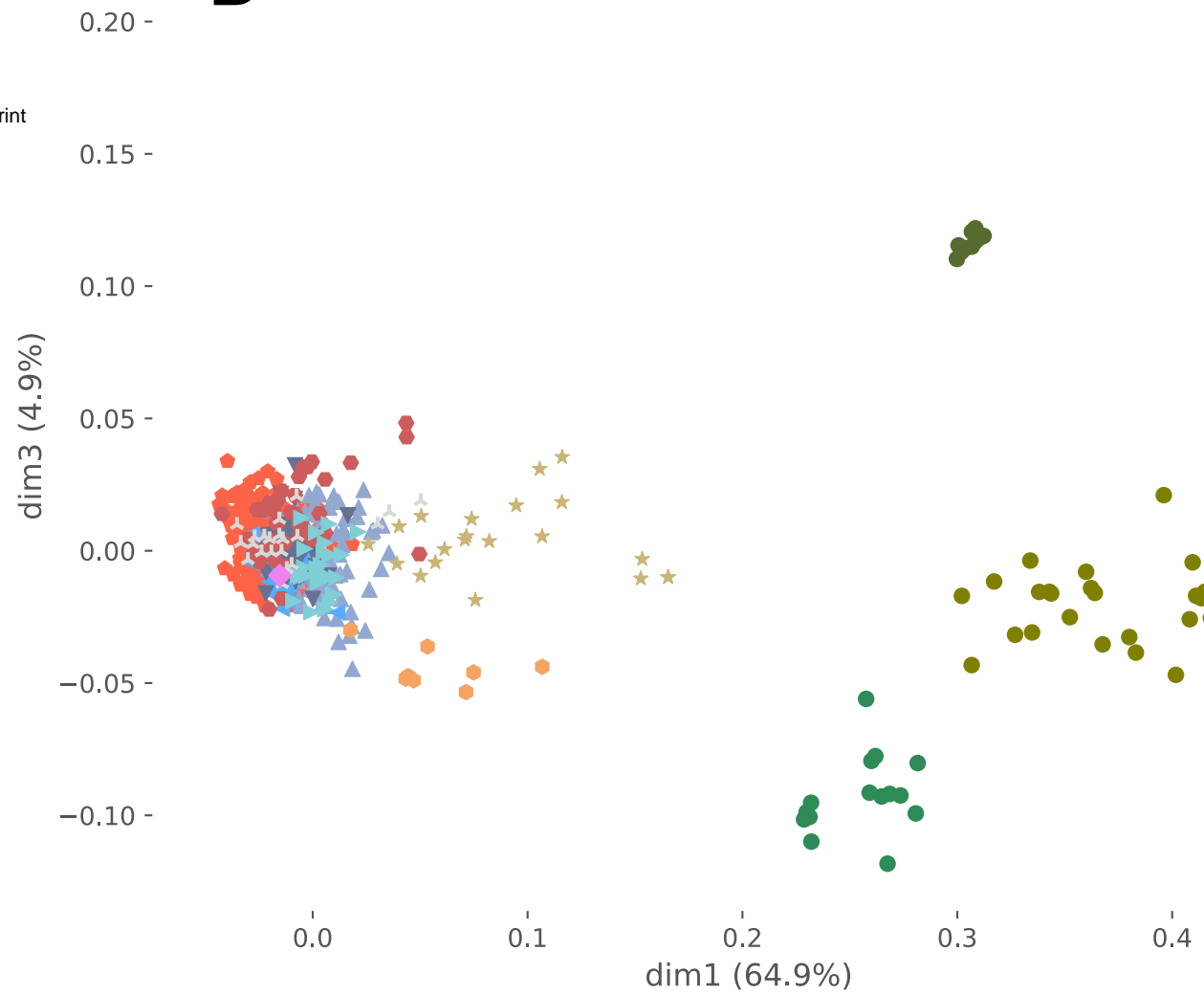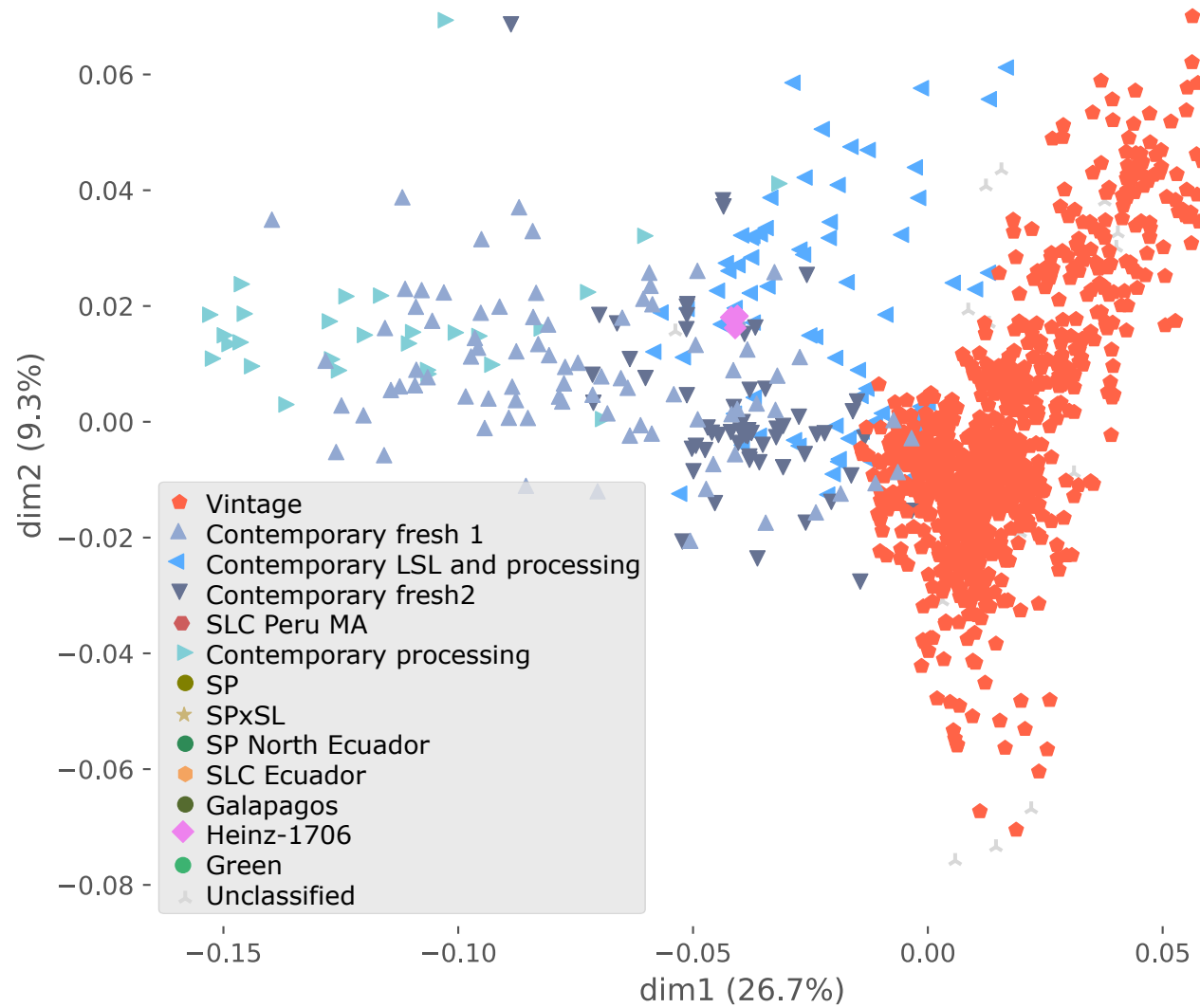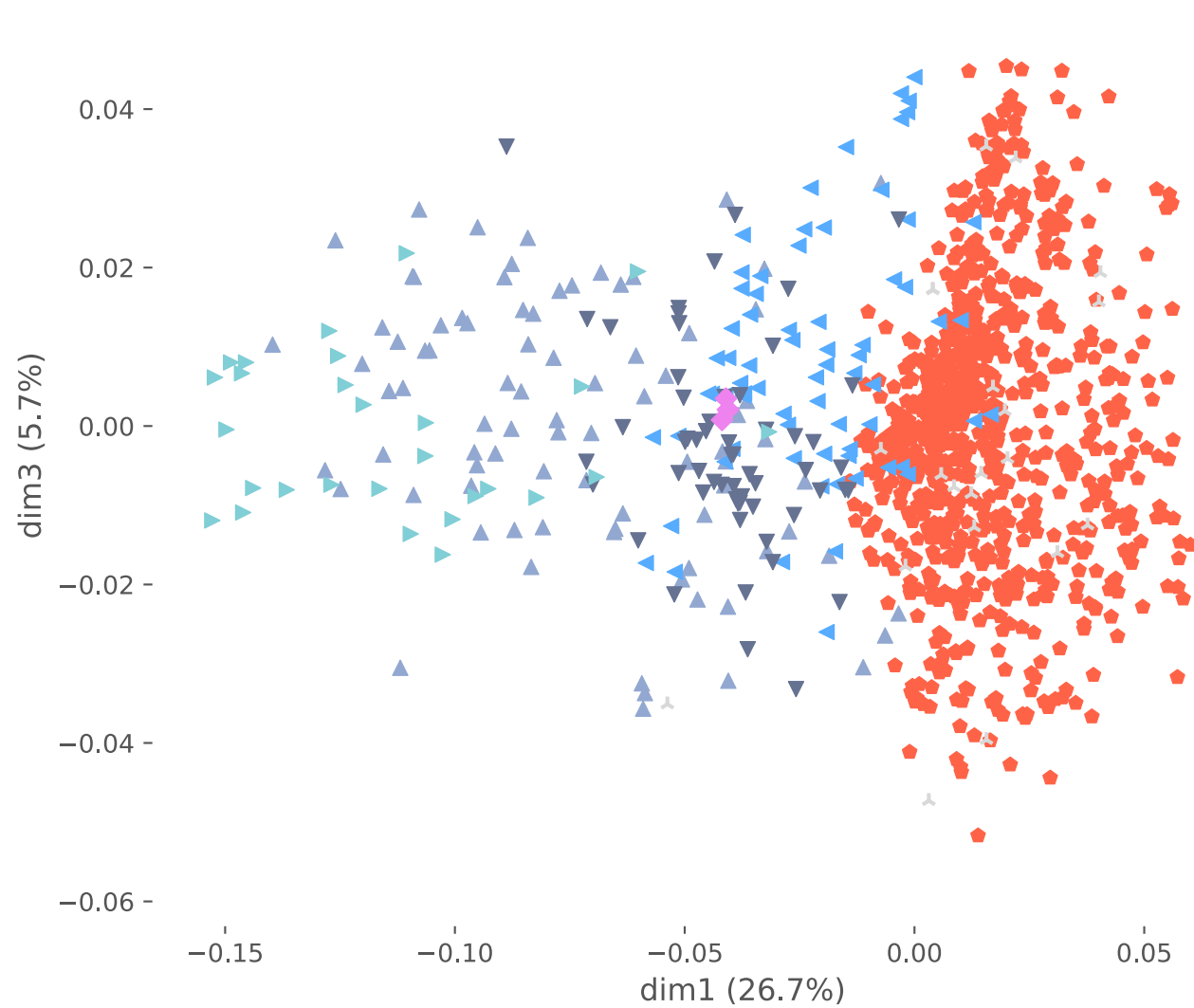
1037

1038    Fig. 4. Allele frequencies across the genome in Vintage genetic groups and their relationship with

1039    phenotypic diversity. (A) Clustering of genetic groups based on allele frequencies. Allele

1040    frequency of the major allele within each genetic group is indicated by a density color according

1041    the legend (blue, frequency=0, to white, frequency=1. (B) Distribution of the different traits within

1042    genetic groups. (C) Statistical significance indicated by a colored gradient of -log(p) values of the

1043    SNP-trait associations by Genome-Wide Association Analysis.

1044

1045    Fig. 5. Evolutionary relationships between vintage European tomato, Modern tomato and Peru

1046    and Mesoamerica *Solanum lycopersicum* var. *cerasiforme* (SLC), *S. pimpinellifolium* (SP) and the

1047    hybrids SPxSL. Split network based on the Dest distances between genetic groups. The country

1048    of origin of accessions within each genetic group is represented by a pie chart depicted in the

1049    bottom left. (A) Zoom only in European modern and vintage tomatoes. (B) Zoom on American

1050    ancestral and wild tomatoes. Each edge of the network represents a split of the accessions based

1051    on one or more characteristics. If there was no conflict, each split was represented by a single

1052    edge, while in the case of contradictory patterns the partition was represented by a set of parallel

1053    edges. The edge lengths represent the weight of each split, which is equivalent to the distance
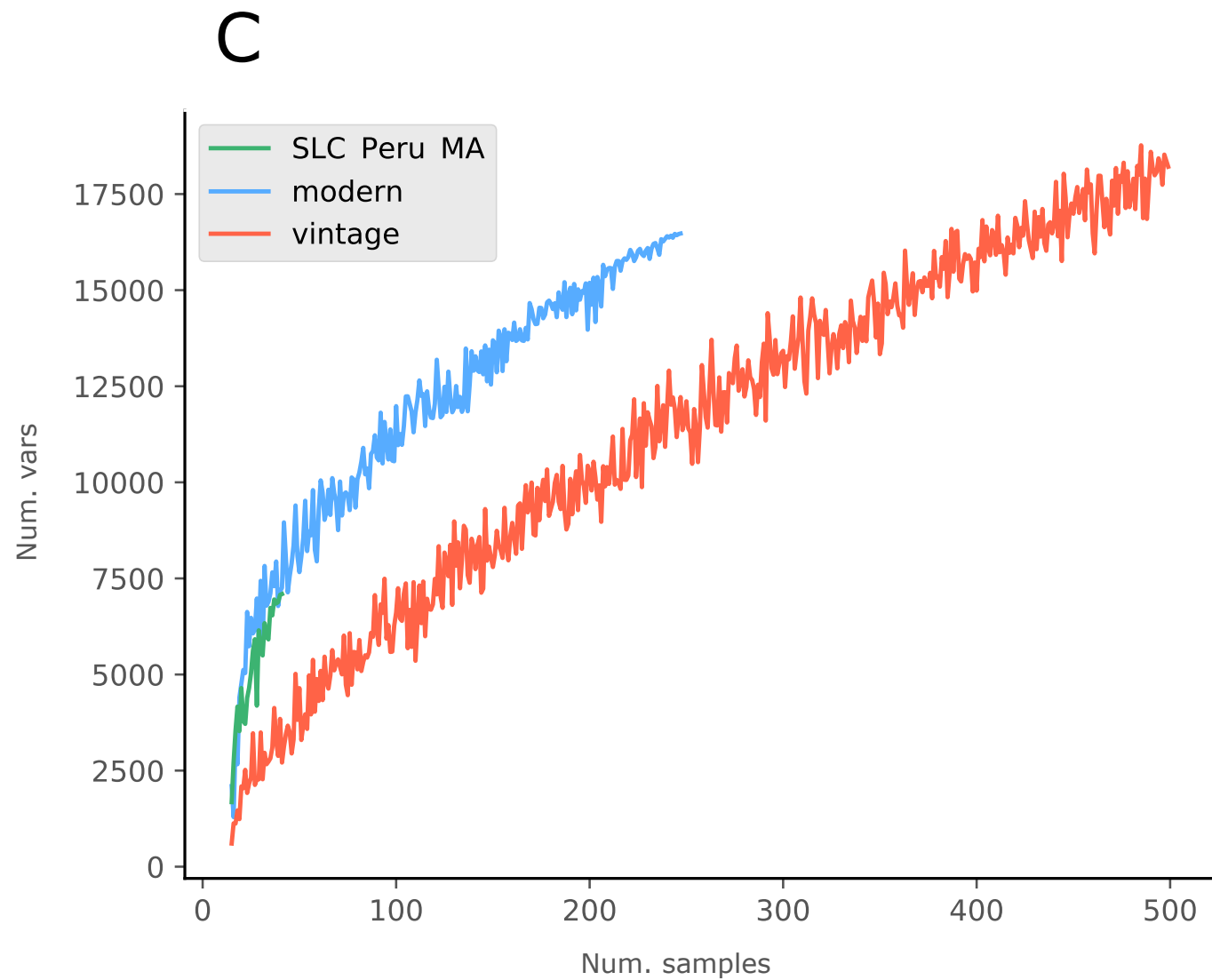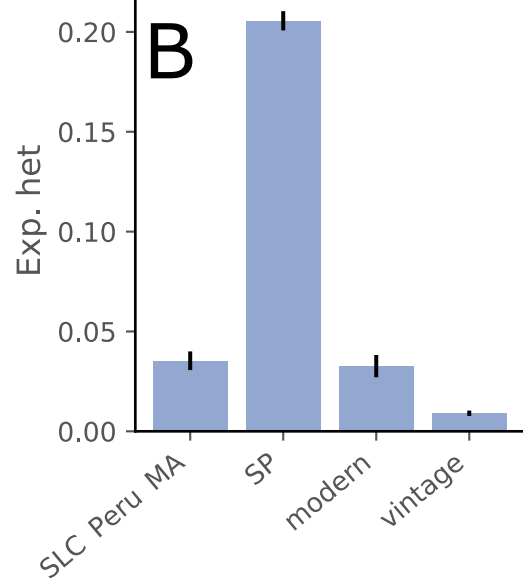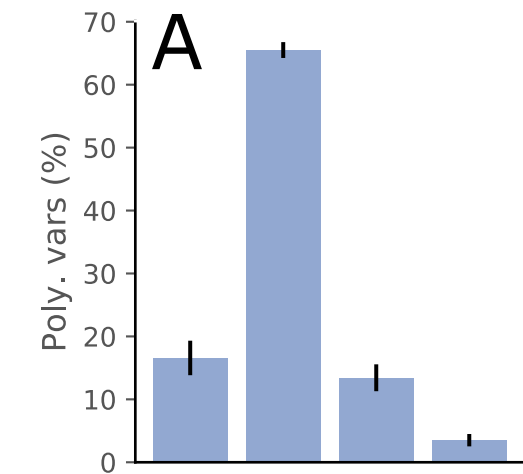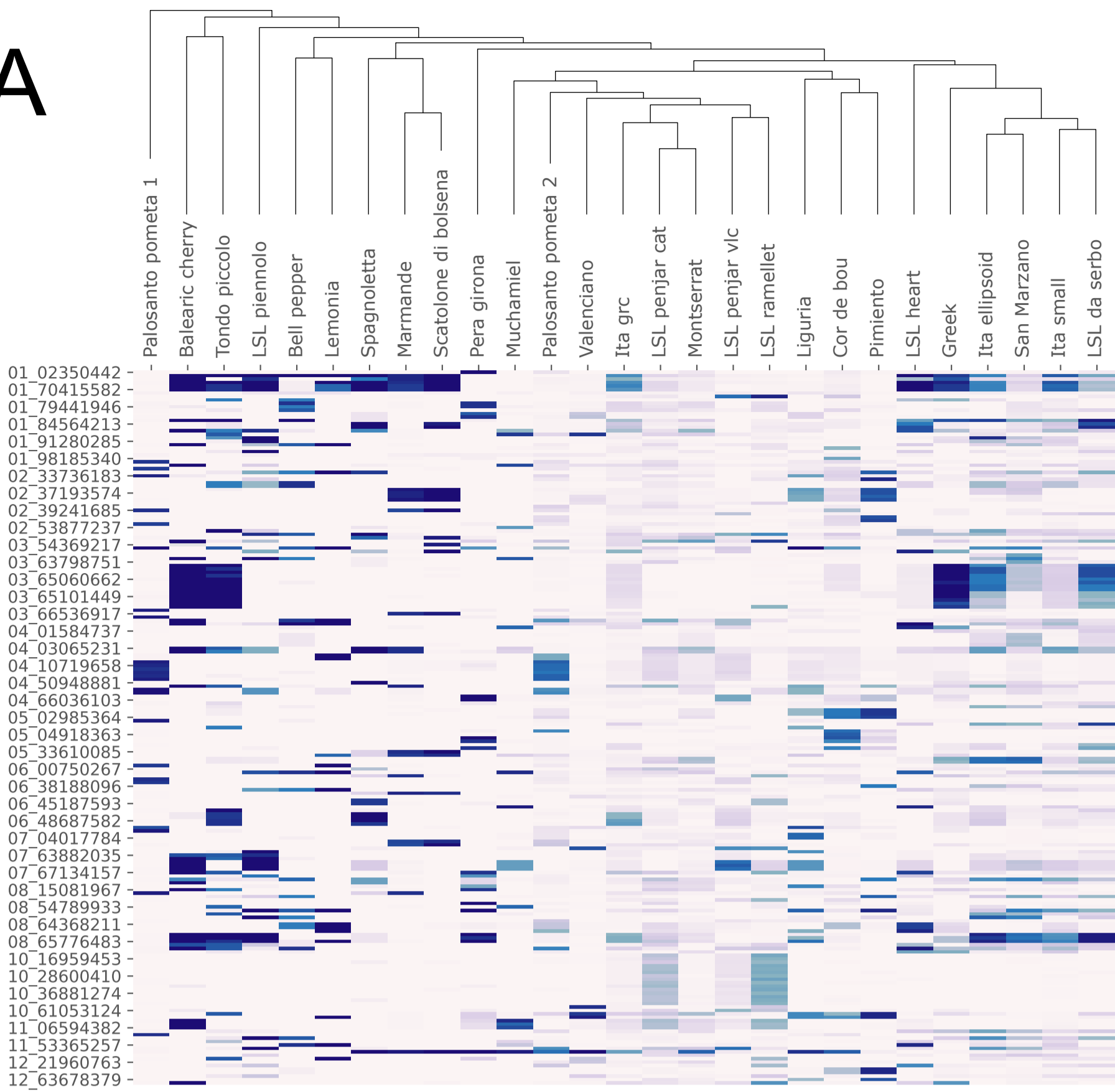
1054    between groups.

1055

Legend:
▲ Contemporary fresh 1
◀ Contemporary LSL and processing
▼ Contemporary fresh2
▶ Contemporary processing
◆ Heinz-1706
⅄ Unclassified

A

LSL ramellet
LSL penjar VLC
modern fresh
Palosanto-pometa 1
Palosanto-pometa 2
Modern LSL & processing
Liguria
LSL penjar cat
montserrat
Muchamiel
Valenciano
Scatolone di Bolsena
Marmande
LSL heart
SLC Peru MA
Cor de bou
Ita Grc
Ita small
Pimiento
San Marzano
Grc
Ita ellipsoid
LSL da serbo

ESP
ITA
GRC
FRA
other

B
SP
SPxSL
SP NEcu
modern vintage SLC Peru MA