



**HAL**  
open science

## **Solanum pan-genomics and pan-genetics reveal paralogs as contingencies in crop engineering**

Matthias Benoit, Katharine Jenike, James Satterlee, Srividya Ramakrishnan, Iacopo Gentile, Anat Hendelman, Michael Passalacqua, Hamsini Suresh, Hagai Shohat, Gina Robitaille, et al.

► **To cite this version:**

Matthias Benoit, Katharine Jenike, James Satterlee, Srividya Ramakrishnan, Iacopo Gentile, et al.. Solanum pan-genomics and pan-genetics reveal paralogs as contingencies in crop engineering. 2024. hal-04846989

**HAL Id: hal-04846989**

**<https://hal.science/hal-04846989v1>**

Preprint submitted on 18 Dec 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 ***Solanum* pan-genomics and pan-genetics reveal paralogs as contingencies in**  
2 **crop engineering**

3  
4 Matthias Benoit<sup>1,\*†</sup>, Katharine M. Jenike<sup>2,3\*</sup>, James W. Satterlee<sup>1,4,o</sup>, Srividya Ramakrishnan<sup>3,o</sup>,  
5 Iacopo Gentile<sup>5,o</sup>, Anat Hendelman<sup>1,4,o</sup>, Michael J. Passalacqua<sup>5</sup>, Hamsini Suresh<sup>5</sup>, Hagai Shohat<sup>4</sup>,  
6 Gina M. Robitaille<sup>1,4</sup>, Blaine Fitzgerald<sup>1,4</sup>, Michael Alonge<sup>3,†</sup>, Xingang Wang<sup>4,†</sup>, Ryan Santos<sup>4,†</sup>,  
7 Jia He<sup>1,4</sup>, Shujun Ou<sup>3,†</sup>, Hezi Golan<sup>6</sup>, Yumi Green<sup>7</sup>, Kerry Swartwood<sup>7</sup>, Gina P. Sierra<sup>8</sup>, Andres  
8 Orejuela<sup>9</sup>, Federico Roda<sup>8</sup>, Sara Goodwin<sup>4</sup>, W. Richard McCombie<sup>4</sup>, Elizabeth B. Kizito<sup>10</sup>, Edeline  
9 Gagnon<sup>11,12,†</sup>, Sandra Knapp<sup>13</sup>, Tiina E. Särkinen<sup>12</sup>, Amy Frary<sup>14</sup>, Jesse Gillis<sup>4,15,#</sup>, Joyce Van  
10 Eck<sup>7,16,#</sup>, Michael C. Schatz<sup>2,3,#</sup>, Zachary B. Lippman<sup>1,4,5,#</sup>

- 11  
12 1. Howard Hughes Medical Institute, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY,  
13 USA  
14 2. Department of Genetic Medicine, Johns Hopkins School of Medicine, Baltimore, MD, USA  
15 3. Department of Computer Science, Johns Hopkins University, Baltimore, MD, USA  
16 4. Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA  
17 5. School of Biological Sciences, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA  
18 6. SiteKicks.ai, Setauket, NY, USA  
19 7. Boyce Thompson Institute, Ithaca, NY, USA  
20 8. Max Planck Tandem Group, Facultad de Ciencias, Universidad Nacional de Colombia,  
21 Bogotá, Colombia  
22 9. Departamento de Biología, Facultad de Ciencias Exactas y Naturales, Universidad de  
23 Cartagena, Cartagena de Indias, Colombia  
24 10. Faculty of Agricultural Sciences, Uganda Christian University, Mukono, Uganda  
25 11. Department of Integrative Biology, University of Guelph, Ontario, Canada  
26 12. Royal Botanic Garden Edinburgh, 20A Inverleith Row, Edinburgh, EH3 5LR, UK  
27 13. Natural History Museum, London, UK  
28 14. Department of Biological Sciences, Mount Holyoke College, South Hadley, MA, USA

29 15. Physiology Department and Donnelly Centre for Cellular and Biomolecular Research,  
30 University of Toronto, Toronto, ON, Canada

31 16. Plant Breeding and Genetics Section, School of Integrative Plant Science, Cornell University,  
32 Ithaca, NY, USA

33

34 \* These authors contributed equally

35 ° These authors contributed equally

36

37 # Corresponding authors

38

39 † Present address: LIPME, Université de Toulouse, INRAE, CNRS, Castanet-Tolosan, France

40 (M.B.); Ohalo Genetics, Aptos, CA, USA (M.A. and X.W.); Verve Therapeutics, Boston, MA,

41 USA (R.S.); Department of Molecular Genetics, Ohio State University, Columbus, OH, USA

42 (S.O.); School of Life Sciences, Technical University of Munich, Freising, Germany (E.G.).

43

44 **Keywords:** pan-genome, *Solanum*, tomato, potato, eggplant, indigenous crops, domestication,

45 paralogs, gene expression, *cis*-regulatory, haplotypes, CRISPR, QTL, breeding, epistasis

46 **ABSTRACT**

47 Pan-genomics and genome editing technologies are revolutionizing the breeding of globally  
48 cultivated crops. A transformative opportunity lies in the reciprocal exchange of genotype-to-  
49 phenotype knowledge of agricultural traits between these major crops and hundreds of locally  
50 cultivated indigenous crops, thereby enhancing the diversity and resilience of our food system.  
51 However, species-specific genetic variants and their interactions with desired natural or engineered  
52 mutations pose barriers to achieving predictable phenotypic effects, even between closely related  
53 crops or genotypes. Here, by establishing a pan-genome of the crop-rich genus *Solanum* and  
54 integrating functional genomics and genetics, we show that gene duplication and subsequent  
55 paralog diversification are a major obstacle to genotype-phenotype predictability. Despite broad  
56 conservation of gene macrosynteny among chromosome-scale references for 22 species, including  
57 13 indigenous crops, hundreds of global and lineage-specific gene duplications exhibited dynamic  
58 evolutionary trajectories in paralog sequence, expression, and function, including among members  
59 of key domestication gene families. Extending our pan-genome with 10 cultivars of African  
60 eggplant and leveraging quantitative genetics and genome editing, we uncovered an intricate  
61 history of paralog emergence and evolution within this indigenous crop. The loss of an ancient  
62 redundant paralog of the classical regulator of stem cell proliferation and fruit organ number,  
63 *CLAVATA3 (CLV3)*, was compensated by a lineage-specific tandem duplication. Subsequent  
64 pseudogenization of the derived copy followed by a cultivar-specific structural variant resulted in  
65 a single fused functional copy of *CLV3* that modifies locule number alongside a newly identified  
66 gene controlling the same trait. Our findings demonstrate that paralog diversifications over short  
67 evolutionary periods are critical yet underexplored contingencies in trait evolvability and  
68 independent crop domestication histories. Unraveling these contingencies is crucial for translating  
69 genotype-to-phenotype relationships across related species.

## 70 INTRODUCTION

71 Global food production is currently based on fewer than 10 intensively bred commodity  
72 crops from only three plant families<sup>1</sup>: grasses (corn, rice, sugarcane, wheat), legumes (soybean),  
73 and nightshades (potato, tomato). In contrast, indigenous crops comprise a large, heterogeneous  
74 group of hundreds of species which could contribute to agricultural biodiversity and resilience<sup>2</sup>.  
75 Many indigenous crops belong to the same families as the major crops but are differentiated by  
76 their narrower range of cultivation and scale of production<sup>3</sup>. For instance, the grasses millet  
77 (*Eleusine coracana*) and teff (*Eragrostis tef*) and the legumes cowpea (*Vigna unguiculata*) and  
78 pigeonpea (*Cajanus cajan*) are locally adapted and important to diets in specific regions of Africa  
79 and Asia<sup>4-6</sup>. Within the nightshade (*Solanaceae*) family, the genus *Solanum* alone contains dozens  
80 of crop and wild species cultivated in specific regions of Africa and South America for their leaves  
81 and/or, fruits, including African eggplant (*S. aethiopicum*), naranjilla (*S. quitoense*), African black  
82 nightshade (*S. scabrum*) and pepino (*S. muricatum*)<sup>7,8</sup>.

83 Indigenous crops are viewed through several different lenses—agricultural,  
84 ethnobotanical, and scientific—each with its own unique biases and objectives<sup>2,3,9,10</sup>. Bridging and  
85 harmonizing these viewpoints offers an opportunity to better serve local communities and  
86 encourage broader adoption for industrialization. Breeding of indigenous crops has been limited  
87 relative to global commodity crops. It is widely assumed that decades of research on major crops,  
88 along with advances in genome sequencing and genome editing technologies, can be leveraged to  
89 address residual undesirable ancestral traits that limit productivity of indigenous, locally adapted  
90 crops<sup>11,12</sup>. Engineering beneficial mutations could help rapidly expand the diversity of food species  
91 beyond our current genetically narrow, industrialized agricultural systems<sup>2,13</sup>. Despite great  
92 progress in genome engineering technologies, however, background dependencies—species-  
93 specific genetic modifiers that lead to unpredictable phenotypic outcomes even between closely  
94 related species or varieties—remain underappreciated barriers<sup>14</sup>. Indeed, plant breeders have long  
95 lamented that beneficial alleles and quantitative trait loci (QTLs) often underperform when  
96 transferred to different backgrounds due to interactions among variants<sup>15,16</sup>—a challenge that will  
97 persist with genome editing<sup>17,18</sup>.

98 Our recent tomato pan-genome and associated functional genetics have demonstrated that  
99 gene duplications can be potent sources of background modifiers<sup>19,20</sup>. Duplications initially result  
100 in genetic redundancy which permits the accumulation of mutations in coding and *cis*-regulatory

101 sequences through genetic drift. Consequently, paralog redundancy can degrade, leading to three  
102 canonical outcomes over long evolutionary time: gene loss (pseudogenization), partitioning of  
103 ancestral functions (subfunctionalization) or gain of new functions (neofunctionalization)<sup>21,22</sup>.  
104 However, the dynamics of how paralogs diverge over shorter time frames, in their sequences,  
105 expression patterns, and functions, is less well understood. Genomic and functional dissections of  
106 paralogs have largely been limited to within individual species or between widely diverged  
107 lineages, and thus have not captured more intermediate trajectories and variable functional  
108 consequences of paralog divergence. A deeper understanding of paralog histories and their  
109 potentially interdependent relationships could provide greater predictability of phenotypic  
110 outcomes when translating genetic knowledge between closely related species. Here, we present a  
111 *Solanum* pan-genome and leverage this resource in conjunction with pan-genetics, forward and  
112 reverse genetics across species, to comprehensively analyze paralog evolutionary dynamics,  
113 demonstrating the value of resolving these underexplored contingencies as we strive to improve  
114 indigenous crops for local and climate change adapted agriculture.

115

## 116 RESULTS

### 117 A chromosome-scale pan-genome of the genus *Solanum*.

118 *Solanum* is one of the most species-rich, ecologically diverse and economically important  
119 plant genera<sup>7,8</sup>. The genus includes the major crops eggplant (*S. melongena*), potato (*S.*  
120 *tuberosum*), and tomato (*S. lycopersicum*) and at least 24 indigenous crops, including African  
121 eggplant (*S. aethiopicum*), naranjilla (*S. quitoense*) and pepino (*S. muricatum*)<sup>23</sup>. Spanning  
122 approximately 16-44 Ma of evolution<sup>24,25</sup>, the diversity of the genus *Solanum*, along with existing  
123 genomic and genetic tools in specific species<sup>26,27</sup>, makes it a leading system to study paralog  
124 evolution over short evolutionary time scales. We selected 22 species encompassing a broad  
125 phylogenetic sample of the ecological (**Fig. 1a**), phenotypic (**Fig. 1b, Extended Data Fig. 1a**),  
126 and taxonomic (**Fig. 1c, Supplementary Table 1**) diversity within *Solanum*, including regionally  
127 important indigenous crop and ornamental species and several of their wild progenitors. These  
128 species are grouped into four main categories that reflect the spectrum of plant use and  
129 domestication: wild (W); locally-important, consumed (C); ornamental (O); domesticated food  
130 crop (D) (**Fig. 1a,b**). Using PacBio HiFi sequencing and other long-range scaffolding data, we  
131 assembled chromosome-scale genomes for all 22 species, including phased haplotypes of the

132 clonally-propagated and highly heterozygous pepino, for a total of 23 assemblies all reaching  
133 reference quality (average QV>53, average N50=65.8Mbp) (**Extended Data Fig. 1b,c,**  
134 **Supplementary Table 2**). Final genome sizes ranged from ~713 Mbp (*S. etuberosum*) to ~2.5  
135 Gbp (*S. robustum*), with members of the *Lasiocarpa* subclade having four of the five largest  
136 genomes. An integrated gene prediction strategy for annotation based on RNA-seq and liftover,  
137 allowed us to identify 825,493 high-confidence gene models across the pan-genome (**Extended**  
138 **Data Fig. 1d, Supplementary Table 3, and Methods**). Of these, 495,429 (~60%) were shared  
139 across all samples, reflecting these species' relatively close evolutionary relationships.

140 An ortholog-based phylogenetic tree divided the 22 species into two major clades,  
141 consistent with previous studies<sup>23,24</sup>. Using existing nomenclature<sup>23</sup>, Grade I included the major  
142 crops tomato and potato while Clade II contained all prickly species, including the three cultivated  
143 eggplant species: *S. melongena* (Brinjal eggplant), *S. aethiopicum* (African eggplant), and *S.*  
144 *macrocarpon* (Gboma eggplant) (**Fig. 1c**). Whereas gene content was largely uniform across  
145 species, transposable element content and distribution varied widely (**Supplementary Table 4**).  
146 Consistent with other plant pan-genomes<sup>28,29</sup>, species-specific increases in repetitive content,  
147 driven primarily by a rapid expansion of retrotransposon families, correlated strongly with genome  
148 size expansion (**Fig. 1d**). The pan-genomic k-mer content – illustrating the genomic diversity  
149 within a species relative to the rest of the pan-genome – varied by clade, with 11 species containing  
150 more than 25% species-specific sequences (**Fig. 1d**). Finally, ortholog-based analysis revealed  
151 broad conservation of gene macrosynteny throughout the pan-genome, with the highest  
152 conservation on chromosomes 1, 2, 6, and 9 (**Fig. 1e**). This analysis also revealed large structural  
153 rearrangements across the genus and predominantly within sub-clades of clade II, including, for  
154 example, megabase-scale inversions and translocations involving chromosomes 3, 5, 10, and 12  
155 (**Fig. 1e**). These high-quality genomes provided a foundation for capturing genetic diversity across  
156 the *Solanum* from the clade to the species level, setting the stage for an analysis of paralog  
157 evolutionary dynamics and their impacts on genotype-to-phenotype relationships across this  
158 species-rich, ecologically and economically important plant genus.

159

## 160 **Pan-genome analysis reveals a complex landscape of gene duplications in *Solanum***

161 To develop a comprehensive view of gene evolutionary dynamics across *Solanum*, we  
162 reconstructed the genus-wide history of orthogroup expansion and contraction events across the

163 22 species, anchored on tomato (*S. lycopersicum*) (**Fig. 2a**). From the 44,962 total orthogroups  
164 identified across the *Solanum* pan-genome, we identified several of them were involved in  
165 expansion (26,284) or contraction (37,267) events, with the majority of the evolutionary events  
166 occurring at inner nodes involving orthogroup contractions. Functional enrichment analysis  
167 revealed that expanding and contracting orthogroups are predominantly linked to environmental  
168 response and secondary metabolism, with species- and clade-specific features (**Extended Data**  
169 **Fig. 2a, Supplementary Table 5, Supplementary Table 6**). We then characterized orthogroups  
170 based on their representation in the pan-genome, and classified these orthogroups as core (present  
171 in 100% of the species), near core (present in >70% of genomes), dispensable (present in 5-70%  
172 of species), and private (found in one species only) (**Fig. 2b**). Most orthogroups are core (60.6%)  
173 or near core (20.2%), while smaller proportions are dispensable (14.3%) or private (0.8%). Finally,  
174 75% of pairs of orthologous genes (designated paragroups) are dispensable or private, suggesting  
175 derived paralogs are more genetically flexible than orthologs (**Extended Data Fig. 2b**).

176 Across all orthogroups, gene duplications were widespread, with 70% (575,464 duplicates)  
177 of all genes having a paralog (**Fig. 2c**). We classified the duplications based on their genomic  
178 context as whole-genome (WGD) or single gene duplication, including tandem, proximal,  
179 transposed, or dispersed duplications<sup>30</sup> (**Fig. 2c**). Paralogs most frequently originate from WGDs  
180 from events many millions of years ago; however, single gene duplications, which typically are  
181 more recent and lineage-specific events, collectively dominate the duplication landscape in  
182 *Solanum* (**Extended Data Fig. 2c**). While most of the WGD-derived duplications belong to core  
183 orthogroups, single gene duplications show increased bias towards near core and dispensable  
184 orthogroups (**Fig. 2c**). Analysis of duplication types differentiated according to biological function  
185 using a GO enrichment analysis show that WGD-derived paralog pairs are most strongly  
186 associated with dosage-sensitive processes, such as DNA transcription and DNA replication, as  
187 well as hormone-mediated signal transduction and response (**Fig. 2d**), consistent with previous  
188 reports<sup>31,32</sup>. In contrast, and as already shown in many systems<sup>30,33</sup>, tandem and proximal  
189 duplications are most associated with defense and specialized metabolite biosynthesis, along with  
190 diverse functional roles related to environmental responses (**Fig. 2d**).

191 Paralogous genes functionally diverge through changes in both coding and *cis*-regulatory  
192 sequences<sup>34,35</sup>; however, it is unclear if the relative contributions of these changes are associated  
193 with specific duplication types. To test this, we first used our previously developed algorithm,



194 Conservatory, which simultaneously allows quantification of *cis*-regulatory conservation and  
195 improved calling of paralog pairs based on both protein and *cis*-regulatory conservation<sup>36</sup>  
196 **(Extended Data Fig. 2d and Methods)**. We then incorporated Ka/Ks ratios, as a measure of  
197 coding sequence selection, with both protein and *cis*-regulatory conservation to determine  
198 relationships in coding and regulatory sequence evolution across the duplication types. As  
199 expected, for all five types of duplications, protein similarity decreases with higher Ka/Ks values  
200 **(Fig. 2e, Extended Data Fig. 2e)**. However, two striking patterns of *cis*-regulatory conservation  
201 distinguish different duplication types: tandem and proximal duplicates maintain high *cis*-  
202 regulatory conservation across all levels of selection, whereas WGD, dispersed, and transposed  
203 duplicates show higher levels of *cis*-regulatory sequence similarity with increasing Ka/Ks. This  
204 observation suggests a greater degree of expression pattern conservation among non-locally  
205 duplicated paralogs undergoing functional diversification at the protein level.

206

### 207 **Multi-tissue transcriptomics uncovers the fate of retained paralogs**

208 Research in yeast and other systems suggests that duplicated genes can have negative  
209 fitness effects due to increased expression dosage, leading to stoichiometric imbalances in  
210 macromolecular complexes<sup>37,38</sup>. Consequently, early diversification of *cis*-regulatory sequences of  
211 paralogs may serve to restore ancestral single-copy gene dosage levels in a process called  
212 compensatory drift<sup>21,39</sup>. To explore constraints on total expression dosage from retained paralogs,  
213 we established two broad categories of paralog pairs as Dosage constrained, or Dosage  
214 unconstrained across species and on a per tissue basis **(Fig. 3a)**. We defined dosage constrained  
215 orthogroups as paralog pairs that exhibited similar total expression levels in a given tissue across  
216 species, whereas unconstrained orthogroups did not maintain the same summed expression  
217 **(Extended Data Fig. 3a)**. To assign paralog pairs to these categories, we generated a pan-*Solanum*  
218 gene expression resource comprising 271 samples from 22 species, 15 of which had data from two  
219 or more distinct tissues **(Extended Data Fig. 3b)**. Principal component analysis (PCA) on the  
220 TPM-normalized expression of 5,146 singleton genes showed that the vast majority of samples  
221 clustered by tissue type **(Fig. 3b)**. As in yeast<sup>40</sup>, our data show that paralog pairs typically evolved  
222 under total dosage constraint across tissues and species **(Fig. 3c)**. These pairs also exhibited much  
223 lower rates of non-synonymous mutations and were less likely to be tissue-specific than  
224 unconstrained pairs.

225 Dosage relationships between paralog pairs can be influenced by different evolutionary  
226 trajectories resulting in divergent expression patterns. Among retained paralog pairs within a given  
227 species we considered four groups of common patterns of expression relationships following gene  
228 duplication (**Fig. 3d, Extended Data Fig. 3c**): Group I, Dosage balanced: selection on total dosage  
229 remains high, and pairs retain similar expression profiles and levels across tissues; Group II,  
230 Paralog dominance: Substantial divergence in expression levels that are proportional across  
231 tissues; Group III, Specialization: Expression profiles no longer showing a purely global shift and  
232 instead exhibiting tissue-specific changes; Group IV, Divergence: Paralog pairs are fully diverged  
233 in both expression profile and level. Applying the definitions to our paralog gene expression  
234 dataset showed 58,130 (~8%) of the paralog pairs to a specific group, leaving over 92%  
235 undetermined as they do not yet exhibit strong trajectories (**Fig. 3e,f, Extended Data Fig. 3d**).

236 While these groups were defined by the expression profiles across tissues within a species,  
237 the data also allowed us to evaluate if the groups were associated with distinct genetic features.  
238 We compared protein sequence similarity between the groups, as well as gene family function,  
239 size, expression status, the number of tissues where expressed, and transcription levels (**Fig. 3g,**  
240 **Extended Data Fig. 3e**). We observed that pairs in Group I showed higher sequence similarity,  
241 smaller gene family size, broader expression across tissues, and higher transcription levels than  
242 groups undergoing paralog dominance, specialization and divergence (Groups II-IV) (**Fig. 3g**).  
243 Functional enrichment analysis showed that Groups I-II are enriched in dosage-sensitive processes  
244 such as transcription and translation, while Groups III-IV are enriched in defense response genes  
245 (**Extended Data Fig. 3e**). Moreover, consistent with their conserved expression patterns, paralog  
246 pairs in Groups I and Group II maintained greater *cis*-regulatory sequence conservation than those  
247 in Groups III and IV (**Fig. 3h, Extended Data Fig. 3f**). We further reasoned that the type of  
248 duplications from which gene pairs originated might impact their expression relationships. We  
249 found that the most conserved expression groups—paralog pairs in Groups I and II that also capture  
250 more ancient duplications—were more likely to have originated from WGDs, whereas gene pairs  
251 in Groups III and IV were enriched in small-scale duplications (SSDs) (**Fig. 3i**). Although all four  
252 of our defined Groups have the potential to complicate crop engineering, the 60% of pairs with  
253 correlated expression patterns likely pose the greatest challenge due to interdependent redundant,  
254 compensatory or partially sub-functionalized relationships, which could reflect a continuum of  
255 lineage- or specific-specific variations in these relationships.

256

## 257 **Genetic dissection of lineage-specific paralog diversification and compensatory relationships**

258 The *Solanum* pan-genome provided an opportunity to study the extent to which paralog  
259 diversifications have shaped key genes that influence genotype-phenotype relationships across the  
260 genus. Based on prior characterization and cloning of QTL and developmental genes affecting 16  
261 domestication and breeding traits, we compiled a set of 148 genes and associated paralogs (where  
262 relevant) from primarily the three model *Solanum* crops (eggplant, potato, tomato)  
263 (**Supplementary Table 7**). Our pan-genome revealed widespread variation in these genes between  
264 and within clades, with many cases of gene presence-absence variation (PAV), copy number  
265 variation (CNV), and gene truncation/pseudogenization across the pan-genome. Prominent among  
266 these were 17 orthogroups containing genes, harboring variants that contribute to the three major  
267 components of the crop domestication syndromes (flowering time & plant architecture;  
268 inflorescence architecture & flower number; and fruit size) (**Fig. 4a**). For example, in tomato and  
269 many other species, variation in the dosage-sensitive florigen-antiflorigen family members (*SP*,  
270 *SP5G*, *FTL1a*, *FTL1b*, *SP6D*, *SP6A*, *SFT*) enabled selection for accelerated flowering and short  
271 stature (determinate) plants, key traits that facilitated mechanical harvesting<sup>41–43</sup>. We identified  
272 numerous CNVs and loss-of-function mutations affecting paralogous genes in our pan-genome,  
273 suggesting these variants modulate flowering and growth habit across *Solanum*. In the genetics of  
274 inflorescence architecture, mutations in the MADS-box transcription factor-encoding gene *J2*  
275 allowed mechanical harvesting of tomato by eliminating the abscission zone of fruit stems<sup>44,45</sup>.  
276 However, co-occurring mutations in its ancestral paralog *EJ2* result in undesirable inflorescence  
277 branching<sup>46</sup>. We found one CNV and at least three ancestral losses of *J2* in our pan-genome, with  
278 most losses occurring in the Eastern Hemisphere Spiny eggplant clade (**Fig. 1c**). These species  
279 may therefore be sensitized to changes in inflorescence branching from natural or engineered *EJ2*  
280 mutations.

281 The increase of fruit size in tomato domestication was driven in large part by a promoter  
282 structural variant in the stem-cell signaling peptide gene, *CLAVATA3 (CLV3)*<sup>47</sup>. *CLE9*, a partially  
283 redundant ancestral paralog, falls into Group II (paralog dominance) and partially compensates for  
284 the effect of the *CLV3* domestication allele<sup>48,49</sup>. We previously showed *CLE9* was pseudogenized  
285 or completely lost in several Solanaceae species, which eliminated partial redundancy with  
286 *CLV3*<sup>48</sup>. Notably, except for tomato and *S. americanum*, all species in our pan-genome contain a

287 pseudogenized *CLE9* or lack it entirely. Meanwhile, a subset of the Eastern Hemisphere Spiny  
288 eggplant clade possess locally duplicated intact and pseudogenized copies of *CLV3* (**Fig. 4a, b**).  
289 Our chromosome-scale references revealed complex haplotypes involving these duplications, with  
290 species-specific transposable element invasions and disease resistance genes interspersed between  
291 the paralogs. For example, whereas *S. prinophyllum* carries two intact copies of *CLV3*, one intact  
292 and one to three pseudogenized copies exist in *S. aethiopicum* (African eggplant, 1 pseudogenized  
293 copy), its progenitor *S. anguivi* (1 pseudogenized copy), and *S. linnaeanum* (3 pseudogenized  
294 copies), with extreme variation in transposable element and disease resistance gene content and  
295 structure (**Fig. 4b, Extended Data Fig. 4a,b**). Comparing haplotypes and observing identical  
296 breakpoints in pseudogene structure across these species suggested at least two independent *CLV3*  
297 duplication events in the Eastern Hemisphere Spiny clade where one ancestral duplication was  
298 followed by pseudogenization in the last common ancestor of *S. insanum*, *S. linnaeanum*, *S.*  
299 *anguivi*, and *S. aethiopicum*, whereas a more recent *CLV3* duplication emerged in the lineage  
300 leading to *S. prinophyllum* (**Fig. 4b**). However, we cannot exclude the possibility of three  
301 independent duplications, as *S. violaceum* carries only one *CLV3* copy.

302 The independent duplication resulting in two intact copies of *CLV3* in *S. prinophyllum*  
303 suggests redundancy was re-established in this species (Group I), whereas in species where one  
304 *CLV3* paralog was pseudogenized, redundancy was again lost. We tested this by using  
305 CRISPR/Cas9 to inactivate *CLV3* in three spiny *Solanum* species: *S. cleistogamum* (desert raisin -  
306 *ScleCLV3* single copy), *S. aethiopicum* (African eggplant - one functional (*SaetCLV3a*) and one  
307 pseudogenized (*SaetCLV3b*)), and *S. prinophyllum* (intact copies of *SpriCLV3a* and *SpriCLV3b*)  
308 (**Fig. 4c, Extended Data Fig. 4c,d**). As expected, mutations in the one intact copy of *CLV3* in *S.*  
309 *cleistogamum* and *S. aethiopicum* resulted in extreme fasciation phenotypes, matching tomato *clv3*  
310 *cle9* double mutants (**Fig. 4c**). Similarly, knocking out both copies of *CLV3* in *S. prinophyllum*  
311 (*SpriCLV3a* and *SpriCLV3b*) replicated this severe phenotype.

312 *SpriCLV3a* and *SpriCLV3b* in *S. prinophyllum* are identical in their coding and *cis*-  
313 regulatory sequences, except for a single nucleotide variant in the 3' untranslated region (UTR) of  
314 the ancestral copy. Such high sequence identity suggested that eliminating one copy would be fully  
315 compensated for by the remaining functional copy, similar to the near complete compensation  
316 between *PgriCLV3* and *PgriCLE9* in the *Solanaceae* species *Physalis grisea* (groundcherry)<sup>48</sup>.  
317 Our previously generated transcriptomic data of meristems from *S. prinophyllum*<sup>50</sup> showed both

318 paralogs are expressed to similar levels (**Fig. 4d**) and supported this prediction. Surprisingly, we  
319 found that engineered mutations in either *SpriCLV3* paralog resulted in a subtle shift to more  
320 trilocular fruits compared to wild type (5% trilocular in WT compared to 30% trilocular in single  
321 mutants), suggesting one paralog cannot fully compensate for the other, most likely because of a  
322 gene expression dosage effect (**Fig. 4e,f, Supplementary Table 8**).

323 Taken together, these data suggest that, following the loss of the ancestral *CLE9* paralog,  
324 subsequent tandem duplication events in three spiny *Solanum* lineages would have reestablished  
325 *CLV3* compensation. However, this compensation was then lost again in at least one lineage due  
326 to pseudogenization of the derived *CLV3* duplicate. Finally, despite retention of both nearly  
327 identical copies of *CLV3* in *S. prinophyllum*, complete compensation was not fully maintained.  
328 Similar to *CLV3*, dynamic duplication histories and resulting paralog relationships affecting  
329 meristem proliferation and other gene families critical for domestication and trait improvement  
330 may reveal the species-specific contingencies that impact outcomes in genome engineering.

331

### 332 **African eggplant pan-genomics reveals widespread introgression and paralog diversification**

333 African eggplant (*S. aethiopicum*) is a major crop indigenous to sub-Saharan Africa and  
334 cultivated across the continent on hundreds of thousands of acres. Transported by enslaved  
335 Africans, it is also grown extensively in Brazil, but outside of these regions it remains largely  
336 unknown (**Fig. 5a**)<sup>51,52</sup>. Diverse cultivars are grown in Africa for their edible fruits or leaves, as  
337 well as for the ornamental appeal of specific fruit types<sup>53</sup>. These disparate uses are reflected in the  
338 species' broad intraspecific diversity in vegetative and fruit phenotypes, including fruit shape,  
339 color, and size (**Fig. 5b**). Breeding in African eggplant has primarily focussed on improving  
340 adaptation to abiotic stress conditions<sup>54,55</sup>, with less progress on improving yield or productivity.  
341 Re-engineering or mimicking the effects of known beneficial mutations from tomato and other  
342 *Solanum* model crops could advance these goals, but genomic and genetic resources are limited.

343 To address this, we first phenotyped eight representative accessions (**Supplementary**  
344 **Table 9**) from the Gilo (fruit production), Aculeatum (ornamental), and Shum (leaf production)  
345 cultivar groups in field conditions (**Fig. 5a**) along with one accession of *S. anguivi*. Based on the  
346 observed phenotypic variation, we extended our selection to 10 diverse accessions belonging to  
347 the three cultivar groups (**Supplementary Table 9**) and assembled a long-read based African  
348 eggplant pan-genome that included its wild progenitor *S. anguivi*. The reference African eggplant

349 accession (PI 424860) belongs to the Gilo group, and was used as the representative genotype in  
350 the wider *Solanum* pan-genome (**Fig. 1**). To assess genetic relationships, we computed an  
351 ortholog-based phylogenetic tree (**Fig. 5b**), which indicated two major clades, one comprising the  
352 three Gilo accessions and a second containing the five Aculeatum accessions. Interestingly, the  
353 two Shum accessions did not form a monophyletic group, suggesting that accessions cultivated for  
354 leaf production might have different genetic origins. Protein-coding genes were primarily clustered  
355 at chromosome ends throughout the African eggplant pan-genome, a pattern similar to other  
356 *Solanum* and flowering plant species (**Fig. 5c**). Transposable element distribution complemented  
357 this pattern, with more elements accumulating in the gene-poor pericentromeric regions.

358 Comparing the African eggplant genomes against the reference showed that, at the  
359 sequence level, most of the genome is highly conserved. Over 250,000 structural variants (SVs:  
360 defined as variants at least 50 bp in size) were found across all African eggplant samples, mainly  
361 towards chromosome ends (**Fig. 5d, Extended Data Fig. 5a**). Similar to our tomato pan-genome<sup>19</sup>,  
362 over 68% of SVs were located within 5 kbp upstream or downstream of genes, in addition to 7,234  
363 SVs overlapping exons and therefore likely to disrupt gene function (**Fig. 5e, Extended Fig. Data**  
364 **5b**). While average SV length was similar across accessions, their absolute number varied between  
365 groups, with Gilo possessing the fewest SVs, an expected pattern since the reference African  
366 eggplant belongs to the Gilo group. Notably, the SV distribution showed clade-specific SVs and  
367 SV clusters shared with the wild ancestor *S. anguivi*, suggesting a history of introgression (**Fig.**  
368 **5f**). Using a window-based Jaccard similarity analysis, we found multiple introgressions from *S.*  
369 *anguivi* in the Aculeatum accessions, most evident on chromosomes 3, 4, 11, and 12. Such  
370 widespread introgression suggests recent gene flow from the wild species in the course of African  
371 eggplant breeding, and likely explaining the origin of the Aculeatum ornamental types (**Fig. 5b, f,**  
372 **g**).

373 Similar to tomato, African eggplant cultivar groups exhibit extreme variation in fruit size,  
374 based in large part on variation in locule number (**Fig. 5b**). We reasoned that, beyond interspecific  
375 paralog dynamics observed throughout the pan-genome, recent diversification of key regulators of  
376 fruit locule number, such as *SaetCLV3*, might have favored intraspecific phenotypic diversity. The  
377 *SaetCLV3* locus, located on chromosome 10, is nested in dense SV clusters (**Fig. 5h**). Interestingly,  
378 one Aculeatum accession (804750136) has only a single intact copy of *SaetCLV3*, suggesting the  
379 ancestral pseudogenized copy was eliminated (**Fig. 5i, Extended Data Fig. 5c**). Microsynteny

380 analysis revealed broad rearrangements at the *CLV3* locus between African eggplant and *S.*  
381 *anguivi*, as well as intraspecific diversity (**Fig. 5j**). Notably, we detected two deletions within the  
382 *SaetCLV3* locus in two *S. aethiopicum* accessions (804750136 and PI 247828), including a ~300  
383 kbp deletion between the second exon of *SaetCLV3a* and the first exon of *SaetCLV3b* (**Fig. 5k**).  
384 Remarkably, the large deletion resulted in a fusion between the intact and pseudogenized  
385 *SaetCLV3* copies, resulting in a single functional copy, designated *SaetCLV3<sup>DEL</sup>* (**Fig. 5k**).  
386

### 387 **Paralog contingencies impact fruit locule number step changes in African eggplant**

388 We next sought to understand if *SaetCLV3* haplotype and paralog dynamics influenced  
389 locule number variation. Using our African eggplant genomes, we performed QTL-seq to map loci  
390 controlling locule number (**Supplementary Tables 10, 11, 12**). We generated F2 mapping  
391 populations between the high-locule count reference accession (PI 424860) belonging to the Gilo  
392 group and low- and high-locule count parents belonging to the Shum (804750187) and Aculeatum  
393 (804750136) groups, respectively (**Fig. 6a, Extended Data Fig. 6a**). In contrast to tomato, the  
394 major step change in locule number between the Gilo and Shum groups mapped to a QTL in a 3.9  
395 Mbp region on chromosome 2, which conspicuously did not include *CLV3* or any other known  
396 *CLV* pathway components (**Fig. 6b**). Instead, we identified a candidate gene encoding a serine  
397 carboxypeptidase (*SaetSCPL25-like*, named after its best BLAST hit in *Arabidopsis*<sup>56</sup>) harboring  
398 a 5 bp exonic frameshift deletion in the Gilo parent. Serine carboxypeptidases function in C-  
399 terminal peptide processing, and such control of CLE peptide processing has been demonstrated  
400 in *Arabidopsis*, where mutation of the Zn<sup>2+</sup> carboxypeptidase-encoding gene *SOLI* (*Suppressor of*  
401 *LLPI*) represses CLE-dependent root meristem size-related defects<sup>57</sup>. The mutation in  
402 *SaetSCPL25-like* in African eggplant was associated with the development of approximately two  
403 additional locules (**Fig. 6c**). We validated this association by mutating the orthologs of this gene  
404 in both tomato and *S. prinophyllum* using CRISPR/Cas9, which caused quantitatively similar  
405 locule number increases as the natural mutation in African eggplant (**Fig. 6d**).

406 We also identified two minor effect QTLs from the Aculeatum group, which we mapped  
407 to a 1.8 Mbp region on chromosome 5 and a 4.9 Mbp region on chromosome 10. The latter  
408 encompasses the *SaetCLV3<sup>DEL</sup>* haplotype harboring the reconstituted single copy functional  
409 *SaetCLV3* (**Fig. 5e** and **Fig. 6c**). We found that these two minor effect QTLs interact, with the  
410 homozygous *SaetCLV3<sup>DEL</sup>* genotype masking the increase in locule number conferred by the

411 chromosome 5 haplotype derived from the Aculeatum parent (**Extended Data Fig. 6b**). Though  
412 the specific gene and variant underlying the chromosome 5 QTL remains to be characterized, these  
413 results indicate that multiple interacting loci, two of which affect the CLV3 signaling pathway,  
414 gave rise to increases in locule number in African eggplants.

415 We then asked how these QTLs shaped the domestication history of African eggplant by  
416 examining the alleles present at the three identified loci within the phylogenetic context of our  
417 African eggplant pan-genome (**Fig. 6c**). The Gilo accessions contained the *SaetSCPL25-like*  
418 mutant allele, while all surveyed Aculeateum accessions and one of the Shum accessions harbored  
419 the chromosome 5 minor effect QTL's haplotype. Meanwhile, a single Aculeateum accession  
420 (804750136) contained all three identified alleles, including the minor effect *SaetCLV3<sup>DEL</sup>*  
421 structural variant (**Fig. 6c**). The SV at *SaetCLV3* probably occurred secondarily to variants at  
422 *SaetSCPL25-like* and the chromosome 5 QTL. *SaetCLV3<sup>DEL</sup>* causes a subtle reduction in locule  
423 number, and was perhaps selected to attenuate the locule count increases conferred by the  
424 combined synergistic effect of *SaetSCPL25-like* and the chromosome 5 QTL (**Extended Data Fig.**  
425 **6b**). This contrasts with tomato, where previous studies identified the *SaetCLV3* structural variant  
426 *SlycCLV3<sup>fas</sup>* as a widespread and major effect QTL variant yielding increased fruit locule number,  
427 modified by other minor effect QTLs, including the paralog *SICLE9<sup>49</sup>*. Thus, while QTLs affecting  
428 *CLV* signaling are shared drivers of increased locule number in both tomato and African eggplant,  
429 the specific genes, alleles, and interactions, as well as the magnitude and directionality of these  
430 individual and combined effects, are distinct (**Fig. 6e**). The recurrence of QTLs at *SaetCLV3* in  
431 two independent domestication histories underscores the major contribution of structural variation  
432 on paralog evolutionary dynamics as key contingencies shaping parallel trajectories of crop  
433 domestication and improvement.

434

## 435 DISCUSSION

436 Plant pan-genome resources are emerging at an incredible pace. A widespread assumption  
437 is that implementing genome editing technologies on these foundational resources will be the  
438 panacea to translating genotype-to-phenotype knowledge between related crops and also their wild  
439 relatives<sup>11,12</sup>. Decades of work by plant breeders demonstrates, however, that additive and epistatic  
440 effects from background genetic modifiers are a barrier to predicting desirable outcomes<sup>14–16,58</sup>.  
441 While sequencing high-quality plant references at scale, including potentially telomere-to-



442 telomere genomes<sup>59</sup>, combined with forward genetics, can readily uncover background variation,  
443 identifying orthologs and paralogs and tracing their evolutionary trajectories remains an unsolved  
444 challenge, particularly given the exceptionally complex history of ancient whole-genome  
445 duplications and more recent smaller-scale duplications across flowering plants. This is especially  
446 problematic in pan-genomes spanning broader taxonomic scales, where more extreme amounts of  
447 sequence variation are found.

448 We approached the challenge of resolving orthologs, paralogs, and their diversification  
449 histories using an integrated approach. We used existing tomato and eggplant annotations, multi-  
450 tissue RNA-seq annotations, and manual curation to expose and compare ancient paralogs and  
451 recent tandem duplications across our pan-genome. We mapped core and dispensable genes in the  
452 pan-genome, and among the tens of thousands of paralog pairs identified, expression analyses  
453 revealed a continuum of redundancy relationships, driven by drifting expression patterns,  
454 pseudogenization, or gene loss. Most dramatically, we showed that paralogs of the fruit size gene  
455 *CLV3* captured all three possible scenarios, reflected in independent tandem duplication events,  
456 extreme haplotype shuffling, and pseudogenization, accounting for variation in this domestication  
457 trait within and between species. Our approaches showcase how leveraging knowledge from major  
458 crops to indigenous crops and wild species can reveal previously unknown factors involved in trait  
459 variation, opening the door to reciprocal knowledge gain and new paths to improving all crop  
460 species.

461 Similarly complex paralog evolutionary histories undoubtedly affect other traits in  
462 nightshades, grasses, legumes, and beyond. Assembling widely and deeply sampled species and  
463 genotypes into super pan-genomes<sup>29,60</sup> offers watershed opportunities to both better understand  
464 origins and frequencies of genome fragility within and between species and mobilize advances in  
465 machine learning for *de novo* genetic and genomic predictions at scale. As more accurate machine  
466 learning models are developed, the micro-level analysis (e.g. read-level basecalling<sup>61</sup> or variant  
467 detection<sup>62</sup>) as well as higher level predictions of epigenomic and regulatory activity<sup>63</sup> have been  
468 and will be greatly improved and expedited. Efforts to predict gene expression changes from *cis*-  
469 regulatory variations are also maturing, although limitations in the modeling frameworks and their  
470 training regimes remain obstacles to achieving high predictive accuracy<sup>64</sup>. Advancing these efforts  
471 to predict trait variation from both coding and *cis*-regulatory variations will undoubtedly be even  
472 more challenging. Our work here shows that such models must explicitly account for paralogs and

473 their diversification dynamics over both short and long evolutionary times. Nevertheless, our  
474 ability to predict genotype-to-phenotype relationships, a holy grail for genetics and biology, will  
475 inevitably be enhanced by developing a foundation model trained on ever-increasing catalogs of  
476 molecular, cellular, and organismal data within and across species, to aid in both plant breeding  
477 and understanding natural diversity.

478 We also recognize that real-world implementation of pan-genomic and pan-genetic  
479 resources, tools, and technologies requires a deeper understanding of, and sensitivity to, the central  
480 role that indigenous knowledge and cultures have played in botany and agriculture<sup>10,65,66</sup>. Within  
481 this project, ethnobotanical knowledge from local breeders provided essential expertise in  
482 choosing the lineages, species, and cultivars to give our pan-genome immediate impact in  
483 agriculture. This includes the potential to rescue traits of agronomic interest that may have been  
484 lost during domestication, such as stress resistance and specialized metabolism<sup>67,68</sup>. This is most  
485 exemplified through the inclusion of African eggplant, one of the most economically and culturally  
486 important crops in tropical sub-Saharan Africa. Our integrated genomic, transformation, and  
487 genome editing pipeline complements the rich genetic and phenotypic diversity available in the  
488 African eggplant germplasm, offering new and more predictable avenues for breeding. For  
489 example, from dissecting the parallel, but distinct, genetic and epistatic paths towards increased  
490 locule number in tomato and African eggplant, we have more power to predictably increase locule  
491 number, fruit size and yield in this indigenous and regionally important crop.

492 We expect additional advancements will come from resolving paralog histories and  
493 relationships of flowering regulators, which have been central to agricultural revolutions<sup>69</sup>. It is  
494 important to highlight, however, that while industrialized breeding emphasizes yield, the needs of  
495 subsistence farmers can be different<sup>70</sup>. In the case of African eggplant, modifying flowering time  
496 and inflorescence architecture are arguably as important as increasing fruit size. In varieties grown  
497 for fruit production, earlier flowering and more branched genotypes would simultaneously dwarf  
498 plants and accelerate fruit production and total yield, whereas in varieties cultivated for leaf  
499 consumption, late flowering would prolong vegetative growth and vegetative yield<sup>71,72</sup>. We  
500 propose that the florigen-antiflorigen flowering hormone system and its MADS-box gene targets  
501 should be the primary targets to achieve these goals. In particular, our study revealed distinct  
502 diversifications in African eggplant of both florigen and antiflorigen paralogs from tomato, where  
503 there is already deep knowledge of these genes and their functional relationships<sup>69</sup>. Knowledge of

504 these paralogs, their allelic diversity, and epistatic relationships and contingencies will provide  
505 opportunities to accelerate breeding of these traits in African eggplant with natural alleles of these  
506 genes, which can now be characterized through pan-genome-enabled quantitative genetics, and  
507 will facilitate predictable outcomes from genome engineering. Looking forward, the most  
508 promising strategies for improving indigenous crops can only be realized through effective  
509 communication, understanding, and collaboration among local people, scientists, breeders, and  
510 growers.

511

### 512 **Online Content**

513 Any methods, additional references, Nature Portfolio reporting summaries, source data, extended  
514 data, supplementary information, acknowledgements, peer review information; details of author  
515 contributions and competing interests; and statements of data and code availability are available  
516 at <https://doi.org/xxxxxxxxx>

517 **References**

- 518 1. FAO. The State of Food Security and Nutrition in the World 2020.  
519 <https://openknowledge.fao.org/items/08c592f2-1962-4e1a-a541-695f9404b26d> (2020).
- 520 2. Renard, D. & Tilman, D. National food production stabilized by crop diversity. *Nature* **571**, 257–260  
521 (2019).
- 522 3. Ye, C.-Y. & Fan, L. Orphan Crops and their Wild Relatives in the Genomic Era. *Mol. Plant* **14**, 27–  
523 39 (2021).
- 524 4. Woldeyohannes, A. B. *et al.* Data-driven, participatory characterization of farmer varieties discloses  
525 teff breeding potential under current and future climates. *Elife* **11**, (2022).
- 526 5. Varshney, R. K. *et al.* Draft genome sequence of pigeonpea (*Cajanus cajan*), an orphan legume crop  
527 of resource-poor farmers. *Nat. Biotechnol.* **30**, 83–89 (2011).
- 528 6. Devos, K. M. *et al.* Genome analyses reveal population structure and a purple stigma color gene  
529 candidate in finger millet. *Nat. Commun.* **14**, 3694 (2023).
- 530 7. Moonlight, P. W. *et al.* Twenty years of big plant genera. *Proc. Biol. Sci.* **291**, 20240702 (2024).
- 531 8. Hilgenhof, R. *et al.* Morphological trait evolution in *Solanum* (Solanaceae): Evolutionary lability of  
532 key taxonomic characters. *Taxon* **72**, 811–847 (2023).
- 533 9. Shorinola, O. *et al.* Integrative and inclusive genomics to promote the use of underutilised crops.  
534 *Nat. Commun.* **15**, 320 (2024).
- 535 10. Dwyer, W., Ibe, C. N. & Rhee, S. Y. Renaming Indigenous crops and addressing colonial bias in  
536 scientific language. *Trends Plant Sci.* **27**, 1189–1192 (2022).
- 537 11. Fernie, A. R. & Yan, J. De Novo Domestication: An Alternative Route toward New Crops for the  
538 Future. *Mol. Plant* **12**, 615–631 (2019).
- 539 12. Zsögön, A. *et al.* De novo domestication of wild tomato using genome editing. *Nat. Biotechnol.*  
540 (2018) doi:10.1038/nbt.4272.
- 541 13. Gasparini, K., Figueiredo, Y. G., Araújo, W. L., Peres, L. E. & Zsögön, A. De novo domestication in

- 542 the Solanaceae: advances and challenges. *Curr. Opin. Biotechnol.* **89**, 103177 (2024).
- 543 14. Sackton, T. B. & Hartl, D. L. Genotypic Context and Epistasis in Individuals and Populations. *Cell*  
544 **166**, 279–287 (2016).
- 545 15. Liu, R. *et al.* Evaluating the Genetic Background Effect on Dissecting the Genetic Basis of Kernel  
546 Traits in Reciprocal Maize Introgression Lines. *Genes* **14**, (2023).
- 547 16. Lecomte, L. *et al.* Marker-assisted introgression of five QTLs controlling fruit quality traits into  
548 three tomato lines revealed interactions between QTLs and genetic backgrounds. *Theor. Appl. Genet.*  
549 **109**, 658–668 (2004).
- 550 17. Shen, L. *et al.* QTL editing confers opposing yield performance in different rice varieties. *J. Integr.*  
551 *Plant Biol.* **60**, 89–93 (2018).
- 552 18. Ruffley, M. *et al.* Selection constraints of plant adaptation can be relaxed by gene editing. *bioRxiv*  
553 2023.10.16.562583 (2024) doi:10.1101/2023.10.16.562583.
- 554 19. Alonge, M. *et al.* Major Impacts of Widespread Structural Variation on Gene Expression and Crop  
555 Improvement in Tomato. *Cell* **182**, 145–161.e23 (2020).
- 556 20. Soyk, S. *et al.* Duplication of a domestication locus neutralized a cryptic variant that caused a  
557 breeding barrier in tomato. *Nature Plants* **5**, 471–479 (2019).
- 558 21. Birchler, J. A. & Yang, H. The multiple fates of gene duplications: Deletion, hypofunctionalization,  
559 subfunctionalization, neofunctionalization, dosage balance constraints, and neutral variation. *Plant*  
560 *Cell* **34**, 2466–2474 (2022).
- 561 22. Gout, J.-F. *et al.* Dynamics of Gene Loss following Ancient Whole-Genome Duplication in the  
562 Cryptic Paramecium Complex. *Mol. Biol. Evol.* **40**, (2023).
- 563 23. Gagnon, E. *et al.* Phylogenomic discordance suggests polytomies along the backbone of the large  
564 genus *Solanum*. *Am. J. Bot.* **109**, 580–601 (2022).
- 565 24. Messeder, J. V. S. *et al.* A highly resolved nuclear phylogeny uncovers strong phylogenetic  
566 conservatism and correlated evolution of fruit color and size in *Solanum* L. *New Phytol.* **243**, 765–  
567 780 (2024).

- 568 25. Särkinen, T., Bohs, L., Olmstead, R. G. & Knapp, S. A phylogenetic framework for evolutionary  
569 study of the nightshades (Solanaceae): a dated 1000-tip tree. *BMC Evol. Biol.* **13**, 214 (2013).
- 570 26. Satterlee, J. W. *et al.* Convergent evolution of plant prickles by repeated gene co-option over deep  
571 time. *Science* **385**, eado1663 (2024).
- 572 27. Wu, Y. *et al.* Phylogenomic discovery of deleterious mutations facilitates hybrid potato breeding.  
573 *Cell* **186**, 2313–2328.e15 (2023).
- 574 28. Hufford, M. B. *et al.* De novo assembly, annotation, and comparative analysis of 26 diverse maize  
575 genomes. *Science* **373**, 655–662 (2021).
- 576 29. Bozan, I. *et al.* Pangenome analyses reveal impact of transposable elements and ploidy on the  
577 evolution of potato species. *Proc. Natl. Acad. Sci. U. S. A.* **120**, e2211117120 (2023).
- 578 30. Qiao, X. *et al.* Gene duplication and evolution in recurring polyploidization-diploidization cycles in  
579 plants. *Genome Biol.* **20**, 38 (2019).
- 580 31. Zhang, T. *et al.* Phylogenomic profiles of whole-genome duplications in Poaceae and landscape of  
581 differential duplicate retention and losses among major Poaceae lineages. *Nat. Commun.* **15**, 3305  
582 (2024).
- 583 32. Tang, H., Bowers, J. E., Wang, X. & Paterson, A. H. Angiosperm genome comparisons reveal early  
584 polyploidy in the monocot lineage. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 472–477 (2010).
- 585 33. Qiao, X. *et al.* Different Modes of Gene Duplication Show Divergent Evolutionary Patterns and  
586 Contribute Differently to the Expansion of Gene Families Involved in Important Fruit Traits in Pear  
587 (*Pyrus bretschneideri*). *Front. Plant Sci.* **9**, 161 (2018).
- 588 34. Baudouin-Gonzalez, L. *et al.* Diverse Cis-Regulatory Mechanisms Contribute to Expression  
589 Evolution of Tandem Gene Duplicates. *Mol. Biol. Evol.* **34**, 3132–3147 (2017).
- 590 35. Zhong, X., Lundberg, M. & Råberg, L. Divergence in Coding Sequence and Expression of Different  
591 Functional Categories of Immune Genes between Two Wild Rodent Species. *Genome Biol. Evol.* **13**,  
592 (2021).
- 593 36. Hendelman, A. *et al.* Conserved pleiotropy of an ancient plant homeobox gene uncovered by cis-

- 594 regulatory dissection. *Cell* **184**, 1724–1739.e16 (2021).
- 595 37. Veitia, R. A. & Potier, M. C. Gene dosage imbalances: action, reaction, and models. *Trends*  
596 *Biochem. Sci.* **40**, 309–317 (2015).
- 597 38. Diss, G. *et al.* Gene duplication can impart fragility, not robustness, in the yeast protein interaction  
598 network. *Science* **355**, 630–634 (2017).
- 599 39. Thompson, A., Zakon, H. H. & Kirkpatrick, M. Compensatory Drift and the Evolutionary Dynamics  
600 of Dosage-Sensitive Duplicate Genes. *Genetics* **202**, 765–774 (2016).
- 601 40. Gout, J.-F. & Lynch, M. Maintenance and Loss of Duplicated Genes by Dosage  
602 Subfunctionalization. *Mol. Biol. Evol.* **32**, 2141–2148 (2015).
- 603 41. Nakamichi, N. Adaptation to the local environment by modifications of the photoperiod response in  
604 crops. *Plant Cell Physiol.* **56**, 594–604 (2015).
- 605 42. Pnueli, L. *et al.* The SELF-PRUNING gene of tomato regulates vegetative to reproductive switching  
606 of sympodial meristems and is the ortholog of CEN and TFL1. *Development* **125**, 1979–1989  
607 (1998).
- 608 43. Soyk, S. *et al.* Variation in the flowering gene SELF PRUNING 5G promotes day-neutrality and  
609 early yield in tomato. *Nat. Genet.* **49**, 162–168 (2017).
- 610 44. Budiman, M. A. *et al.* Localization of jointless-2 gene in the centromeric region of tomato  
611 chromosome 12 based on high resolution genetic and physical mapping. *Theor. Appl. Genet.* **108**,  
612 190–196 (2004).
- 613 45. Rick, C. M. A new jointless gene from the Galapagos *L. pimpinellifolium*. *TGC Rep* **23**, (1956).
- 614 46. Soyk, S. *et al.* Bypassing Negative Epistasis on Yield in Tomato Imposed by a Domestication Gene.  
615 *Cell* **169**, 1142–1155.e12 (2017).
- 616 47. Rodriguez-Leal, D. *et al.* Evolution of buffering in a genetic circuit controlling plant stem cell  
617 proliferation. *Nat. Genet.* **51**, 786–792 (2019).
- 618 48. Kwon, C.-T. *et al.* Dynamic evolution of small signalling peptide compensation in plant stem cell  
619 control. *Nature Plants* **8**, 346–355 (2022).

- 620 49. Aguirre, L., Hendelman, A., Hutton, S. F., McCandlish, D. M. & Lippman, Z. B. Idiosyncratic and  
621 dose-dependent epistasis drives variation in tomato fruit size. *Science* **382**, 315–320 (2023).
- 622 50. Lemmon, Z. H. *et al.* The evolution of inflorescence diversity in the nightshades and heterochrony  
623 during meristem maturation. *Genome Res.* **26**, 1676–1686 (2016).
- 624 51. Lester, R. N. & Niakan, L. Origin and domestication of the scarlet eggplant, *Solanum aethiopicum*,  
625 from *S. anguivi* in Africa. in *International Symposium on the Biology and Systematics of the*  
626 *Solanaceae* (Columbia University Press, 1986).
- 627 52. Vorontsova, M. & Knapp, S. *A Revision of the Spiny Solanums, Solanum Subgenus Leptostemonum*  
628 *(Solanaceae), in Africa and Madagascar.* (THE AMERICAN SOCIETY OF PLANT  
629 TAXONOMISTS, 2016).
- 630 53. Yang, R.-Y. & Ojiewo, C. African Nightshades and African Eggplants: Taxonomy, Crop  
631 Management, Utilization, and Phytonutrients. in *African Natural Plant Products Volume II:*  
632 *Discoveries and Challenges in Chemistry, Health, and Nutrition* vol. 1127 137–165 (American  
633 Chemical Society, 2013).
- 634 54. Nakanwagi, M. J., Sseremba, G., Kabod, N. P., Masanza, M. & Kizito, E. B. Identification of growth  
635 stage-specific watering thresholds for drought screening in *Solanum aethiopicum* Shum. *Sci. Rep.*  
636 **10**, 862 (2020).
- 637 55. Sseremba, G., Tongoona, P., Eleblu, J., Danquah, E. Y. & Kizito, E. B. Heritability of drought  
638 resistance in *Solanum aethiopicum* Shum group and combining ability of genotypes for drought  
639 tolerance and recovery. *Sci. Hortic.* **240**, 213–220 (2018).
- 640 56. Fraser, C. M., Rider, L. W. & Chapple, C. An expression and bioinformatics analysis of the  
641 *Arabidopsis* serine carboxypeptidase-like gene family. *Plant Physiol.* **138**, 1136–1148 (2005).
- 642 57. Casamitjana-Martínez, E. *et al.* Root-specific CLE19 overexpression and the *sol1/2* suppressors  
643 implicate a CLV-like pathway in the control of *Arabidopsis* root meristem maintenance. *Curr. Biol.*  
644 **13**, 1435–1441 (2003).
- 645 58. Soyk, S., Benoit, M. & Lippman, Z. B. New Horizons for Dissecting Epistasis in Crop Quantitative



- 646 Trait Variation. *Annu. Rev. Genet.* **54**, 287–307 (2020).
- 647 59. Koren, S. *et al.* Gapless assembly of complete human and plant chromosomes using only nanopore  
648 sequencing. *bioRxiv* (2024) doi:10.1101/2024.03.15.585294.
- 649 60. Shi, T. *et al.* The super-pangenome of *Populus* unveils genomic facets for its adaptation and  
650 diversification in widespread forest trees. *Mol. Plant* **17**, 725–746 (2024).
- 651 61. Baid, G. *et al.* DeepConsensus improves the accuracy of sequences with a gap-aware sequence  
652 transformer. *Nat. Biotechnol.* **41**, 232–238 (2023).
- 653 62. Poplin, R. *et al.* A universal SNP and small-indel variant caller using deep neural networks. *Nat.*  
654 *Biotechnol.* **36**, 983–987 (2018).
- 655 63. Sokolova, K., Chen, K. M., Hao, Y., Zhou, J. & Troyanskaya, O. G. Deep Learning Sequence  
656 Models for Transcriptional Regulation. *Annu. Rev. Genomics Hum. Genet.* (2024)  
657 doi:10.1146/annurev-genom-021623-024727.
- 658 64. Huang, C. *et al.* Personal transcriptome variation is poorly explained by current genomic deep  
659 learning models. *Nat. Genet.* **55**, 2056–2059 (2023).
- 660 65. Kimmerer, R. W. & Artelle, K. A. Time to support Indigenous science. *Science* **383**, 243 (2024).
- 661 66. Bartlett, M. E., Moyers, B. T., Man, J., Subramaniam, B. & Makunga, N. P. The Power and Perils of  
662 De Novo Domestication Using Genome Editing. *Annu. Rev. Plant Biol.* **74**, 727–750 (2023).
- 663 67. Singh, J. & van der Knaap, E. Unintended Consequences of Plant Domestication. *Plant Cell Physiol.*  
664 **63**, 1573–1583 (2022).
- 665 68. Alam, O. & Purugganan, M. D. Domestication and the evolution of crops: variable syndromes,  
666 complex genetic architectures, and ecological entanglements. *Plant Cell* **36**, 1227–1241 (2024).
- 667 69. Eshed, Y. & Lippman, Z. B. Revolutions in agriculture chart a course for targeted breeding of old  
668 and new crops. *Science* **366**, (2019).
- 669 70. Nakyewa, B. *et al.* Farmer preferred traits and genotype choices in *Solanum aethiopicum* L., Shum  
670 group. *J. Ethnobiol. Ethnomed.* **17**, 27 (2021).
- 671 71. Plazas, M. *et al.* Conventional and phenomics characterization provides insight into the diversity and

- 672 relationships of hypervariable scarlet (*Solanum aethiopicum* L.) and gboma (*S. macrocarpon* L.)  
673 eggplant complexes. *Front. Plant Sci.* **5**, 318 (2014).
- 674 72. Park, S. J. *et al.* Optimization of crop productivity in tomato using induced mutations in the florigen  
675 pathway. *Nat. Genet.* **46**, 1337–1342 (2014).
- 676 73. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data.  
677 *Bioinformatics* **30**, 2114–2120 (2014).
- 678 74. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- 679 75. Kokot, M., Dlugosz, M. & Deorowicz, S. KMC 3: counting and manipulating k-mer statistics.  
680 *Bioinformatics* **33**, 2759–2761 (2017).
- 681 76. Ranallo-Benavidez, T. R., Jaron, K. S. & Schatz, M. C. GenomeScope 2.0 and Smudgeplot for  
682 reference-free profiling of polyploid genomes. *Nat. Commun.* **11**, 1432 (2020).
- 683 77. Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. Haplotype-resolved de novo assembly  
684 using phased assembly graphs with hifiasm. *Nat. Methods* **18**, 170–175 (2021).
- 685 78. Alonge, M. *et al.* Automated assembly scaffolding using RagTag elevates a new tomato system for  
686 high-throughput genome editing. *Genome Biol.* **23**, 258 (2022).
- 687 79. Rhie, A., Walenz, B. P., Koren, S. & Phillippy, A. M. Merqury: reference-free quality, completeness,  
688 and phasing assessment for genome assemblies. *Genome Biol.* **21**, 245 (2020).
- 689 80. Kovaka, S. *et al.* Transcriptome assembly from long-read RNA-seq alignments with StringTie2.  
690 *Genome Biol.* **20**, 278 (2019).
- 691 81. Mapleson, D., Venturini, L., Kaithakottil, G. & Swarbreck, D. Efficient and accurate detection of  
692 splice junctions from RNA-seq with Portcullis. *Gigascience* **7**, (2018).
- 693 82. Hosmani, P. S. *et al.* An improved de novo assembly and annotation of the tomato reference genome  
694 using single-molecule sequencing, Hi-C proximity ligation and optical maps. *bioRxiv* 767764 (2019)  
695 doi:10.1101/767764.
- 696 83. Li, D. *et al.* A high-quality genome assembly of the eggplant provides insights into the molecular  
697 basis of disease resistance and chlorogenic acid synthesis. *Mol. Ecol. Resour.* **21**, 1274–1286 (2021).

- 698 84. Shumate, A. & Salzberg, S. L. Liftoff: accurate mapping of gene annotations. *Bioinformatics* **37**,  
699 1639–1643 (2021).
- 700 85. Wu, T. D., Reeder, J., Lawrence, M., Becker, G. & Brauer, M. J. GMAP and GSNAP for Genomic  
701 Sequence Alignment: Enhancements to Speed, Accuracy, and Functionality. *Methods Mol. Biol.*  
702 **1418**, 283–334 (2016).
- 703 86. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100  
704 (2018).
- 705 87. Potato Genome Sequencing Consortium *et al.* Genome sequence and analysis of the tuber crop  
706 potato. *Nature* **475**, 189–195 (2011).
- 707 88. Venturini, L., Caim, S., Kaithakottil, G. G., Mapleson, D. L. & Swarbreck, D. Leveraging multiple  
708 transcriptome assembly methods for improved gene structure annotation. *Gigascience* **7**, (2018).
- 709 89. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative genomics.  
710 *Genome Biol.* **20**, 238 (2019).
- 711 90. Li, H. Protein-to-genome alignment with miniprot. *Bioinformatics* **39**, (2023).
- 712 91. Lovell, J. T. *et al.* GENESPACE tracks regions of interest and gene copy number variation across  
713 multiple genomes. *Elife* **11**, (2022).
- 714 92. Hart, A. J. *et al.* EnTAP: Bringing faster and smarter functional annotation to non-model eukaryotic  
715 transcriptomes. *Mol. Ecol. Resour.* **20**, 591–604 (2020).
- 716 93. Van Bel, M. *et al.* PLAZA 5.0: extending the scope and power of comparative and functional  
717 genomics in plants. *Nucleic Acids Res.* **50**, D1468–D1474 (2022).
- 718 94. Apweiler, R. *et al.* UniProt: the Universal Protein knowledgebase. *Nucleic Acids Res.* **32**, D115–9  
719 (2004).
- 720 95. Jones, P. *et al.* InterProScan 5: genome-scale protein function classification. *Bioinformatics* **30**,  
721 1236–1240 (2014).
- 722 96. Van Bel, M. *et al.* TRAPID: an efficient online tool for the functional and comparative analysis of de  
723 novoRNA-Seq transcriptomes. *Genome Biol.* **14**, 1–10 (2013).

- 724 97. Zhang, R.-G. *et al.* TESorter: an accurate and fast method to classify LTR-retrotransposons in plant  
725 genomes. *Hortic Res* **9**, (2022).
- 726 98. Manni, M., Berkeley, M. R., Seppey, M. & Zdobnov, E. M. BUSCO: Assessing Genomic Data  
727 Quality and Beyond. *Curr Protoc* **1**, e323 (2021).
- 728 99. Jiang, N., Gao, D., Xiao, H. & van der Knaap, E. Genome organization of the tomato sun locus and  
729 characterization of the unusual retrotransposon Rider. *Plant J.* **60**, 181–193 (2009).
- 730 100. Ou, S. *et al.* Benchmarking transposable element annotation methods for creation of a streamlined,  
731 comprehensive pipeline. *Genome Biol.* **20**, 275 (2019).
- 732 101. Barchi, L. *et al.* Improved genome assembly and pan-genome provide key insights into eggplant  
733 domestication and breeding. *Plant J.* **107**, 579–596 (2021).
- 734 102. Ou, S. *et al.* Differences in activity and stability drive transposable element variation in tropical and  
735 temperate maize. *bioRxiv* 2022.10.09.511471 (2022) doi:10.1101/2022.10.09.511471.
- 736 103. Ou, S., Chen, J. & Jiang, N. Assessing genome assembly quality using the LTR Assembly Index  
737 (LAI). *Nucleic Acids Res.* **46**, e126 (2018).
- 738 104. Van Eck, J., Keen, P. & Tjahjadi, M. Agrobacterium tumefaciens-Mediated Transformation of  
739 Tomato. in *Transgenic Plants: Methods and Protocols* (eds. Kumar, S., Barone, P. & Smith, M.)  
740 225–234 (Springer New York, New York, NY, 2019). doi:10.1007/978-1-4939-8778-8\_16.
- 741 105. Katoh, K., Misawa, K., Kuma, K.-I. & Miyata, T. MAFFT: a novel method for rapid multiple  
742 sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–3066 (2002).
- 743 106. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the  
744 Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
- 745 107. Zhang, C., Rabiee, M., Sayyari, E. & Mirarab, S. ASTRAL-III: polynomial time species tree  
746 reconstruction from partially resolved gene trees. *BMC Bioinformatics* **19**, 153 (2018).
- 747 108. Junier, T. & Zdobnov, E. M. The Newick utilities: high-throughput phylogenetic tree processing in  
748 the UNIX shell. *Bioinformatics* **26**, 1669–1670 (2010).
- 749 109. Sayyari, E. & Mirarab, S. Fast Coalescent-Based Computation of Local Branch Support from

- 750 Quartet Frequencies. *Mol. Biol. Evol.* **33**, 1654–1668 (2016).
- 751 110. Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T.-Y. Ggtree: An r package for visualization and  
752 annotation of phylogenetic trees with their covariates and other associated data. *Methods Ecol. Evol.*  
753 **8**, 28–36 (2017).
- 754 111. Wang, L.-G. *et al.* Treeio: An R Package for Phylogenetic Tree Input and Output with Richly  
755 Annotated and Associated Data. *Mol. Biol. Evol.* **37**, 599–603 (2020).
- 756 112. Mendes, F. K., Vanderpool, D., Fulton, B. & Hahn, M. W. CAFE 5 models variation in evolutionary  
757 rates among gene families. *Bioinformatics* **36**, 5516–5518 (2021).
- 758 113. Klopfenstein, D. V. *et al.* GOATOOLS: A Python library for Gene Ontology analyses. *Sci. Rep.* **8**,  
759 10872 (2018).
- 760 114. Wang, D., Zhang, Y., Zhang, Z., Zhu, J. & Yu, J. KaKs\_Calculator 2.0: a toolkit incorporating  
761 gamma-series methods and sliding window strategies. *Genomics Proteomics Bioinformatics* **8**, 77–  
762 80 (2010).
- 763 115. Yanai, I. *et al.* Genome-wide midrange transcription profiles reveal expression level relationships in  
764 human tissue specification. *Bioinformatics* **21**, 650–659 (2005).
- 765 116. Takagi, H. *et al.* QTL-seq: rapid mapping of quantitative trait loci in rice by whole genome  
766 resequencing of DNA from two bulked populations. *Plant J.* **74**, 174–183 (2013).
- 767 117. Doyle, J. J. & Doyle, J. L. A rapid DNA isolation procedure for small quantities of fresh leaf tissue.  
768 *Phytochemical bulletin* (1987).
- 769 118. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search  
770 tool. *J. Mol. Biol.* **215**, 403–410 (1990).
- 771 119. Harris, R. S. Improved pairwise alignment of genomic DNA. (The Pennsylvania State University.,  
772 2007).
- 773 120. Charif, D. & Lobry, J. R. SeqinR 1.0-2: A Contributed Package to the R Project for Statistical  
774 Computing Devoted to Biological Sequences Retrieval and Analysis. in *Structural Approaches to*  
775 *Sequence Evolution: Molecules, Networks, Populations* (eds. Bastolla, U., Porto, M., Roman, H. E.

776 & Vendruscolo, M.) 207–232 (Springer Berlin Heidelberg, Berlin, Heidelberg, 2007).

777 doi:10.1007/978-3-540-35306-5\_10.

## 778 **Methods**

### 779 **Tissue collection and high molecular weight DNA extraction**

780 For extraction of high molecular weight DNA, young leaves were collected from 21-day-  
781 old light-grown seedlings. Prior to tissue collection, seedlings were etiolated in complete darkness  
782 for 48 h. Flash-frozen plant tissue was ground using a mortar and pestle and extracted in four  
783 volumes of ice-cold extraction buffer 1 (0.4 M sucrose, 10 mM Tris-HCl pH 8, 10 mM MgCl<sub>2</sub>,  
784 and 5 mM 2-mercaptoethanol). Extracts were briefly vortexed, incubated on ice for 15 min, and  
785 filtered twice through a single layer of Miracloth (Millipore Sigma). Filtrates were centrifuged at  
786 4000 rpm for 20 min at 4°C, and pellets were gently resuspended in 1 ml of extraction buffer 2  
787 (0.25 M sucrose, 10 mM Tris-HCl pH 8, 10 mM MgCl<sub>2</sub>, 1% Triton X-100, and 5 mM 2-  
788 mercaptoethanol). Crude nuclear pellets were collected by centrifugation at 12,000g for 10 min at  
789 4°C and washed by resuspension in 1 ml of extraction buffer 2 followed by centrifugation at  
790 12,000g for 10 min at 4°C. Nuclear pellets were resuspended in 500 µl of extraction buffer 3 (1.7  
791 M sucrose, 10 mM Tris-HCl pH 8, 0.15% Triton X-100, 2 mM MgCl<sub>2</sub>, and 5 mM 2-  
792 mercaptoethanol), layered over 500 µl extraction buffer 3, and centrifuged for 30 min at 16,000g  
793 at 4°C. The nuclei were resuspended in 2.5 ml of nuclei lysis buffer (0.2 M Tris pH 7.5, 2 M NaCl,  
794 50 mM EDTA, and 55 mM CTAB) and 1 ml of 5% Sarkosyl solution and incubated at 60°C for  
795 30 min.

796 To extract DNA, nuclear extracts were gently mixed with 8.5 ml of chloroform/isoamyl  
797 alcohol solution (24:1) and slowly rotated for 15 min. After centrifugation at 4000 rpm for 20 min,  
798 3 ml of aqueous phase was transferred to new tubes and mixed with 300 µl of 3 M NaOAc and  
799 6.6 ml of ice-cold ethanol. Precipitated DNA strands were transferred to new 1.5 ml tubes and  
800 washed twice with ice-cold 80% ethanol. Dried DNA strands were dissolved in 100 µl of elution  
801 buffer (10 mM Tris-HCl, pH 8.5) overnight at 4°C. Quality, quantity, and molecular size of DNA  
802 samples were assessed using Nanodrop (ThermoFisher), Qubit (ThermoFisher), and pulsed-field  
803 gel electrophoresis (CHEF Mapper XA System, Biorad) according to the manufacturer's  
804 instructions.

805

## 806 **Tissue collection, RNA extraction and quantification**

807 All tissues were collected in 3-4 biological replicates from different greenhouse-grown  
808 plants at approximately 09:00-10:00 AM and flash frozen in liquid nitrogen in 1.5 mL microfuge  
809 tubes containing a 5/32 inch (~3.97 mm) 440 stainless steel ball bearing (BC Precision, TN, USA).  
810 Tubes containing tissue were placed in a -80°C stainless steel tube rack and ground using a  
811 SPEX™ SamplePrep 2010 Geno/Grinder™ (Cole-Parmer, NJ, USA) for 1 min at 1440 rpm. For  
812 shoot apices, total RNA was extracted using TRIzol (Invitrogen, MA, USA) according to the  
813 manufacturer's instructions for ground tissue. For all other tissues (cotyledons, hypocotyls, leaves,  
814 flower buds, and flowers), total RNA was extracted using Quick-RNA MicroPrep Kit (Zymo  
815 Research). RNA was treated with DNase I (Zymo Research, CA, USA) according to the  
816 manufacturer's instructions. Purity and concentration of the resulting total RNA was assessed  
817 using a NanoDrop One spectrophotometer (Fisher Scientific, MA, USA). Libraries for RNA-  
818 sequencing were prepared by KAPA mRNA HyperPrep Kit (Roche, Basel, Switzerland). Paired-  
819 end 100-base sequencing was conducted on the NextSeq 2000 P3 sequencing platform (Illumina,  
820 CA, USA). Reads were trimmed using trimmomatic v0.39<sup>73</sup> and then mapped to their respective  
821 genome using STAR v2.7.5c<sup>74</sup> and expression computed in transcripts per million (TPM).

822

## 823 **Genome assembly**

824 Reference quality genome assemblies for each of the 22 species (and two reference quality  
825 genomes for *S. muricatum*) (**Supplementary Table 2** for accession information) were generated  
826 using a combination of long-read sequencing (Pacific Biosciences, CA, USA) for contigging and  
827 optical mapping (Bionano Genomics, CA, USA) for scaffolding. Between 1-4 PacBio Sequel IIe  
828 flow cells (Pacific Biosciences, CA, USA) were used for the sequencing of each sample (average  
829 read N50 = 11,221 bp, average coverage = 53X, average read QV = 83.28). Prior to assembly, we  
830 counted k-mers from raw reads with KMC3<sup>75</sup> (version 3.2.1) and estimated genome size,  
831 sequencing coverage, and heterozygosity with GenomeScope2.0<sup>76</sup>. For 5 samples  
832 (**Supplementary Table 2** for details), low quality reads were filtered out with a custom script  
833 ([github.com/pan-sol/pan-sol-pipelines](https://github.com/pan-sol/pan-sol-pipelines)). Sequencing reads from each sample were assembled with  
834 hifiasm<sup>77</sup> exact parameters and software version varied between samples based on the level of



835 estimated heterozygosity and are reported in **Supplementary Table 2**. Post assembly, the draft  
836 contigs were screened for possible microbial contamination as previously described<sup>19</sup>.

### 837 **Genome assembly scaffolding**

838 Optical mapping (Bionano Genomics, CA, USA) was performed for 17 samples to  
839 facilitate scaffolding. Scaffolding with optical maps was performed using the Bionano solve  
840 Hybrid Scaffold pipeline with the recommended default parameters  
841 (<https://bionano.com/software-downloads/>). Hybrid scaffold N50s ranged from 33,254,022 bp to  
842 219,385,699 bp (see **Supplementary Table 2** for more detail including Bionano molecules per  
843 sample). High-throughput chromosome conformation capture (Hi-C) from Arima Genomics, CA,  
844 USA was performed for 8 samples to finalize scaffolding. With Hi-C, reads were integrated with  
845 the Juicer (v0.7.17-r1198-dirty) pipeline. Next, misjoins and chromosomal boundaries were  
846 manually curated in the Juicebox (v1.11.08) application. Chromosomes were named based on  
847 sequence homology, determined with RagTag<sup>78</sup> scaffold (v2.1.0, default parameters), with the  
848 phylogenetically-closest finished genome (see **Supplementary Table 2** for details), 12 of these  
849 samples (including nine *S. aethiopicum* samples) were scaffolded with Ragtag. Finally, small  
850 contigs (< 50,000 bp) with > 95% of the sequence mapping to a named chromosome were  
851 removed. Additionally, small contigs (< 100,000 bp) with > 80% of the sequence mapping to a  
852 named chromosome that contained one or more duplicated BUSCO genes, but no single BUSCO  
853 genes, were also removed using a python script. Using merquy<sup>79</sup> with the HiFi data, the final  
854 consensus quality of the assemblies was estimated as QV=51.1333 on average and a completeness  
855 of 99.2741% on average.

### 856 **Gene Annotation**

857 The gene annotation pipeline (**Extended Data Fig. 1d**) involved several crucial steps.  
858 Initially, the quality of raw RNASeq reads underwent assessment using FastQC v0.11.9  
859 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). Subsequently, reference-based  
860 transcripts were generated employing STAR v2.7.5c<sup>74</sup> and Stringtie2 v2.1.2<sup>80</sup> workflows. To  
861 refine the data, invalid splice junctions from the STAR aligner were filtered out utilizing Portcullis  
862 v1.2.0<sup>81</sup>. Orthologs with coverage above 50% and 75% identity were lifted from Heinz v4.0<sup>82</sup> and  
863 Eggplant v4.1<sup>83</sup> via Liftoff v1.6.3<sup>84</sup> using parameters --copies,--exclude\_partial and employing  
864 both Gmap version 2020-10-14<sup>85</sup> and Minimap2 v2.17-r941<sup>86</sup> aligners. In addition, protein

865 evidence from several published Solanaceae genomes<sup>82,83,87</sup>, and the UniProt/SwissProt database  
866 were utilized to support gene annotation. Structural gene annotations were generated through the  
867 Mikado v2.0rc2<sup>88</sup> framework, leveraging evidence from the Daijin pipeline. Additionally,  
868 microsynteny and orthology to Heinz v4.0 and Eggplant v4.0 were assessed using Microsynteny  
869 and Orthofinder v2.5.2<sup>89</sup>. Correction of gene models with inframe stop codons utilized Miniprot2<sup>90</sup>  
870 protein alignments from Heinz v4.0 and Eggplant v4.1. Furthermore, gene models lacking start or  
871 stop codons were adjusted by placing them within 300 base pairs of the nearest codon location  
872 using a custom python script ([github.com/pan-sol/pan-sol-pipelines](https://github.com/pan-sol/pan-sol-pipelines)) (**Supplementary Table 3**).  
873 Overall gene synteny was visualized using GENESPACE (v1.3.1)<sup>91</sup>.

874 For functional annotation, ENTAP v0.10.8<sup>92</sup> integrated data from diverse databases such  
875 as PLAZA dicots (5.0)<sup>93</sup>, Uniprot/Swissprot<sup>94</sup>, TREMBL, RefSeq, Solanaceae proteins, and  
876 InterProScan<sup>595</sup> with Pfam, TIGRFAM, Gene Ontology, and TRAPID<sup>96</sup> annotations. Finally, the  
877 annotated data underwent a series of filtering steps, excluding proteins shorter than 20 amino acids,  
878 those exceeding 20 times the length of functional orthologs and transposable element genes, which  
879 were removed using the TESorter<sup>97</sup> pipeline.

880 We assessed the completeness of the gene models by assessing single-copy orthologs  
881 through BUSCO<sup>98</sup> in protein mode, comparing them against the *solanales\_odb10 database*.  
882 Additionally, we examined the presence or absence of a curated set of 180 candidate genes known  
883 to be crucial in QTL studies.

#### 884 **Transposable element annotation**

885 The *S. lycopersicum* chloroplast and mitochondrion sequences were collected from NCBI  
886 reference sequences NC\_007898.3 and NC\_035963.1, respectively. Non-transposable element  
887 repeat sequences including 18S rDNA (OK073663.1), 5S rDNA (X55697.1), 5.8S rDNA  
888 (X52265.1), 25S rDNA (OK073662.1), DNA spacer (AY366528.1), centromeric repeat  
889 (JA176199.1), and telomere sequences (TTTAGGG) were collected from NCBI and further  
890 curated. Transposable element sequences curated in the SUN locus study<sup>99</sup> as well as several other  
891 transposable element sequences from NCBI were also collected. These sequences were combined  
892 as the curated set of tomato repeats.

893 *De novo* transposable element annotation was first performed on each genome using EDTA  
894 v2.1.5<sup>100</sup>, with coding sequences from the ITAG4.0 Eggplant V4 annotation<sup>101</sup> provided (--cds) to  
895 purge gene coding sequences in the transposable element annotation and parameters of --anno 1 -

896 -sensitive 1 for sensitive detection and annotation of repeat sequences. Curated tomato repeats  
897 were supplied to EDTA (--curatedlib) for the *de novo* annotation. Transposable element  
898 annotations of individual genomes were together processed by panEDTA<sup>102</sup> for the creation of  
899 consistent pan-genome transposable element annotation. Summary of whole-genome repeat  
900 annotations were derived from .sum files generated by panEDTA (**Supplementary Table 4**).

901 Evaluation of repeat assembly quality was performed using LAI b3.2<sup>103</sup> with inputs  
902 generated by EDTA and parameters -t 48 -unlock. LAI of *S. aethiopicum* genomes were  
903 standardized based on the HiFi-based reference assembly, with parameters -iden 95.71 -totLTR  
904 49.22 -genome\_size 1102623763 -t 48 -unlock.

905

## 906 **Generation of CRISPR-Cas9-induced mutants**

907 CRISPR guide RNAs to target *CLV3* and *SCPL25* across *Solanum* species were designed  
908 using Geneious. The Golden Gate cloning approach as described in (29) was used to create  
909 multiplexed gRNA constructs. Plant regeneration and *Agrobacterium tumefaciens*-mediated  
910 transformation of *S. prinophyllum* were performed according to our previously published  
911 protocol<sup>104</sup>. For *S. cleistogamum* plant regeneration, the medium was supplemented with 0.5 mg/L  
912 zeatin instead of 2 mg/L and for the selection medium, 75 mg/L kanamycin was used instead of  
913 200 mg/L. For *S. aethiopicum*, the protocol was the same as for *S. cleistogamum*, except the fourth  
914 transfer of transformed plantlets is done onto medium supplemented with 50 mg/L kanamycin.  
915 Seed germination time in culture can vary between species and batches of harvested seeds.  
916 Typically, *S. prinophyllum* germination took 8-10 days, *S. cleistogamum* germinated in 6-8 days,  
917 and *S. aethiopicum* in 7-10 days.

918

## 919 **Distribution maps and species status**

920 Species were categorized into wild, domesticated, locally-important consumed, or  
921 ornamental based on taxonomic literature and expert opinion<sup>8</sup> (PBI *Solanum* Project (2024).  
922 *Solanaceae* Source. <http://www.solanaceaesource.org/>). Native ranges were derived from the same  
923 taxonomic literature and approximate centroids of the ranges were used for the mapping.

924

## 925 **Phylogenomic analyses**

926 *Jaltomata sinuosa* was used an outgroup for the *Solanum* pan-genome tree, whereas the  
927 closely related *S. anguivi*, *S. insanum*, and *S. melongena* were used as an outgroup for the *Solanum*  
928 *aethiopicum* dataset. Orthofinder<sup>89</sup> was used to identify single copy orthologs across all species.  
929 This resulted in 7,825 loci for the *Solanum* pan-genome dataset, and 19,769 loci for the *S.*  
930 *aethiopicum* dataset. To reduce computing time, we randomly subsampled 5,000 loci for the *S.*  
931 *aethiopicum* dataset. To reduce the effect of missing data and long branch attraction, sequences  
932 shorter than 25% of the average length for each loci were eliminated, following Gagnon *et al.*  
933 (2022)<sup>23</sup>. MAFFT<sup>105</sup> was used to align each locus individually. Only loci that had all species in the  
934 alignment were kept. trimAl was also used to remove columns that had more than 75% gaps. IQ-  
935 TREE2<sup>106</sup> was used to generate individual ML trees for each locus. The resulting phylogenies were  
936 used for coalescent analyses with ASTRAL-III version 5.7.3<sup>107</sup>, where tree nodes with <30% BS  
937 values were collapsed using Newick Utilities version 1.5.0<sup>108</sup>. Branch support was assessed using  
938 localPP support<sup>109</sup>, where PP values >0.95 were considered strong, 0.75 to 0.94 weak to moderate,  
939 and ≤0.74 as unsupported. Trees were visualized with R using the packages ggtree<sup>110</sup> and treeio<sup>111</sup>.

940

## 941 **Gene expansion contraction analysis**

942 To analyze gene expansions and contractions, we processed the ultrametric species tree  
943 and gene family counts from OrthoFinder using CAFE5<sup>112</sup>. CAFE5 was run with the gamma model  
944 and parameter 'k=3' to identify changes in gene family size along the species tree while accounting  
945 for rate variation among gene families.

946

## 947 **GO enrichment analysis**

948 Gene Ontology (GO) enrichment analysis was performed using the GOATOOLS  
949 package<sup>113</sup> to investigate the functional implications of genes associated with various duplication  
950 types including whole-genome (WGD), tandem (TD), proximal (PD), transposed (TSD) and  
951 dispersed (DSD) duplications. Genes were classified into these different duplication categories by  
952 DupGenefinder<sup>30</sup>. Additionally, we conducted GO enrichment on gene expansions  
953 (**Supplementary Table 5**) and contractions (**Supplementary Table 6**) identified across all

954 lineages as reported by CAFE5, to explore functional trends related to these gene copy number  
955 changes across the pangenome.

956

### 957 **Synteny analysis**

958 Genomic neighborhood around *CLV3* for selected species was manually inspected to detect  
959 and annotate intact and pseudogenized *CLV3* copies using pairwise sequence comparison with  
960 Exonerate ([www.ebi.ac.uk/about/vertebrate-genomics/software/exonerate](http://www.ebi.ac.uk/about/vertebrate-genomics/software/exonerate)). Synteny plots were  
961 generated from a reciprocal BLASTP table obtained running Clinker (v0.0.29,  
962 [github.com/gamcil/clinker](https://github.com/gamcil/clinker)). Pseudomolecule visualization was generated *via* a custom script  
963 ([github.com/pan-sol/pan-sol-pipelines](https://github.com/pan-sol/pan-sol-pipelines)). Transposable elements and resistance genes annotations  
964 were overlaid as needed using custom scripts ([github.com/pan-sol/pan-sol-pipelines](https://github.com/pan-sol/pan-sol-pipelines)).

965

### 966 **Gene expression analysis**

967 Reads from each tissue sample were aligned to the corresponding species-specific genome  
968 using STAR v2.7.2b<sup>74</sup>, and only samples with more than 50% uniquely mapped reads were  
969 retained for subsequent analysis. For each species with two or more biological replicates per tissue,  
970 we calculated the Spearman correlation between tissue replicates, and removed samples with low  
971 correlation (0.75 or below). This yielded gene expression estimates for 271 samples across 22  
972 species, with 15 species having expression data in two or more tissues. Expression data was TPM-  
973 normalized and genes with zero expression across all samples were excluded from further analysis.  
974 Principal component analysis was performed on the tissue-specific expression profiles of 5,146  
975 singleton genes shared across all 22 species to reveal the global relationships among samples.

976

977 *Is the total dosage of duplicate gene pairs conserved across Solanum?*

978 Survival of a gene after duplication depends on the competition between preservation to  
979 maintain partial or total dosage and mutational degradation rendering one copy with reduced or no  
980 function. Consequently, functional fates of duplicate genes are often characterized by the extent  
981 of selective pressures on total dosage. To assess the relative importance of dosage balance (copies  
982 evolving under strong purifying selection to maintain total dosage) and neutral drift (no selection  
983 on total dosage) in maintaining duplicate genes, we compared the total expression of paralog pairs

984 within each tissue for each pair of species. Note that the prickle tissue from *S. prinophyllum* is not  
985 included in this analysis since it is absent in the other 21 species.

986 In each tissue, gene expression was averaged over the biological replicates for each species.  
987 For each pair of species with expression data in a shared tissue, orthogroups with exactly two  
988 copies in each species with non-zero average expression in the tissue were retained for further  
989 analysis. For each tissue and species pair, we calculated the summed expression of paralog pairs  
990 in each retained orthogroup, and observed that the total “orthogroup-level” expression was highly  
991 correlated across species suggesting a prominent role of dosage balance in shaping the expression  
992 evolution of paralogs. We computed the ratio of the orthogroup-level expression between the  
993 species pair and transformed them into z-scores. For each orthogroup in a species expressed in the  
994 tissue of interest, we averaged the *p*-values from all pairwise species comparisons, adjusted the  
995 average *p*-values using Benjamini-Hochberg correction, and classified orthogroups with adjusted  
996 average *p*-value < 0.05 as dosage-unconstrained orthogroups. All other orthogroups in the species  
997 and tissue were assumed to be evolving under constraint on total dosage.

998 All other orthogroups were assumed to evolve under selective constraint on total dosage.  
999 Note that the high z-score threshold provides a conservative estimate of the number of paralog  
1000 pairs evolving under drift. Sequence evolution rates for paralog pairs (Ka/Ks) were calculated  
1001 using KaKs\_Calculator 2.0<sup>114</sup>.

1002

### 1003 *Different modes of paralog functional evolution*

1004 For each of the 15 species with expression in two or more tissues, the expression data was  
1005 first subset to genes with more-than-median expression in at least one sample. Coexpression  
1006 network for each species was constructed by calculating the Pearson correlation between all pairs  
1007 of genes, ranking the correlation coefficients for each gene (with NAs assigned the median rank),  
1008 and then standardizing the network by the maximum ranked correlation coefficient. Coexpression  
1009 for each pair of paralogs in each orthogroup was obtained from this rank-standardized network.  
1010 For each paralog pair with non-zero expression in two or more samples, we also computed the  
1011 fold-change of expression across samples and used the absolute values of mean and standard  
1012 deviation (SD) of log<sub>2</sub>-transformed fold-change across samples to summarize the degree of  
1013 expression divergence between the two copies.

1014 We classified the paralog pairs within each species into different retention categories based  
1015 on their variation in expression levels and correlated expression across samples. We selected these  
1016 two axes of variation since they intuitively represent average expression difference (fold-change)  
1017 and specific pattern of difference (coexpression) between gene pairs. We classified paralog pairs  
1018 into four broad groups as follows:

- 1019 I. Dosage-balanced: coexpression  $> 0.9$ , mean  $\log_2$  fold-change  $< 1$ , SD of  $\log_2$  fold-change  
1020  $< 1$
- 1021 II. Paralog dominance: coexpression  $> 0.9$ , mean  $\log_2$  fold-change  $\geq 1$ , SD of  $\log_2$  fold-  
1022 change  $< 1$
- 1023 III. Specialized: coexpression  $> 0.9$ , mean  $\log_2$  fold-change  $\geq 1$ , SD of  $\log_2$  fold-change  $\geq 1$
- 1024 IV. Diverged: coexpression  $< 0.5$ , mean  $\log_2$  fold-change  $\geq 1$ , SD of  $\log_2$  fold-change  $\geq 1$

1025

1026 Paralogs originating from whole genome (WGD), tandem and proximal duplications were  
1027 obtained using the DupGen\_finder pipeline<sup>30</sup>. WGD pairs with  $K_s$  ranging from 0.2 to 2.5, and  
1028 tandem and proximal duplicates with  $K_s$  ranging from 0.05 to 2.5 were used to generate the stacked  
1029 bar plots corresponding to whole genome and small-scale duplications, respectively, in **Fig. 3i**.

1030 Gene family size for each classified paralog pair within a species corresponds to the  
1031 number of genes in its orthogroup. The expression breadth of a gene corresponds to the number of  
1032 tissues (among apices, cotyledon, hypocotyl, inflorescence, leaves) where the gene has an average  
1033 expression greater than 3 TPM. Number of shared tissues expressing a paralog pair is computed  
1034 by intersecting the expression breadths of both copies, and ranges from 0 to 5. A gene was  
1035 considered non-functional if it was annotated as a pseudogene or had an average expression below  
1036 3 TPM. Tissue-specific genes for each tissue were identified as genes with the highest expression  
1037 in the tissue of interest, tissue-specificity score<sup>115</sup> greater than 0.7 and with expression greater than  
1038 5 TPM in the relevant tissue.

1039

#### 1040 **Mapping of loci controlling *S. aethiopicum* locule number**

1041 The high-locule count parent and reference accession PI 424860, and low- and higher-  
1042 locule count parents 804750187 and 804750136, respectively, were selected as founding parents  
1043 to map QTLs and their causative variants affecting fruit locule number. Resulting F1 progeny were  
1044 selfed to generate F2 mapping populations, which were sown in the greenhouse and then

1045 transplanted to a field site at Lloyd Harbor, New York, USA, during the summer of 2022.  
1046 Approximately 10 fruits were collected from each F2 individual and the number of locules exposed  
1047 by slicing the fruit transversely and counting. In the 804750187 x PI 424860 and 804750136 x PI  
1048 424860 derived F2 populations, 144 and 135 individuals were phenotyped, respectively. For each  
1049 population, DNA from 30 random individuals at the low and high ends of the phenotypic  
1050 distribution for locule number were pooled for bulk-segregant QTL-Seq analysis. The DNA from  
1051 8 individuals of the common parental accession PI 424860 were also pooled to capture parental  
1052 polymorphisms.

1053 DNA from 30 randomly selected low- and high-locule count individuals was extracted  
1054 from young leaf tissue using a DNeasy Plant Pro Kit (Qiagen, Hilden, Germany) according to the  
1055 manufacturer's instructions for high-polysaccharide content plant tissue. Tissue used for extraction  
1056 was ground using a SPEX™ SamplePrep 2010 Geno/Grinder™ (Cole-Parmer, NJ, USA) for 2  
1057 min at 1440 rpm. Sample DNA (1 µL assay volume) concentrations were assayed using Qubit 1X  
1058 dsDNA HS buffer (ThermoFisher, MA, USA) on a Qubit 4 fluorometer (ThermoFisher, MA,  
1059 USA) according to the manufacturer's instructions. Separate pools were made for the parents, the  
1060 bulked high-locule count F2 individuals, and the bulked low-locule count F2 individuals, with an  
1061 equivalent mass of DNA pooled from each individual to yield a final pooled mass of 3 µg in each  
1062 bulk. DNA pools were purified using 1.8X volume of AMPure XP beads (Beckman Coulter, CA,  
1063 USA) and the DNA concentration and purity assayed by Qubit and a NanoDrop One  
1064 spectrophotometer (Fisher Scientific, MA, USA), respectively.

1065 Paired-end sequencing libraries for QTL-Seq analysis were prepared with >1 µg of DNA  
1066 using a KAPA HyperPrep PCR-free kit (Roche, Basel, Switzerland) according to the  
1067 manufacturer's instructions. Indexed libraries were pooled for sequencing on a NextSeq 2000 P3  
1068 chip (Illumina, CA, USA). Mapping was performed using the end-to-end pipeline implemented in  
1069 the QTL-Seq software package<sup>116</sup> (v2.2.4, [github.com/YuSugihara/QTL-seq](https://github.com/YuSugihara/QTL-seq)) with reads aligned  
1070 against the *S. aethiopicum* (Saet3, PI 424860) genome assembly.

1071 To determine the effects of the two identified QTL on locule number in the 804750136 x  
1072 PI 424860 derived populations, co-segregation analysis was performed on the full F2 populations  
1073 by genotyping *SaetCLV3* and the minor effect locus on chromosome 5. For *SaetCLV3*, a cleaved  
1074 amplified polymorphic sequence (CAPS) assay was used to genotype a variant in the promoter  
1075 region of *SaetCLV3* linked to the identified *CLV3* SV haplotypes. A 1258 bp region surrounding



1076 an *AseI* restriction fragment length polymorphism (RFLP) in the *SaetCLV3* promoter was  
1077 amplified using KOD One™ PCR Master Mix (Toyobo, Osaka, Japan) on template DNA extracted  
1078 by the cetyltrimethylammonium bromide (CTAB) method<sup>117</sup> (see **Supplementary Table 13** for  
1079 primers 5431 & 4681). To 5 µL of the resulting PCR product, a 10 µL reaction containing 0.2 µL  
1080 *AseI* (New England BioLabs, MA, USA) and 1 µL CutSmart™ r3.1 Buffer (New England  
1081 BioLabs, MA, USA) was incubated for 2 hours at 37 °C. The reactions were then loaded onto a  
1082 1% agarose gel and electrophoresed in an Owl™ D3-14 electrophoresis box (Thermo Scientific,  
1083 MA, USA) containing 1X TBE buffer for 30 min at 180 V delivered from an Owl™ EC 300 XL  
1084 power supply (Thermo Scientific, MA, USA). The electrophoresis results were visualized under  
1085 UV light using a Bio-Rad ChemiDoc™ XRS+ (Bio-Rad, CA, USA) imaging platform and  
1086 ImageLab™ (Bio-Rad, CA, USA) software. Resulting banding patterns were then used to assign  
1087 genotypes. For the chromosome 5 QTL, primers (see **Supplementary Table 13** for primers 5883  
1088 & 5884) were used to amplify a 425 bp region harboring a 1 bp deletion occurring near the summit  
1089 of the QTL peak using KOD One™ PCR Master Mix. The resulting PCR products were purified  
1090 using Ampure 1.8X beads and served as template for Sanger sequencing (Azenta Genewiz, NJ,  
1091 USA). The sequencing results were then used to assign genotype calls at chromosome 5.

1092

### 1093 **Conservatory analysis**

1094 The Conservatory algorithm (V2.0)<sup>36</sup> was employed to identify conserved noncoding  
1095 sequences (CNSs) within the *Solanaceae* family (**Extended Data Figure 2d**)  
1096 (<https://conservatorycns.com/dist/pages/conservatory/about.php>). A total of 26 genomes,  
1097 including 23 *Solanum* genomes, two tomato genomes (Heinz and M82) and one groundcherry  
1098 (*Physalis grisea*), were used as references to enable the identification of CNSs irrespective of  
1099 structural variations among references. Protein similarity was scored using Bitscore<sup>118</sup>, while *cis*-  
1100 regulatory similarity was assessed using LastZ<sup>119</sup> score. Homologous gene pairs were required to  
1101 share at least one CNS. For orthogroup calling, all orthologous genes shared at least one CNS with  
1102 the reference gene. Gene pairs with a conservation score exceeding 90% of the highest score were  
1103 classified as paralogs (**Extended Data Figure 2b**). A total of 844,525 paralogs were identified  
1104 across the *Solanum* pan-genome. Sequence evolution pressure rates (Ka/Ks) for paralog pairs were  
1105 calculated using the R seqinR package (v4.2-36)<sup>120</sup>. Gene duplication events were classified using  
1106 DupGenefinder<sup>30</sup>, identifying whole-genome (WGD) and transposed (TSD) duplications for gene

1107 pairs recognized by both Conservatory and DupGenefinder tools. Tandem (TD) and proximal (PD)  
1108 duplications were defined based on gene positioning: adjacent genes were considered TD, and  
1109 genes up to 10 genes apart were defined as PD. All other duplicated gene pairs were categorized  
1110 as dispersed (DSD) duplications (**Extended Data Figure 2c**). Of the identified paralogs, 23,730  
1111 were associated with expression groups and were used to compare relationships between sequence  
1112 evolution pressure rates and protein and *cis*-regulatory divergence across different expression  
1113 groups. Homologs, orthogroups, and paragroups were identified, and relationships between protein  
1114 and *cis*-regulatory elements were visualized using custom scripts, which are available on GitHub  
1115 ([github.com/pan-sol/pan-sol-pipelines](https://github.com/pan-sol/pan-sol-pipelines)).

1116

### 1117 **Statistical analysis**

1118 All statistical tests were performed in R. For the quantitative analysis of fruit locule  
1119 numbers in Figures 4f, 6c, 6d, and Extended Data Figure 6b, n represents the "number of fruits  
1120 quantified." Pairwise comparisons were conducted using Dunnett's T3 test for multiple  
1121 comparisons with unequal variances, with default parameters (see **Supplementary Tables 14-17**).

1122

### 1123 **Reporting summary**

1124 Further information on research design is available in the Nature Portfolio Reporting Summary  
1125 linked to this article.

1126

### 1127 **Data availability**

1128 All data are available within this Article and its Supplementary Information. Raw sequencing data  
1129 are available in the SRA under BioProject PRJNA1073673. Genome, expression, and phenotypic  
1130 data are available at Solpangenomics website ([www.solpangenomics.com](http://www.solpangenomics.com)). Paralog expression  
1131 analysis scripts are available at [github.com/gillislabs/pansol\\_expression\\_analysis](https://github.com/gillislabs/pansol_expression_analysis). Other analysis  
1132 scripts are available within [github.com/pan-sol/pan-sol-pipelines](https://github.com/pan-sol/pan-sol-pipelines).

1133

1134 **Acknowledgements**

1135 We thank members of the Lippman laboratory and critical friends M. Bartlett, Y. Eshed, and I.  
1136 Efroni for discussions and feedback. We thanks B. Seman from the Lippman lab for technical  
1137 support, and T. Mulligan, K. Schlecht, and S. Qiao for assistance with plant care. We thank E.  
1138 Cruickshank and T. Jenike for Pepino dulce fruit images. We thank S. Muller, R. Wappel, S.  
1139 Mavruk-Eskipehliyan, and E. Ghiban from the CSHL Genome Center for sequencing support. MB  
1140 is grateful for financial support from the Plant Health and Environment department of the French  
1141 National Institute for Agriculture, Food, and Environment (INRAE). JWS is supported by an NSF  
1142 Postdoctoral Fellowship in Biology (IOS-2305651). MJP is funded by the William Randolph  
1143 Hearst Scholarship from the Cold Spring Harbor School of Biological Sciences. SO is supported  
1144 by the OSU STEM Education Faculty Startup Award and the Global Gateways Initiative Grant.  
1145 Funded by the Convenio 566 of 2014 between Universidad Nacional de Colombia and Minciencias  
1146 (GPS and FR). UCU Research Funds (EBK). WRM is the Davis Family Professor of Human  
1147 Genetics. Work on *Solanum* taxonomy, morphology and phylogenetics was funded by NSF  
1148 Planetary Biodiversity Initiative grant "PBI Solanum: a worldwide initiative" (DEB-0316614 to  
1149 SK), Sibbald Trust (RH), Fonds de recherche du Québec - Nature et Technologies postdoctoral  
1150 fellowship (EG), and National Geographic Society Northern Europe Award GEFNE49-12 (TS).  
1151 National Institutes of Health grant R01MH113005 (JG). National Science Foundation Plant  
1152 Genome Research Program grant IOS-2216612 (AF, JG, JVE, MCS, ZBL). The Howard Hughes  
1153 Medical Institute (ZBL). Finally, we thank the indigenous peoples of Australia on whose ancestral  
1154 lands *S. cleistogamum* and *S. prinophyllum* grow.

1155

1156 **Author Contributions**

1157 Conceptualization: MCS, ZBL

1158 Data curation: MB, KMJ, JWS, SR, AH, MJP, HS1, MA, XW, SO

1159 Formal analysis: MB, KMJ, JWS, SR, IG, AH, MJP, HS1, MA, XW, SO, EG, TS, AF, JG, JVE,  
1160 MCS, ZBL

1161 Funding acquisition: MB, JWS, MJP, SO, GPS, FR, EBK, EG, SK, TS, AF, JG, JVE, MCS,  
1162 ZBL

1163 Investigation: MB, KMJ, JWS, SR, IG, AH, MJP, HS1, HS2, MA, XW, SO, JG, JVE, MCS,  
1164 ZBL

1165 Methodology: MB, KMJ, JWS, SR, AH, MJP, HS1, GMR, MA, XW, SO, YG, KS, EG, SK, TS,  
1166 JG, JVE, MCS, ZBL

1167 Project administration: MB, KMJ, AF, JG, JVE, MCS, ZBL

1168 Software: KMJ, SR, AH, MA, SO, MCS

1169 Resources: MB, KMJ, JWS, HS2, BF, MA, XW, RS, JH, HG, YG, KS, GPSR, AO, FR, SG,  
1170 WRM, EG, SK, TS, AF, MCS

1171 Supervision: JG, JVE, MCS, ZBL

1172 Validation: MB, KMJ, JWS, SR, AH, SO, EBK, EG, SK, TS, AF, JG, JVE, MCS, ZBL

1173 Visualization: MB, KMJ, JWS, SR, AH, GMR, EBK, JG, JVE, MCS, ZBL

1174 Writing – original draft: MB, KMJ, MCS, ZBL

1175 Writing – review & editing: MB, KMJ, JWS, SR, IG, AH, MJP, HS1, HS2, SO, EBK, EG, SK,  
1176 TS, AF, JG, JVE, MCS, ZBL

1177

#### 1178 **Competing Interests**

1179 WRM is a founder and shareholder in Orion Genomics, a plant genomics company. Z.B.L. is a  
1180 consultant for and a member of the Scientific Strategy Board of Inari Agriculture.

1181

#### 1182 **Additional information**

1183 Supplementary information The online version contains supplementary material available at  
1184 <https://doi.org/xxxxxx>

1185

1186 Correspondence and requests for materials should be addressed to Jesse Gillis, Joyce Van Eck,  
1187 Michael C. Schatz, or Zachary B. Lippman.

1188 **FIGURE LEGENDS**

1189

1190 **Figure 1: *Solanum* pan-genome captures the phenotypic, ecologic, agricultural, and genomic**  
1191 **diversity of this crop-rich genus. (a)** Approximate centroid of the native range for the 22 selected  
1192 *Solanum* species, grouped by type of agricultural use: wild (W), locally-important consumed (C),  
1193 ornamental (O), and domesticated (D). **(b)** Phenotypic diversity of shoots and fruits from a subset  
1194 of *Solanum* species in the pan-genome. Scale bars: 5 cm (shoots) and 1 cm (fruits). **(c)** Orthogroup-  
1195 based phylogeny of the *Solanum* pan-genome recapitulates the major clades, Grade I and Clade II.  
1196 Branch lengths reflect coalescent units. **(d)** Genomic features of each species of the *Solanum* pan-  
1197 genome. Genome size (Gbp) and representation of non-repetitive (light grey) and repetitive (dark  
1198 grey) sequences (left). Percentage of pan-k-mers shared across the pan-genome in each reference  
1199 (middle). Contribution of the different transposable element families in the total repeat landscape  
1200 of each genome (right). **(e)** GENESPACE plot showing gene macrosynteny across the pan-genome  
1201 relative to tomato. Scale bar: 9000 genes.

1202

1203 **Figure 2: Pan-genomic analysis of orthogroup conservation and diversity of gene**  
1204 **duplications. (a)** Orthogroups expansions and contractions across the pan-genome. The  
1205 orthogroup-based phylogeny is adapted from **Fig. 1c**. The estimated expansion (blue) and  
1206 contraction (orange) rates of orthogroups are shown at each node. **(b)** Cumulative curves showing  
1207 detection of the four orthogroup conservation groups as a function of the number of species  
1208 available in the pan-genome. **(c)** Schematic of the potential mechanisms underlying different gene  
1209 duplication categories (left). Stacked bar chart showing the number of genes derived from the  
1210 different types of duplication sorted by orthogroup conservation groups (right). WGD: whole-  
1211 genome duplication; TD: tandem duplication; PD: proximal duplication; TRD: transposed  
1212 duplication; DSD: dispersed duplication; SC: single copy. **(d)** Functional enrichment of gene  
1213 duplication types detected across the pan-genome. The top five enriched GO terms per duplication  
1214 type are shown. Circle size represents gene ratio. **(e)** Divergence of protein and *cis*-regulatory  
1215 sequences across increasing evolutionary pressure, as measured by Ka/Ks values, for the indicated  
1216 types of gene duplications. BLASTP (protein sequence conservation) and LASTZ (*cis*-regulatory  
1217 sequence conservation from the Conservatory algorithm) normalized alignment scores were used  
1218 to plot the predicted mean and 95% confidence interval.

1219

1220 **Figure 3: Widespread paralogous diversification across *Solanum* revealed by multi-tissue**  
1221 **gene expression analysis. (a)** Schematic of dosage-constrained and dosage-unconstrained  
1222 orthogroups reflecting different degrees of selection on the total dosage of paralog pairs across  
1223 species. **(b)** PCA of the normalized expression matrix from 5,146 singleton genes shared across  
1224 all 22 species. The expression matrix consists of the summed expression of paralog pairs. Tissue  
1225 samples are colored by tissue identity. **(c)** Bar plots showing that paralog pairs under constrained  
1226 total dosage across species are less tissue-specific (left) than unconstrained paralogs (right). **(d)**  
1227 Schematic of four categories of functional expression groups of retained paralogs: Group I: Dosage  
1228 balance; Group 2: Paralog dominance; Group III: Specialization; Group IV: Divergence. **(e)**  
1229 Scatter plots showing the distribution of paralog pairs according to their co-expression level and  
1230 mean log<sub>2</sub> fold-change (top) or standard deviation (S.D.) log<sub>2</sub> fold-change (bottom) in expression.  
1231 The four derived paralog expression groups are shown. **(f)** Representatives of paralog pairs  
1232 capturing the different patterns of expression delimited across the pan-genome. **(g)** Genes included  
1233 in the four paralog expression groups display contrasting protein sequence similarity (top left),  
1234 gene family size (top right), number of shared expression domains (tissues) (bottom left), or  
1235 propensity to undergo gene loss for orthogroups in different dosage quartiles (bottom right). **(h)**  
1236 Effect of *cis*-regulatory sequence conservation on the different expression groups in relation to  
1237 increased selection on protein sequence. For each expression group the predicted mean and 95%  
1238 confidence interval of the normalized LastZ score is shown. **(i)** Stacked bar plots showing the  
1239 proportion of each paralog expression group attributed to paralog pairs derived from either whole-  
1240 genome duplication (WGD) or small-scale duplication (SSD).

1241

1242 **Figure 4: Functional dissection of lineage-specific paralog diversification through pan-**  
1243 **genetics reveals modified compensatory relationships in a major fruit size regulator. (a)** Pan-  
1244 genome-wide gene presence-absence and copy number variation in 17 orthogroups containing  
1245 genes known to regulate three major domestication and improvement traits in tomato. Stars  
1246 indicate gene truncation or pseudogenization. **(b)** Haplotype diversification at the *CLV3* locus  
1247 across the eggplant clade. Presence-absence of *CLV3* paralogs is shown. Lineage-specific *CLV3*  
1248 duplications are marked with asterisks. Full circles denote functional *CLV3* copies and half circles  
1249 denote truncated/pseudogenized copies. Grey lines illustrate conservation, while blue lines

1250 represent loss of synteny. **(c)** CRISPR/Cas9 genome editing of *CLV3* orthologs in three species of  
1251 the eggplant clade. Engineered loss-of-function mutations in *S. cleistogamum* (*ScleCLV3*, top), *S.*  
1252 *aethiopicum* (*SaetCLV3a/b*, middle), and *S. prinophyllum* (*SpriCLV3a/b*, bottom) resulted in  
1253 severely fasciated stems and flowers in all three species. Scale bars: 1 cm. **(d)** Quantification of  
1254 *SpriCLV3* paralog-specific transcripts by RNA-seq. **(e)** Locules per fruit after paralog-specific  
1255 gene editing of *SpriCLV3a* and *SpriCLV3b* in *S. prinophyllum*. Single paralog mutants cause a  
1256 subtle shift from bilocular to trilocular fruits; inactivation of both paralogs results in highly  
1257 fasciated fruits. Arrowheads mark locules. Scale bars: 1 cm. **(f)** Quantification of locule number  
1258 in single and double *SpriCLV3a* and *SpriCLV3b* mutants. Proportion of each locule number per  
1259 genotype is shown.

1260

1261 **Figure 5: Pan-genome of African eggplant reveals widespread structural variation, wild**  
1262 **species introgression, and *CLV3* paralog diversification.** **(a)** Images of field-grown African  
1263 eggplant in Mukuno, Uganda (left) and New York, USA (right). **(b)** Ortholog-based phylogeny of  
1264 10 African eggplant accessions covering three main cultivar groups (Gilo, Shum, and Aculeatum)  
1265 and the wild progenitor *S. anguivi*. Representative shoots and fruits are shown for each accession.  
1266 Scale bars: 5 cm (shoots), 2 cm (fruits). Branch lengths reflect coalescent units. **(c)** Pan-genomic  
1267 features across African eggplant reference genome. Frequencies of: (i) sequences private to the  
1268 reference, (ii) core sequence, (iii) genes, (iv) transposable elements, and (v) SVs. **(d)** Average SV  
1269 lengths (bp) for deletions (dotted lines) and insertions (solid lines) across the three African  
1270 eggplant cultivar groups. **(e)** Number of SVs overlapping with genomic features across accessions.  
1271 **(f)** Jaccard similarity of SVs across the African eggplant pan-genome measured against *S. anguivi*  
1272 in 2 Mbp windows. Putative introgression from *S. anguivi* on chromosomes 3, 4, 11, and 12 are  
1273 highlighted by red boxes. **(g)** Close-up of chromosome 4 introgression shown by SV density. **(h)**  
1274 SV density surrounding the *SaetCLV3* locus across the pan-genome. Genomic positions of  
1275 *SaetCLV3a* and *SaetCLV3b* are shown. Window size: 10 kbp. **(i)** Presence-absence and copy  
1276 number variation of *CLV3* across the pan-genome. *CLE9* is absent in all genotypes. *S. aethiopicum*  
1277 and *S. anguivi* are shown for reference. **(j)** Conservation of exonic microsynteny (grey bars)  
1278 between *SangCLV3*, *SaetCLV3<sub>REF</sub>*, and *SaetCLV3<sub>DEL</sub>* haplotypes. Scale: 100 kb. **(k)** Long-reads  
1279 pile-up at the *SaetCLV3* locus identifies a deletion structural variation and distinct *SaetCLV3*  
1280 haplotype in accession 804750136. **(l)** Diagram of deletion-fusion allele of *CLV3* (*SaetCLV3<sub>DEL</sub>*)

1281 arose in accession 804750136. The 7 bp indel and SNPs were used as markers to validate the  
1282 deletion-fusion scenario.

1283

1284 **Figure 6: Pan-genetic dissection of fruit locule variation in African eggplant.** (a) Intraspecific  
1285 crosses between representative accessions of each of the three main cultivated groups of African  
1286 eggplant were used to generate F2 mapping populations for QTL-Sequencing (QTL-seq). (b)  
1287 Major (1) and minor (2) effect QTLs affecting locule number identified by bulk-segregant QTL-  
1288 Seq.  $\Delta$ SNP-indices for three identified QTL on chromosomes 2, 5, and 10 indicate the relative  
1289 abundance of parental variants in bulked pools of F2 individuals (low and high locule classes)  
1290 calculated in 2000 kbp sliding windows. (c) Stacked bar plots showing fruit locule number from  
1291 phylogenetically-arranged African eggplant accessions. Presence of the three mapped QTL alleles  
1292 (different intensity green bars) in each accession are indicated on the phylogenetic tree. (d)  
1293 CRISPR/Cas9 engineered mutant alleles of *SCPL25* serine carboxypeptidase orthologs in tomato  
1294 (*SlycSCPL25*) and *S. prinophyllum* (*SpriSCPL25*) (left), along with representative images of  
1295 transverse fruit sections from mutant plants (right) and quantification of fruit locule number  
1296 (bottom). Scale bars: 1 cm. (e) Schematics comparing the genetic basis of step changes underlying  
1297 increased locule number and fruit size in tomato and African eggplant. Arrowheads in transverse  
1298 fruit depictions indicate locules. Average fruit locule number ( $\mu$ ) and fruit number (n) are indicated  
1299 to the right of stacked bar plots.



1300 **EXTENDED FIGURE LEGENDS**

1301

1302 **Extended Data Figure 1: *Solanum* pan-genome species (selected images), *de novo* assemblies,**  
1303 **and gene annotation pipeline. (a)** Phenotypic diversity of shoots and fruits (where available)  
1304 from a subset of the species selected for the *Solanum* pan-genome. Scale bars: 5 cm (shoots) and  
1305 1 cm (fruits). **(b)** Total sizes of the pan-*Solanum* genome assemblies evaluated by cumulative  
1306 sequence length. Genomes of tomato (*S. lycopersicum*, Heinz SL4.0 and M82) and Brinjal  
1307 eggplant (*S. melongena*, V3) are shown as references. **(c)** Hi-C contact map from *S. candidum*  
1308 shown as a representative example of data used to generate chromosome-scale assemblies. **(d)**  
1309 Flow chart depicting the gene annotation pipeline used in this study, noting the required input data  
1310 (RNA-seq data, protein alignments, and genome sequences), tools, and custom scripts.  
1311 Preprocessing, annotation, homology, functional annotation, and packaging steps are detailed.

1312

1313 **Extended Data Figure 2: Comparative genomic analysis of orthogroup dynamics and**  
1314 **Conservatory analysis of paralogous gene pairs across pan-*Solanum* species. (a)** Functional  
1315 enrichment for orthogroup expansions and contractions in tomato, eggplant, and major *Solanum*  
1316 clades. The top five enriched GO terms per species/clade are shown. Circle size represents gene  
1317 ratio. **(b)** Comparison of orthogroups conservation group size and the subsequent paragroups,  
1318 defined by the number of species having paralogous genes. Note that ~60% of duplicated gene  
1319 orthogroups are conserved across all *Solanum* pan-genome species (Core), while less than 1% of  
1320 the paragroups are Core. **(c)** Duplicated gene pairs classification of the pan-genome species  
1321 according to duplication type. **(d)** Flow chart of the Conservatory tool used to define conserved  
1322 non-coding sequences (CNSs) across pan-genome orthogroups and paragroups. **(e)** Divergence of  
1323 protein and *cis*-regulatory sequences across increasing evolutionary pressure, as measured by  
1324  $K_a/K_s$  values, for the indicated types of gene duplications. For each duplication type the predicted  
1325 mean, residuals, and 0.95 confidence interval of the normalized BLASTP and LastZ scores are  
1326 shown.

1327

1328 **Extended Data Figure 3: Paralog pairs expression analysis. (a)** Schematic of dosage-  
1329 constrained and dosage-unconstrained orthogroups reflecting different degrees of selection on the  
1330 total dosage of paralog pairs across species. Orthogroup 1 has paralog pairs with identical total

1331 dosage across species, whereas orthogroup 2 has different total dosages in each species. For each  
1332 tissue, orthogroup and species, the total dosage of two paralogs is compared with that of the two  
1333 homologs in each of the remaining species, and deviations from the expected ratio of total dosages  
1334 are classified as “unconstrained”. This is repeated for all species that share the orthogroup and  
1335 expressed in the tissue of interest, and the majority classification across species is taken as the  
1336 classification for the entire orthogroup. Therefore, orthogroup 1 is classified as “dosage-  
1337 constrained” while orthogroup 2 is classified as “dosage-unconstrained”. **(b)** The fraction of  
1338 uniquely mapped reads for each tissue sample and species (left), and the average gene expression  
1339 correlation with other samples from the same tissue and species (right). Red arrows in both cases  
1340 point to the five outlier samples excluded from further analysis. **(c)** Sankey plot shows the  
1341 concordance between classification of paralog pairs based on two independent approaches (total  
1342 dosage conservation and conservation of expression levels and profiles). Thickness of lines  
1343 connecting each pair of groups shows the odds ratio of enrichment. **(d)** Line plots showing  
1344 examples of paralog pairs in each of the four groups of paralog expression patterns. **(e)** Functional  
1345 enrichment for paralog pairs from the different groups. The top five enriched GO terms per  
1346 expression group is shown. Circle size represents gene ratio. **(f)** Relationship of protein and *cis*-  
1347 regulatory sequence conservation on the different paralog expression groups over increasing  
1348 evolutionary pressure. For each expression group the predicted mean, 95% confidence interval,  
1349 and residuals of the normalized LastZ score are shown.

1350

1351 **Extended Data Figure 4: Extreme variation in transposable elements and resistant gene**  
1352 **content at the *CLV3* locus across *Solanum*.** **(a)** Gene and transposable element compositions are  
1353 highly variable at the *CLV3* locus across the eggplant clade. While most of the gene content shows  
1354 collinearity, the transposable element profile and density varies considerably. Stacked bars show  
1355 the absolute number and type of transposable element for the window of three genes. **(b)**  
1356 Microsyntenic relationships at the *CLV3* locus across the eggplant clade show dynamic expansions  
1357 and contractions of resistance genes. Resistance genes are identified by blue dots. Presence-  
1358 absence of *CLV3* paralogs is shown as in Figure 4. Lineage-specific *CLV3* duplications denoted  
1359 with asterisks. Window sizes range from 397,829 bp (*S. torvum*) to 634,079 bp (*S. aethiopicum*)  
1360 and are centered on the *CLV3* locus. Functional *CLV3* copies are denoted by full circles while  
1361 truncated/pseudogenized copies are shown as half circles, as in Figure 4. Grey lines illustrate

1362 conservation, while blue lines represent loss of synteny. **(c)** CRISPR/Cas9 gene-edited loss-of-  
1363 function null alleles of *CLV3* genes in *S. prinophyllum* and *S. cleistogamum*. **(d)** CRISPR/Cas9  
1364 gene-edited loss-of-function null alleles of African eggplant *SaetCLV3a/b*. Numbers represent the  
1365 proportion of cloned and sequenced *SaetCLV3a/b* alleles as a ratio of the total number of clones  
1366 sequenced in the three first-generation transgenic (T0) plants showing fasciation phenotypes.

1367

1368 **Extended Data Figure 5: Structural variants and gene copy number variation in the African**  
1369 **eggplant pan-genome. (a)** Structural variant density across all chromosomes in African eggplant  
1370 and its wild progenitor *S. anguivi* in 2 Mbp windows. **(b)** Percentage of structural variants  
1371 overlapping with different genomic features. **(c)** Gene presence-absence and copy number  
1372 variation in 17 orthogroups containing known genes regulating three major domestication traits in  
1373 tomato across the African eggplant pan-genome and *S. anguivi*. Stars mark gene truncation or  
1374 pseudogenization.

1375

1376 **Extended Data Figure 6: Interactions between the *CLV3* and Chr5 African eggplant locule**  
1377 **number QTLs in F2 populations. (a)** Averaged fruit locule number counts for plants from the  
1378 804750136 x PI 424860 (top) and 804750187 x PI 424860 (bottom) segregating F2 populations.  
1379 Average locule counts for the parental genotypes are also shown. **(b)** Stacked bar plots showing  
1380 fruit locule number from ranked F2 population-derived genotypes segregating the reference (REF)  
1381 and alternative (ALT) alleles of *SaetCLV3* and the chromosome 5 QTLs. P: parents.

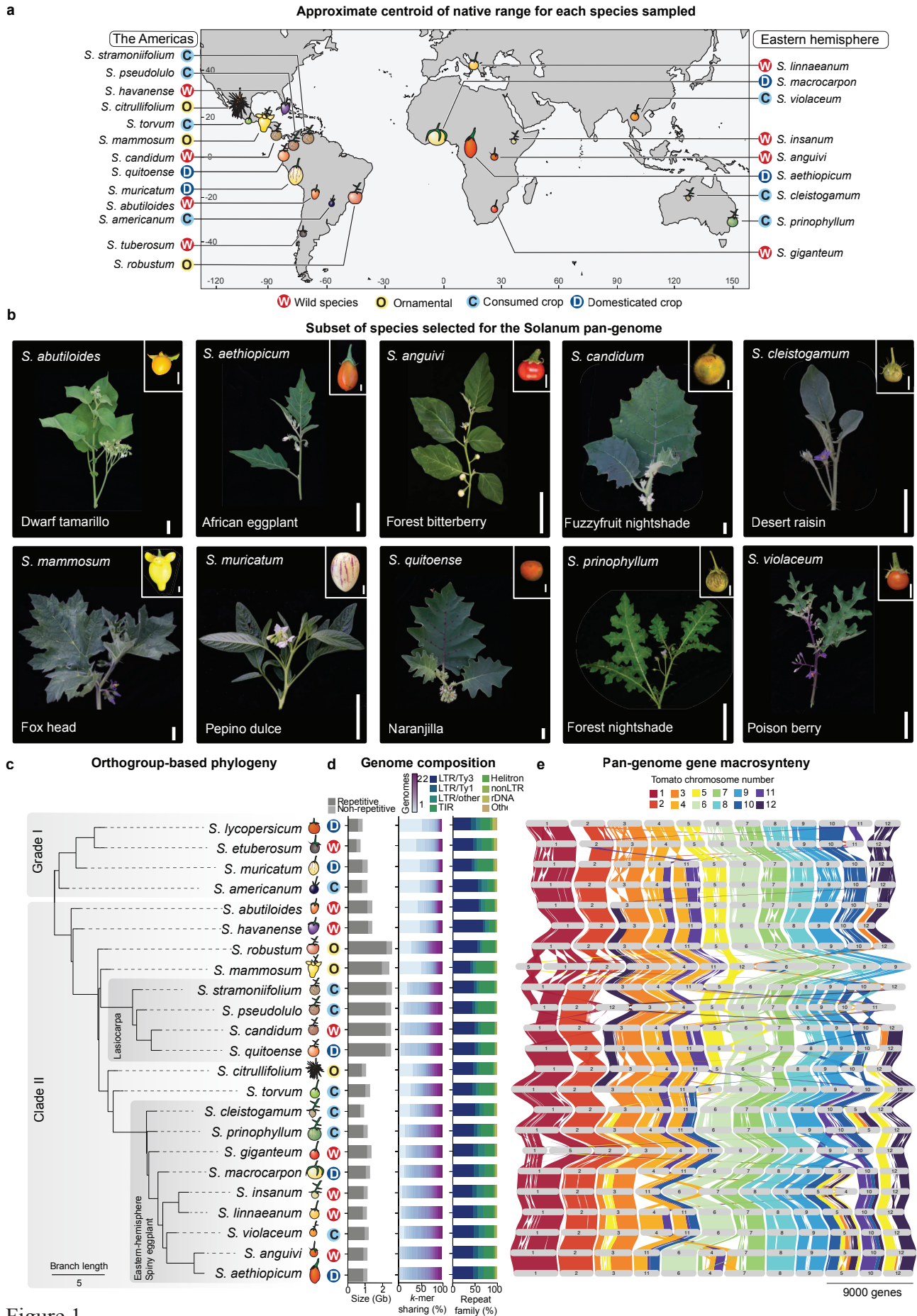


Figure 1

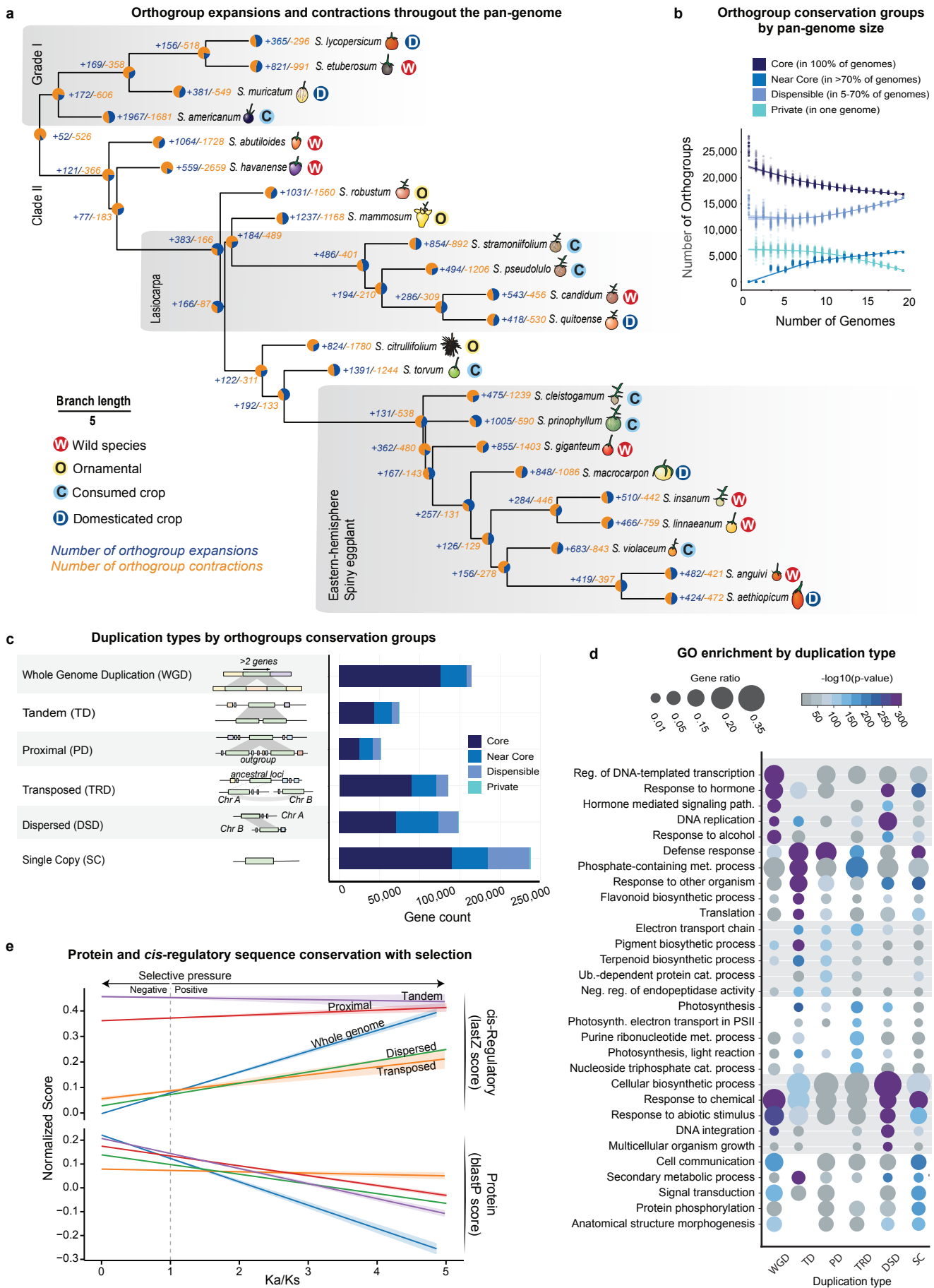


Figure 2

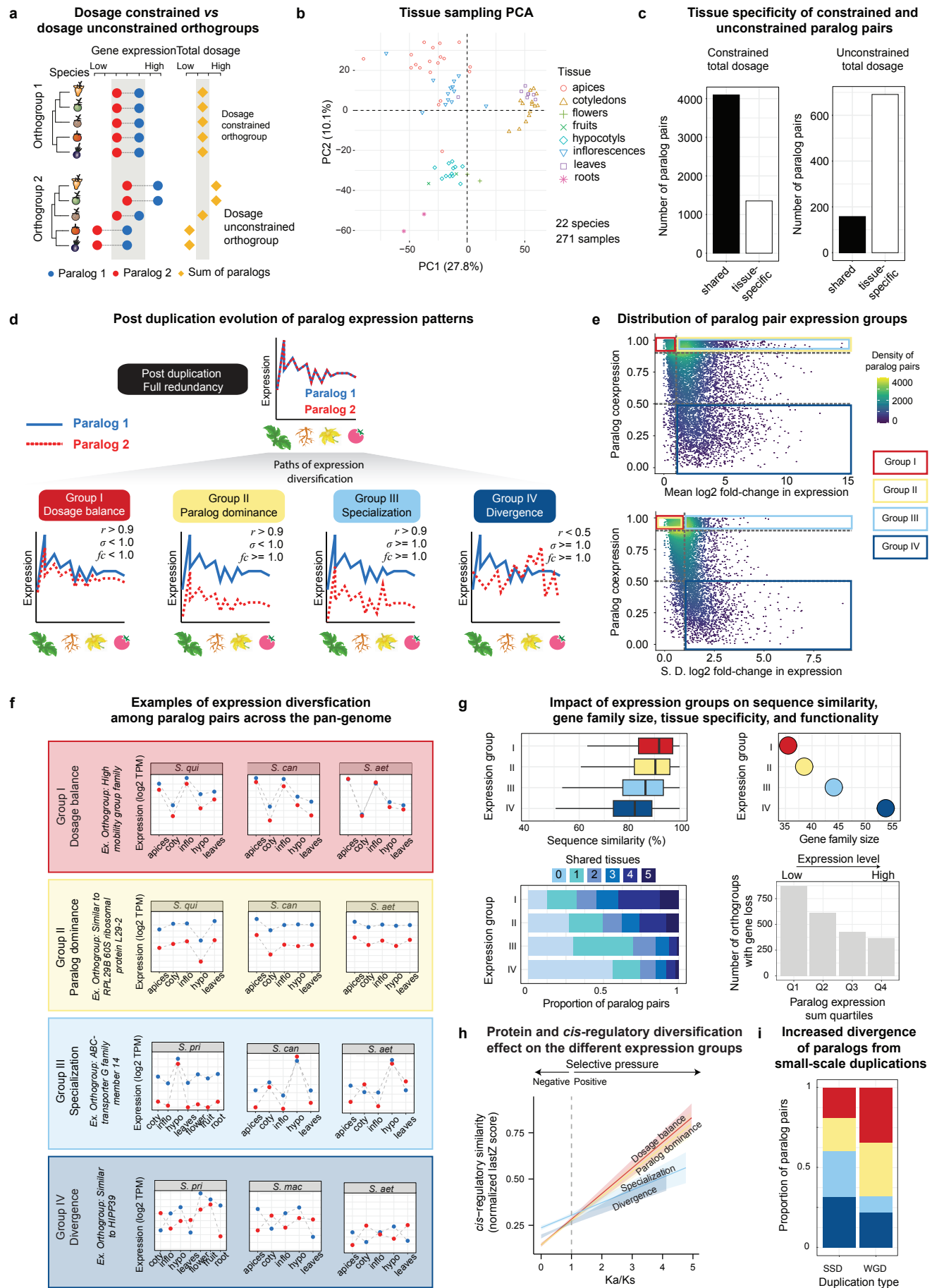
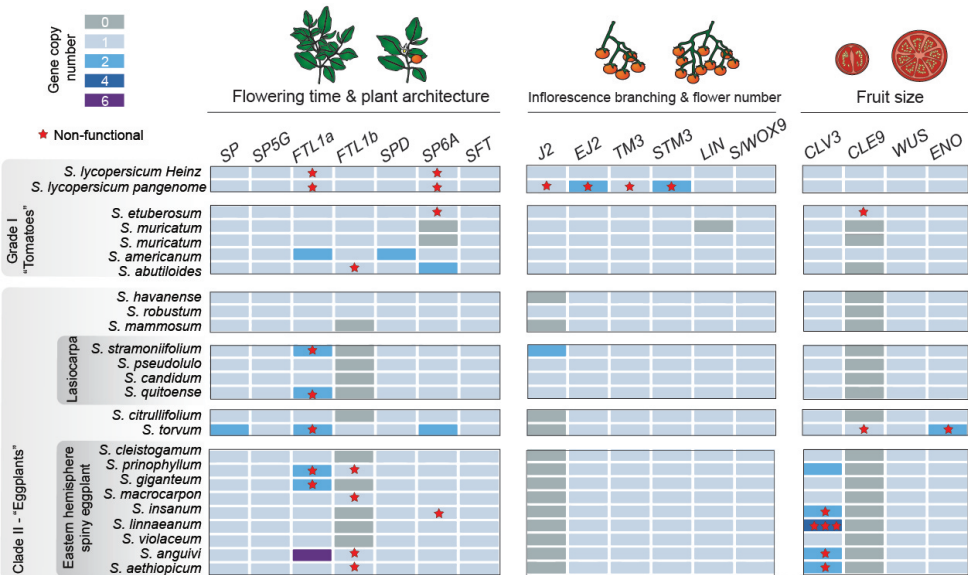
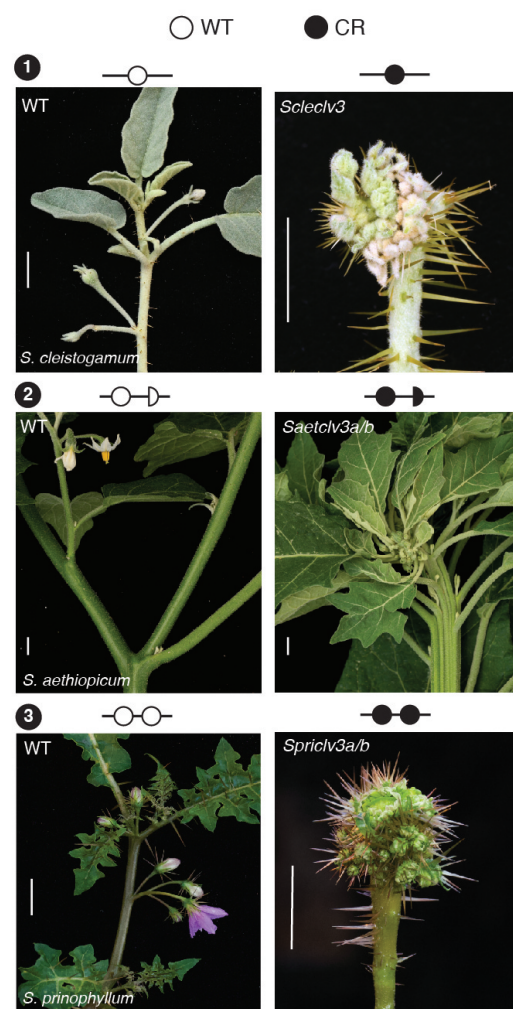


Figure 3

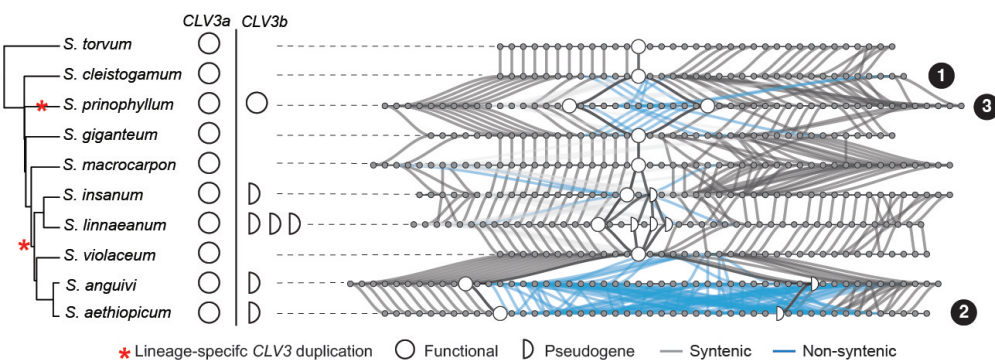
**a Variation in developmental genes and their paralogs underlying major domestication traits**



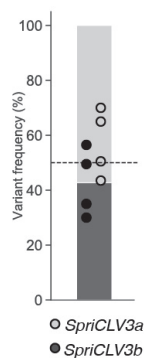
**c Genome editing of *Solanum* CLV3 orthologs**



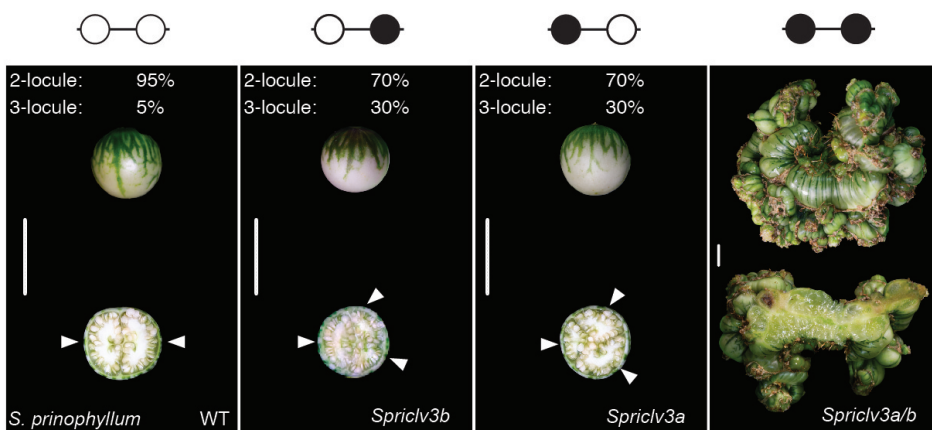
**b Extreme haplotype diversity at the CLV3 locus**



**d CLV3 paralog specific transcript quantification by RNA-seq**



**e Fruit locule number in *SpricCLV3* paralog-specific edited plants**



**f Quantification of paralogous CLV3 dosage relationships in *S. prinophyllum***

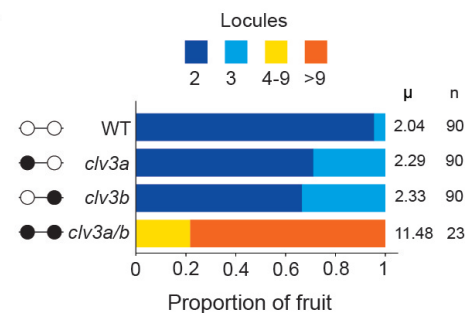
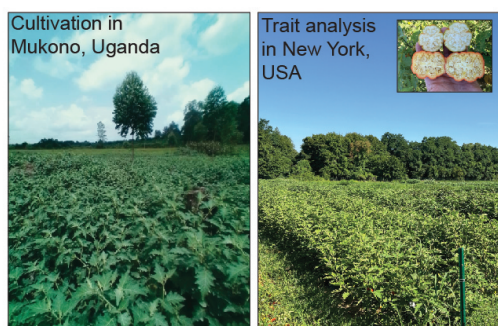
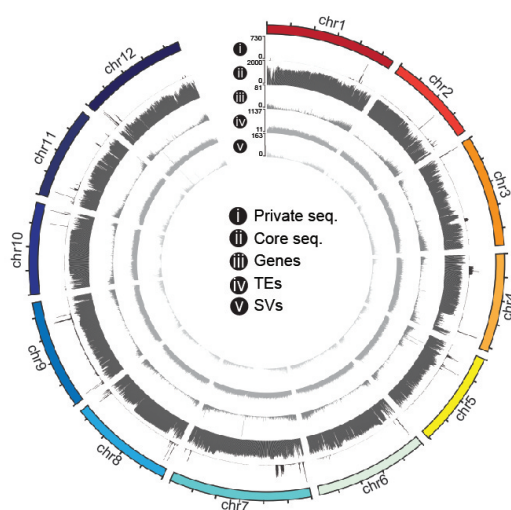


Figure 4

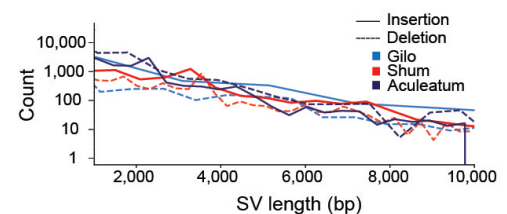
**a Pan-genome enabled genotype-to-phenotype analysis of field grown African eggplant**



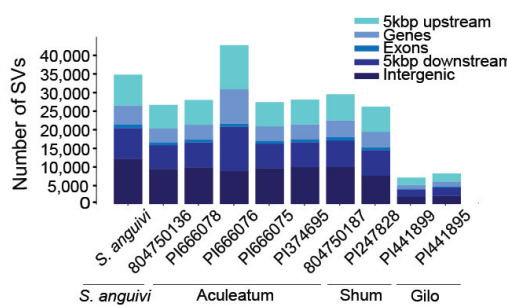
**c Genomic features in reference African eggplant**



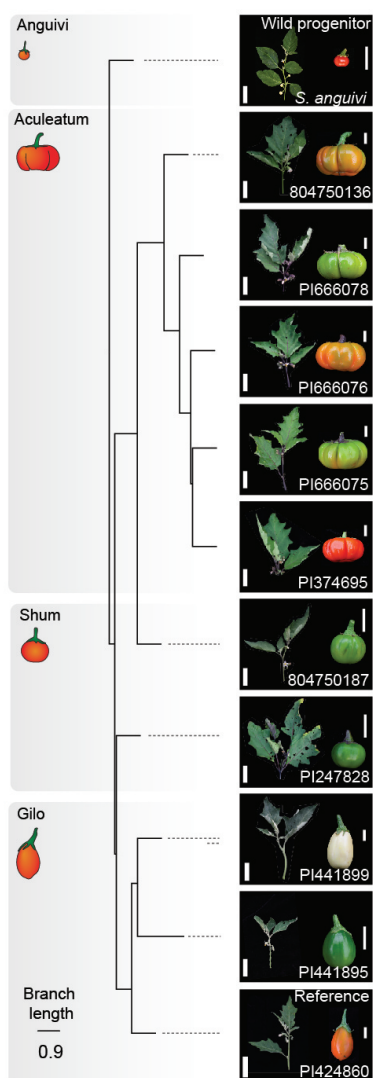
**d Average SV lengths**



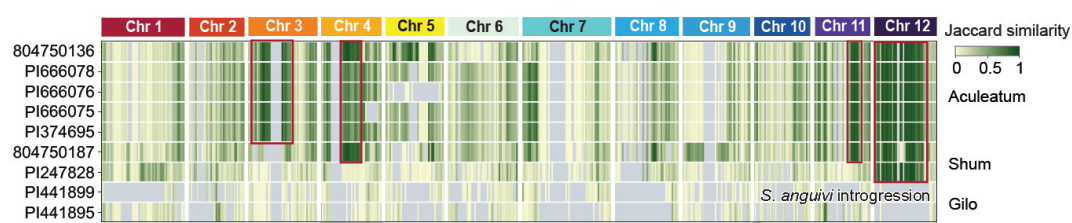
**e Features overlapping SVs by accession**



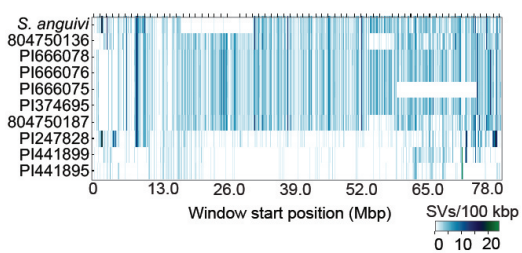
**b 10 African eggplant accessions and wild progenitor (*S. anguivi*)**



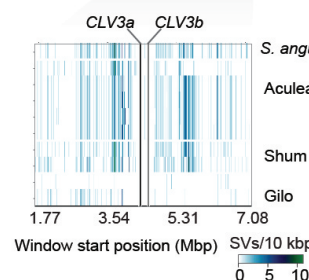
**f Jaccard similarity to *S. anguivi***



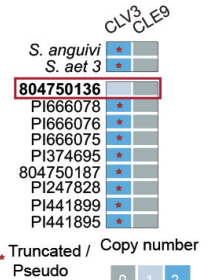
**g SV density surrounding suspected introgression with *S. anguivi* (chr4)**



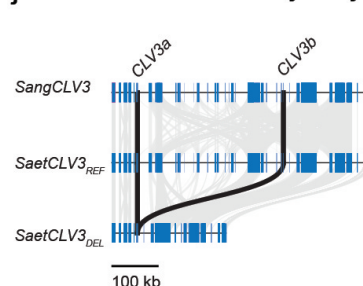
**h SV density surrounding *CLV3* locus (chr10)**



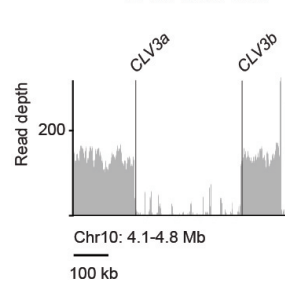
**i Variation in *CLV3* and *CLE9***



**j *CLV3* locus microsynteny**



**k Deletion at *CLV3* locus in 804750136**



**l SNP and INDEL based validation of deletion**

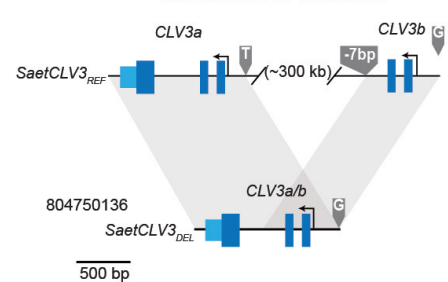
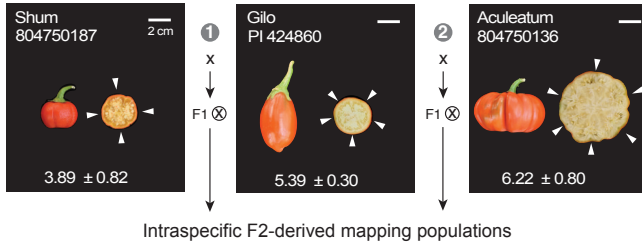


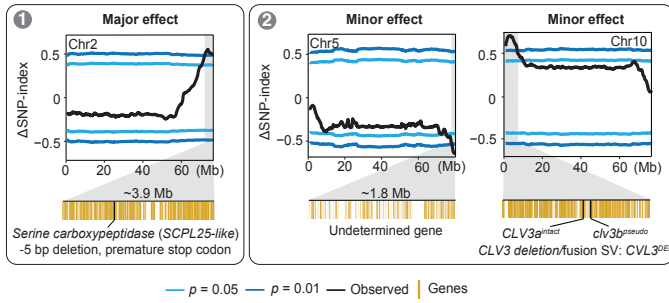
Figure 5



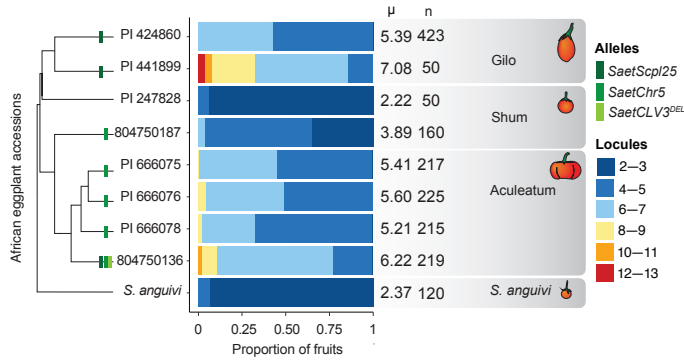
**a Generation of African eggplant locule number QTL mapping populations**



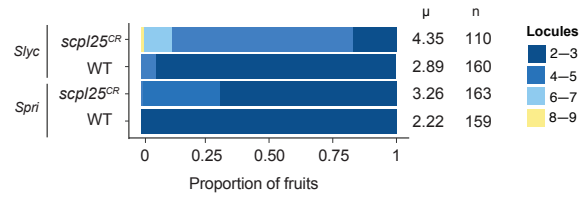
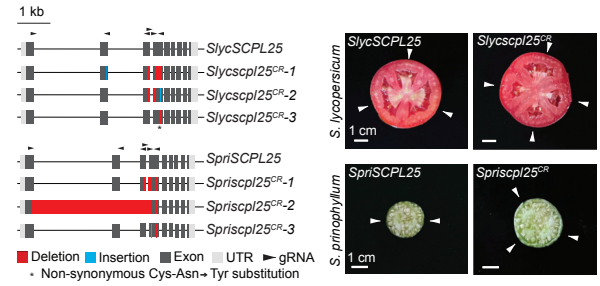
**b QTL-seq identifies major and minor effect loci controlling locule number**



**c Phylogenetic context of fruit locule QTL alleles among African eggplant accessions**



**d Genome editing of SCPL25 yields consistent increase in fruit locule number across species**



**e Paralog contingencies influence repeatability of step-changes in Solanum locule number**

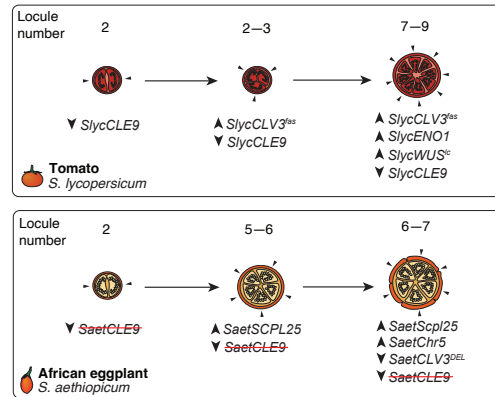


Figure 6