


Super-pangenome analysis highlights genomic diversity and functional association across wild and cultivated tomato species

Received: 1 December 2021

Accepted: 21 February 2023

Published online: 6 April 2023

 Check for updates

Ning Li^{1,11}, Qiang He^{2,11}, Juan Wang¹, Baike Wang¹, Jiantao Zhao³, Shaoyong Huang^{1,4}, Tao Yang¹, Yaping Tang¹, Shengbao Yang¹, Patiguli Aisimutuola¹, Ruiqiang Xu^{1,4}, Jiahui Hu^{1,4}, Chunping Jia^{1,5}, Kai Ma¹, Zhiqiang Li⁶, Fangling Jiang⁷, Jie Gao⁴, Haiyan Lan⁵, Yongfeng Zhou³, Xinyan Zhang³, Sanwen Huang³, Zhangjun Fei^{8,9}, Huan Wang¹⁰✉, Hongbo Li³✉ & Qinghui Yu¹✉

Effective utilization of wild relatives is key to overcoming challenges in genetic improvement of cultivated tomato, which has a narrow genetic basis; however, current efforts to decipher high-quality genomes for tomato wild species are insufficient. Here, we report chromosome-scale tomato genomes from nine wild species and two cultivated accessions, representative of *Solanum* section *Lycopersicon*, the tomato clade. Together with two previously released genomes, we elucidate the phylogeny of *Lycopersicon* and construct a section-wide gene repertoire. We reveal the landscape of structural variants and provide entry to the genomic diversity among tomato wild relatives, enabling the discovery of a wild tomato gene with the potential to increase yields of modern cultivated tomatoes. Construction of a graph-based genome enables structural-variant-based genome-wide association studies, identifying numerous signals associated with tomato flavor-related traits and fruit metabolites. The tomato super-pangenome resources will expedite biological studies and breeding of this globally important crop.

Tomato (*Solanum lycopersicum* L.) is among the most important vegetable crops in terms of global production (<http://www.fao.org/faostat/en/#data/QCL>), also serving as a classic model system for genetic, developmental and physiological studies of fleshy fruits¹. It belongs

to the genus *Solanum* in the nightshade family Solanaceae. Cultivated tomatoes have lost substantial genetic diversity owing to a domestication bottleneck and intensive artificial selection in pursuit of bigger fruits and higher yield², which has impeded tomato improvement.

¹The State Key Laboratory of Genetic Improvement and Germplasm Innovation of Crop Resistance in Arid Desert Regions (Preparation), Key Laboratory of Genome Research and Genetic Improvement of Xinjiang Characteristic Fruits and Vegetables, Institute of Horticultural Crops, Xinjiang Academy of Agricultural Sciences, Urumqi, China. ²Institute of Crop Sciences, Chinese Academy of Agricultural Sciences, Beijing, China. ³Shenzhen Branch, Guangdong Laboratory for Lingnan Modern Agriculture, Shenzhen Key Laboratory of Agricultural Synthetic Biology, Genome Analysis Laboratory of the Ministry of Agriculture and Rural Affairs, Agricultural Genomics Institute at Shenzhen, Chinese Academy of Agricultural Sciences, Shenzhen, China. ⁴College of Horticulture, Xinjiang Agricultural University, Urumqi, China. ⁵College of Life Science and Technology, Xinjiang University, Urumqi, China. ⁶Adsen Biotechnology Co., Ltd., Urumqi, China. ⁷College of Horticulture, Nanjing Agricultural University, Nanjing, China. ⁸Boyce Thompson Institute, Cornell University, Ithaca, NY, USA. ⁹US Department of Agriculture-Agricultural Research Service, Robert W. Holley Center for Agriculture and Health, Ithaca, NY, USA. ¹⁰Biotechnology Research Institute, Chinese Academy of Agricultural Sciences, Beijing, China. ¹¹These authors contributed equally: Ning Li and Qiang He. ✉e-mail: wanghuan@caas.cn; lihongbo_solab@163.com; yueqinghui@xaas.ac.cn

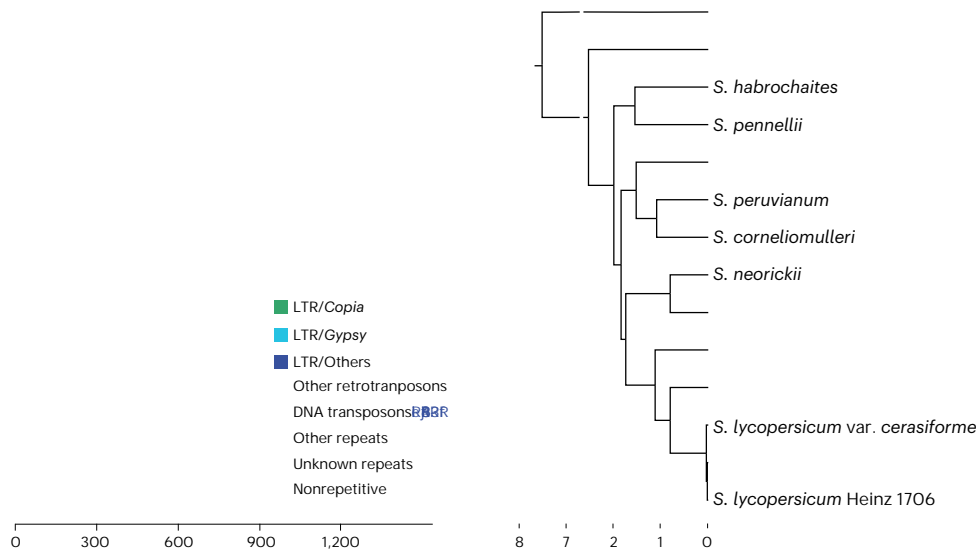
By contrast, wild tomatoes in *Solanum* section *Lycopersicon*, which have adapted to various ecological environments in western South

and utilization of genetic variants in tomato wild relatives. Recently, it was highlighted that a super-pangenome that includes genomic information of many diverse species, especially wild relatives within a genus, could expedite crop improvement²³. Therefore, it is necessary to assemble additional reference genomes for tomato wild relatives to accelerate biological studies and genetic improvement in tomato. In this study, we construct a section-wide super-pangenome by de novo assembling 11 chromosome-level genomes from ten tomato species, representing major clades of tomato wild relatives and their cultivated counterparts in *Lycopersicon*. Comparative analyses reveal the panorama of genomic content, evolutionary history and structural variation across tomato species, empowering the discovery of a wild tomato gene that has the potential to increase yield in modern cultivated tomatoes. These results will provide insight for the construction and exploitation of super-pangenomes in other crop species.

Results

Eleven wild and cultivated tomato reference genomes

To represent the diversity of wild and cultivated tomato species, we selected nine wild tomatoes (eight species from *Solanum* section *Lycopersicon*: *S. habrochaites*, *Solanum chilense*, *Solanum peruvianum*, *Solanum corneliomulleri*, *Solanum neorickii*, *Solanum chmielewskii*, *S. pimpinellifolium* and *S. galapagense*; and one from *Solanum* section *Lycopersicoideae*: *S. lycopersicoideae*) and two diverse domesticated tomatoes (*S. lycopersicum* var. *cerasiforme* and *S. lycopersicum* var. *lycopersicum* cv. M82; Table 1). We assembled a high-quality chromosome-scale reference genome of wild tomato *S. galapagense* 'LA0436', using a hybrid assembly approach integrating Pacific Biosciences (PacBio) sequencing, optical genome mapping (Bionano Genomics) and high-throughput chromosome conformation capture (Hi-C; Supplementary Note and Supplementary Tables 1–5). The 802-Mb final assembly had a contig N50 length of 15.5 Mb, and more than 99.5% of sequences in the final assembly were anchored to the 12 chromosomes, higher than the corresponding percentages for the three existing reference genomes 'LA2093' (99.0%), 'Heinz 1706' (97.5%) and 'LA716' (93.6%) (Table 1). The ten other tomato genomes were also assembled at chromosome level using the above-mentioned strategy, except that Bionano data were not generated. These



genomes had monoploid assembly lengths ranging between 770.0 Mb (*S. chmielewskii*) and 1.2 Gb (*S. lycopersicoides*), close to their predicted genome sizes (Table 1 and Supplementary Tables 6 and 7). More than 99% of Illumina short reads and 95.7% of ESTs could be mapped to the 11 tomato genome assemblies, and 94.0% of embryophyte Benchmarking Universal Single-Copy Orthologs (BUSCO)²⁴ were captured in these assemblies, indicative of their high completeness (Supplementary Tables 8–10).

We combined *ab initio* prediction, homology search and transcriptome mapping approaches for protein-coding gene prediction (Methods), resulting in gene numbers ranging from 31,613 (*S. chmielewskii*) to 34,375 (*S. chilense*), similar to that of Heinz 1706 (35,768) but fewer than that of LA716 (44,965) (Table 1). A total of 81.7% to 89.5% of exons of the predicted genes were supported by transcript data, suggesting the high quality of gene predictions. All assembled genome sequences and annotations are publicly accessible through a web-based database (<http://caastomato.biocloud.net>).

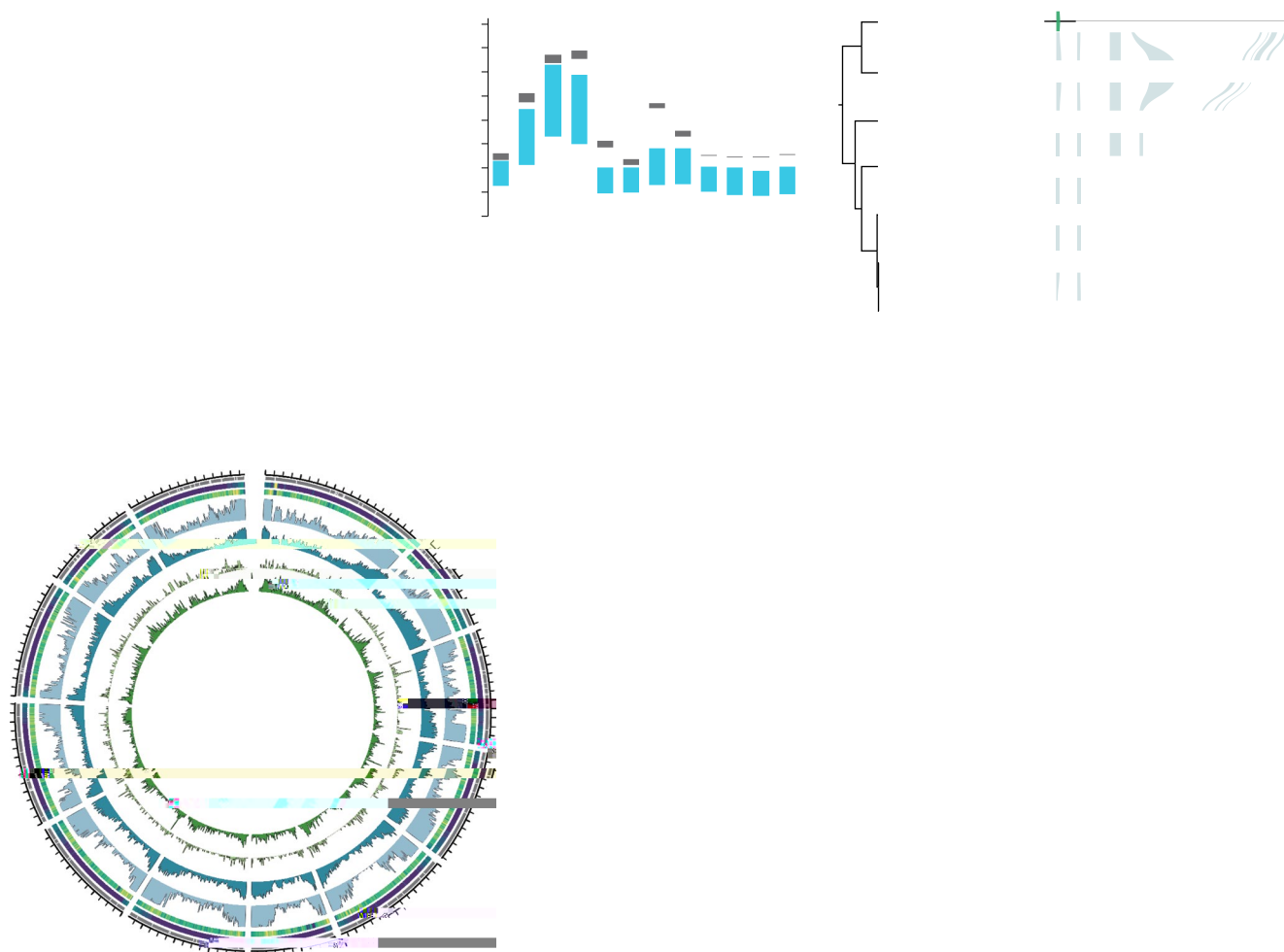
Eukaryotic genomes are rich in transposable elements (TEs), which shape genome evolution through expansions, eliminations and transpositions²⁵. The TE contents of the 11 tomato genomes ranged from 64.3% to 74.5%, with long terminal repeat retrotransposons (LTR-RTs) representing the most abundant class of TE (Fig. 1a). A higher abundance of *Gypsy* LTR-RTs was found in *S. lycopersicoides*, which possibly contributed to it having the largest assembled genome size (1.2 Gb) among the tomato species²² (Fig. 1a). To trace the evolutionary history of the expanded TEs in *S. lycopersicoides*, we estimated insertion times of 162,216 intact LTR-RTs and detected a lineage-specific burst of *Gypsy* LTR-RTs occurring c. 2 million years ago (Ma) in *S. lycopersicoides*, after its divergence from potato, probably leading to its large extant genome (Supplementary Fig. 3). Notably, we observed recent amplification of *Gypsy* and *Copia* LTR-RTs in four wild tomato species (*S. lycopersicoides*, *S. corneliomulleri*, *S. peruvianum* and *S. chilense*; Supplementary Fig. 3), implying that these wild species may have increasing degrees of genomic diversity and environmental adaptability compared with cultivated tomatoes. These results provide insight into the role of TEs in genome evolution of the *Solanum* genus.

Phylogeny of *Lycopersicon* and neighboring species

Reconstructing the phylogeny of *Lycopersicon* species has been problematic owing to the conflict between gene trees and morphological trees, especially for the wild tomato clade³. The phylogenetic relationship between *S. pennellii* and other tomatoes remains unresolved⁶, owing largely to limited available genomic data, despite *S. pennellii* being considered to be a unique group based on morphological classification. Using 9,343 single-copy orthologous genes, we inferred the phylogeny of ten wild and three domesticated tomatoes, using potato (*Solanum tuberosum*) as an outgroup; the results indicated that section *Lycopersicoides* (including *S. lycopersicoides*) was sister to section *Lycopersicon* (Fig. 1b), consistent with previous research³. Based on the phylogeny, we resolved the polytomy issue in *Lycopersicon* and unambiguously classified *Lycopersicon* species into four main clades. Clade I encompassed two species, *S. pennellii* and *S. habrochaites*, which diverged from the common ancestor of the other wild and cultivated tomatoes (except *S. lycopersicoides*) c. 1.97 Ma. Clade IV, which comprised domesticated tomatoes and two closely related wild species (*S. galapagense* and *S. pimpinellifolium*), divided from the ancestor of clade III (*S. neorickii* and *S. chmielewskii*) approximately 1.73 Ma. Similar to a recent study of *Oryza* genus evolution²⁶, a few conflicts were observed between the phylogeny constructed using genes from one chromosome and that built using whole-genome genes (Fig. 1b and Supplementary Fig. 4). For example, within *Lycopersicon*, phylogenetic analyses using genes from chromosomes 1, 2, 9 and 11 showed that *S. pennellii* was sister to other wild and cultivated tomato species, rather than clustering into a monophyletic group with *S. habrochaites* as inferred from the genome-wide phylogeny (Supplementary Fig. 4), suggesting possible incomplete lineage sorting and/or hybridization events. These results enhance our understanding of the evolutionary history within *Solanum* section *Lycopersicon*.

Super-pangenome of tomato

Although pangenomes for cultivated tomato and its close wild relatives have been reported¹³, the gene pool of *Lycopersicon*, which contains wild and cultivated tomato species, remains largely inaccessible. Here, we extended the tomato pangenome that integrates genomes from



three *Solanum* species¹³ to a super-pangenome covering 11 species in the *Solanum* genus. We defined 40,457 pangenome families by clustering protein-coding genes of the 11 chromosome-scale genomes assembled herein and two previously released genomes^{8,18}; this number of gene families was higher than that of the *Oryza* genus²⁶ but lower than that of soybean²⁷. The number of gene families increased rapidly when including more genomes, suggesting that the 13 genomes are diverse and that a single reference genome cannot capture the full genetic diversity in tomato (Fig. 2a). Only 54.0% of gene families were conserved among the 13 tomato genomes (core gene families), and the number of core genes (23,839) was lower than that of the previously reported pangenome of 540

used in this study. The dispensable gene families (present in two to 12 accessions) occupied 38.4% of gene families, and 7.6% of pangenome families were categorized as accession-specific.

Gene ontology (GO) enrichment analysis showed that core genes were enriched for biological processes including carboxylic acid, lipid or organic substance metabolic process, RNA modification or processing and amide transport, consistent with the results of a previous study¹³ (Supplementary Tables 11 and 12), whereas the dispensable genes were enriched for terpenoid biosynthesis, telomere maintenance, mitochondrial electron transport and photosynthesis (Supplementary Fig. 6). Expression levels of core genes were significantly

previous tomato pangenome¹³ were captured in our super-pangenome (Supplementary Table 13), and we also identified 9,320 nonredundant genes absent from the reported tomato pangenome¹³ (Supplementary Note and Supplementary Table 14), indicating the rich diversity of the 13 wild and domesticated tomatoes. This super-pangenome dataset lays a foundation for exploration and exploiting of genes or alleles in wild tomato species.

Extensive variation among wild and cultivated tomatoes

Despite efforts to characterize genetic variants among cultivated tomatoes and their proposed progenitor species *S. pimpinellifolium*^{12,13,16}, the genetic diversity among distantly related wild tomato species, for example, *S. peruvianum*, *S. habrochaites* and *S. chilense*, remains poorly explored. We identified 2.0–8.1 million SNPs and 0.6–1.5 million small indels (< 50 base pairs (bp) in size) in the 12 tomato genomes, relative to the reference *S. galapagense* genome. The total number of SNPs and small indels (42.4 M) was much higher than that of each accession (Supplementary Tables 15 and 16), suggesting a diverse nature among the 12 wild and cultivated tomato accessions (Supplementary Note). Leveraging genome alignments, we identified 103,333 insertions, 119,794 deletions, 41,960 CNVs, 23,516 translocations and 1,320 inversions (<1 Mb in length) in the 12 tomato accessions compared with the *S. galapagense* genome (Supplementary Tables 17 and 18). Species in clade II (*S. chilense*, *S. peruvianum* and *S. corneliomulleri*) contained markedly varied numbers of SVs (Fig. 2b), possibly associated with the recent proliferation of LTR-RTs in those genomes (Supplementary Fig. 3). The majority of insertions, deletions and CNVs were shorter than 2 kb, 2 kb and 8 kb, respectively, and most of the translocations had lengths shorter than 20 kb, whereas some inversions were longer than 300 kb (Supplementary Fig. 17). We found that insertions and deletions were more likely to be found at both ends of the chromosomes, consistent with previous studies^{12,20}, whereas inversions and translocations were randomly distributed along the 12 chromosomes (Fig. 2c). SVs were more likely to occur at repeat regions than nonrepeat genomic regions (Student's *t* test, $P = 1.03 \times 10^{-4}$). We further identified 5,186 large indels (>50 bp) fixed either in all wild or all domesticated tomato genomes investigated in this study, some of which led to insertions of protein-coding genes present only in certain wild tomato genomes (Supplementary Table 19 and Fig. 2d). Further functional characterization of these variants may enable a better understanding of the genetic basis of phenotypic divergence between domesticated tomatoes and their wild relatives.

Previous studies have identified several SVs responsible for phenotypic variation, including a 1.4-kb deletion in the *CSR* gene resulting in increased fruit weight²⁸, a 7.1-kb deletion in the *LNK2* locus responsible for a light-conditional clock deceleration²⁹, an 85-bp deletion in the promoter of *ENO* that regulates floral meristem activity³⁰ and a CNV affecting *NSGT* associated with biosynthesis of a fruit flavor volatile guaiacol¹². These SVs were all accurately detected in this study (Supplementary Figs. 18–21), indicating the broad diversity of our collection. Two different alleles (4,724 bp and 4,151 bp) have been identified at 149 bp upstream of *TomLoxC*, a gene encoding a 13-lipoxygenase; the 4,151-bp allele was reported to contribute to desirable fruit flavor and is rare in cultivated tomatoes¹³. We found that *S. pennellii*, *S. habrochaites*, *S. chilense* and *S. neorickii* carried the 4,151-bp allele upstream of *TomLoxC* (Supplementary Fig. 22), suggesting that these wild species have the potential to improve fruit flavor in cultivated tomato by backcrossing. The extensive variation among wild and cultivated tomato species presented herein provides access for further harnessing of the genetic diversity of distantly related wild tomato species in genomic-based breeding.

Hidden genetic diversity of tomato wild species

Large inversions have been reported to suppress recombination by reducing crossing-over^{31,32}, resulting in severe linkage drag when

conducting backcross breeding. To overcome this, it is necessary to choose donor lines without inverted segments harboring targeted genes. However, a holistic view of genome-wide inversions is not available, owing to the lack of chromosome-scale wild tomato genomes. Based on the 11 high-quality tomato genomes, we identified 12 (*S. lycopersicum* var. *lycopersicum* cv. Heinz 1706) to 42 (*S. chmielewskii*) megabase-scale inversions compared with the *S. galapagense* genome (Supplementary Table 20). Notably, a 7.1-Mb inversion on chromosome 3, carrying 55 genes, was present in all clade IV tomato accessions compared with other wild species (except *S. pennellii*) and was supported by clear chromatin interactions around the breakpoints when Hi-C reads of *S. neorickii* and *S. chmielewskii* were mapped to the *S. galapagense* genome (Fig. 2e). This inversion might occur after the divergence between species from clade IV and other clades. Given that *S. pennellii* does not carry this inversion within this region, this wild tomato species would be an ideal donor parent to introduce possibly favored genes within this 7.1-Mb segments into elite cultivars by backcrossing.

Previous research reported a tomato pan-SV map, which was built by long-read sequencing of 100 cultivated and closely related wild tomato accessions¹². Compared with this pan-SV map, 180,314 out of the 224,447 SVs were exclusively identified in this study, of which 4,124 (2.3%) were localized within coding regions (CDS) of 3,515 genes (Supplementary Note and Supplementary Table 21), suggesting that the majority of SVs found in this study were captured owing to the inclusion of distantly related wild tomato species. Integrating our identified SVs with the pan-SV dataset generated 153,873 insertions, 203,364 deletions, 2,952 inversions and 45,987 duplications in 112 tomato accessions (12 in this study and 100 in the pan-SV map), allowing us to investigate the divergence of SVs during tomato evolution. We divided these 112 accessions into four groups: wild (19 non-*S. pimpinellifolium* wild accessions), SP (22 *S. pimpinellifolium* accessions), SLC (24 *S. lycopersicum* var. *cerasiforme* accessions) and SLL (47 big-fruited *S. lycopersicum* var. *lycopersicum* accessions; Supplementary Fig. 24a). The vast majority of SVs displayed relatively low frequencies (<0.25) in all four groups, and accessions from the wild group contained a higher proportion of SVs with presence frequency between 0 and 0.25 (Supplementary Fig. 24b). We observed that 8,094 SVs exhibited significant frequency changes between the wild and cultivated (SLL and SLC) groups (Fisher's exact test, false discovery rate (FDR) < 0.01; Supplementary Fig. 24c), affecting upstream regions and exons of 2,585 genes. Functional analyses indicated that these genes were mainly enriched for biological processes such as meristem development and ammonium transport (Supplementary Fig. 24d). We further identified 388 highly divergent SVs between wild and cultivated tomatoes, which disrupted CDS of 328 genes by causing frameshift, loss of exons or in-frame insertions (Supplementary Table 22). These results suggest that SVs in these distantly related wild tomatoes have undergone distinct evolutionary trajectories compared with cultivated tomatoes and their progenitors. Our analyses also provide a candidate dataset for further characterizing genes underlying phenotypes with great divergence between wild and cultivated tomatoes.

A wild tomato cytochrome P450 gene that increases yield

A major goal of tomato breeding is to increase yield by developing varieties with larger fruit size and/or more effective shoot branches. Regulation of shoot architecture is thus of great interest to the tomato research community³³. Wild tomato species usually display a markedly greater number of lateral fruit-bearing branches than their domesticated counterpart; however, whether we can introduce this trait into cultivated tomatoes, especially modern processing tomato varieties, remains elusive. Among the 388 highly divergent SVs between wild and cultivated tomatoes that greatly affected gene CDS, a 244-bp deletion, showing the second most significant frequency change (FDR = 1.43×10^{-8} ; Supplementary Fig. 24c), was present in the first exon of *Sgall2gO15720* (Fig. 3a,b). This gene encodes a protein belonging to the cytochrome

P450 superfamily, which has been reported to play important parts in plant growth, development and secondary metabolite biosynthesis³⁴. The 244-bp deletion was found in 22.22% of the 19 wild accessions and 100% of cultivated tomatoes, which represented the derived state, as this deletion was absent from all the nine wild tomato species used in this study (Fig. 3a,b and Supplementary Fig. 25). *Sgal12g015720* was

mapping of short reads from SV regions and thus SV genotyping^{38,39}. We then genotyped these SVs in a tomato population comprising 321 accessions⁷ and performed SV-based GWAS for 32 flavor-related compounds² and 362 fruit metabolites⁴⁰. For comparison, we also called SNPs and indels from the 321 accessions and employed SNP-based GWAS.

Significantly associated signals were detected for 17 flavor volatiles and 249 fruit metabolites. Surprisingly, we observed that only 5.2% (161) of peaks (quantitative trait loci) overlapped (800-kb flanking region) between SV-based and SNP-based GWAS results, and 21.3% (658) could only be identified by SVs. The remaining 2,263 (73.4%) were exclusively detected by SNPs (Fig. 4a,b and Supplementary Table 23). Examples included a peak at 65.2 Mb on chromosome 10 that could only be detected using SVs, which was strongly associated with the content of geranylacetone ($P = 7.91 \times 10^{-9}$), one of the important tomato flavor volatiles contributing a leafy flavor to fruits (Fig. 4c and Supplementary Fig. 29). The leading SV was a 347-bp deletion, and the content of geranylacetone in tomato fruits significantly differed

between accessions carrying the reference allele and those carrying the alternative allele (Student's *t* test, $P = 3.7 \times 10^{-8}$, Fig. 4c). Similarly, we detected significantly associated SVs for the content of additional metabolites (Fig. 4d–f, Supplementary Figs. 30–32 and Supplementary Table 23). Tomato accessions carrying alleles of corresponding leading SVs showed significantly increased content of these metabolites -

species contains valuable breeding materials. However, the availability of only a few wild tomato genomes has hampered the exploration and utilization of alleles and gene repertoire in those wild species. The chromosome-scale reference genomes for nine wild tomato species presented here offer valuable resources for not only comparative genomics but also biological studies and molecular breeding in tomato. Notwithstanding, our dataset still lacks three wild tomato species in *Solanum* section *Lycopersicon* (*Solanum cheesmaniae*, *Solanum huaylasense* and *Solanum arcanum*). *S. cheesmaniae* is endemic to the Galápagos island with yellow to orange fruits⁴¹, whereas *S. huaylasense* and *S. arcanum* are wild tomatoes segregated from *S. peruvianum*⁴². Development of their genome sequences and annotation will further enrich our understanding of the biodiversity and evolutionary trajectory within *Lycopersicon*.

Although pangenomes for numerous crops have been reported, most of them incorporated one or a few species⁴³. Here, we constructed a super-pangenome by analyzing 11 distinct tomato species, representative of major wild and cultivated tomato clades. Coupling this with an existing dataset¹², we identified a wild tomato gene that could increase fruit yield by an average of 67.1% in OE transgenic lines (Fig. 3d–f). As both OE lines and ILs carrying this gene had higher numbers of fruit-bearing branches (Fig. 3d and Supplementary Fig. 27), we anticipate its use in modern processing tomatoes. According to tomato population resequencing data, this gene was predominantly found in wild tomato accessions (52% of *S. pimpinellifolium*, 80% of *S. cheesmaniae* and 100% of *S. galapagense*), in contrast to a mere 6% and 19% in cultivated tomato forms *S. lycopersicum* var. *lycopersicum* and *S. lycopersicum* var. *cerasiforme*, respectively (Supplementary Table 24). These results indicate that this gene, although potentially important, has not been widely utilized in tomato breeding programs. Backcrossing would be an ideal approach to introduce this gene into cultivated tomatoes from wild species. However, hybridization between wild and cultivated crops may lead to severe repression of genetic recombination, owing largely to the presence of large-scale genomic divergence, such as large inversions^{10,31}. This may ultimately result in the introduction of exotic genomic fragments carrying unfavorable alleles that are hard to purge⁴⁰. We did not observe chromosomal rearrangements between the genome of Heinz 1706 and those of eight out of the nine wild species surrounding this gene (Supplementary Table 25), suggesting that introgression of this gene by backcrossing, when the donor parent is properly selected, would be less likely to cause linkage drag.

To facilitate the utilization of genetic diversity from our super-pangenome, we constructed a graph-based genome reference for wild and cultivated tomatoes by integrating SV information for 112 tomatoes from 11 *Solanum* species into the linear reference sequence, offering a powerful platform for population-level SV genotyping. As previous research has suggested that SVs are more likely to be causal variants in tomato¹⁶, further studies could use this graph-based genome to perform SV-based association analyses to identify additional signals responsible for agronomically important traits. However, the current graph-based tomato genome is only capable of storing certain types of SVs: insertions, deletions and inversions. Other SVs of relatively high complexity, such as inverted duplications and translocations, cannot yet be integrated. Furthermore, SVs with multiple alleles are not represented in the graph, as downstream analytic pipelines can only handle biallelic variants. It is possible that an insertion with distinct inserted fragments in various individuals contributes to different phenotypic outcome. We anticipate further implementation of relevant tools and algorithms that could tackle these issues. This research will accelerate comparative genomics and biological studies in tomato and shed light on the utilization of super-pangenomes in crop improvement.

Online content

Any methods, additional references, Nature Portfolio reporting summaries, source data, extended data, supplementary information,

acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41588-023-01340-y>.

References

- Giovannoni, J. J. Genetic regulation of fruit development and ripening. *Plant Cell* **16**, S170–S180 (2004).
- Tieman, D. et al. A chemical genetic roadmap to improved tomato flavor. *Science* **355**, 391–394 (2017).
- Peralta, I. E., Spooner, D. M. & Knapp, S. Taxonomy of wild tomatoes and their relatives (*Solanum* sect. *Lycopersicoides*, sect. *Juglandifolia*, sect. *Lycopersicon*; Solanaceae). *Syst. Bot. Monogr.* **84**, 1–186 (2008).
- Rick, C. M. Perspectives from plant genetics: the Tomato Genetics Stock Center. In *Genetic resources at risk: scientific issues, technologies, and funding policies. Proceedings of a symposium, American Association for the Advancement of Science annual meeting, San Francisco, California, USA, 16 January 1989* (Eds McGuire, P. E. & Qualset, C. O.) 11–19 (Genetic Resources Conservation Program, University of California, 1990).
- Mutschler, M. A. et al. QTL analysis of pest resistance in the wild tomato *Lycopersicon pennellii*: QTLs controlling acylsugar level and composition. *Theor. Appl. Genet.* **92**, 709–718 (1996).
- Spooner, D. M., Peralta, I. E. & Knapp, S. Comparison of AFLPs with other markers for phylogenetic inference in wild tomatoes [*Solanum* L. section *Lycopersicon* (Mill.) Wettst.]. *TAXON* **54**, 43–61 (2005).
- Beckles, D. M., Hong, N., Stamova, L. & Luengwilai, K. Biochemical factors contributing to tomato fruit sugar content: a review. *Fruits* **67**, 49–64 (2012).
- The Tomato Genome Consortium. The tomato genome sequence provides insights into fleshy fruit evolution. *Nature* **485**, 635–641 (2012).
- Hosmani, P. S. et al. An improved de novo assembly and annotation of the tomato reference genome using single-molecule sequencing, Hi-C proximity ligation and optical maps. Preprint at *bioRxiv* <https://doi.org/10.1101/767764> (2019).
- Lin, T. et al. Genomic analyses provide insights into the history of tomato breeding. *Nat. Genet.* **46**, 1220–1226 (2014).
- Aflitos, S. et al. Exploring genetic variation in the tomato (*Solanum* section *Lycopersicon*) clade by whole-genome sequencing. *Plant J.* **80**, 136–148 (2014).
- Alonge, M. et al. Major impacts of widespread structural variation on gene expression and crop improvement in tomato. *Cell* **182**, 145–161.e23 (2020).
- Gao, L. et al. The tomato pan-genome uncovers new genes and a rare allele regulating fruit flavor. *Nat. Genet.* **51**, 1044–1051 (2019).
- Sherman, R. M. & Salzberg, S. L. Pan-genomics in the human genome era. *Nat. Rev. Genet.* **21**, 243–254 (2020).
- Della Coletta, R., Qiu, Y., Ou, S., Hu, M. B. & Hirsch, C. N. How the pan-genome is changing crop genomics and improvement. *Genome Biol.* **22**, 3 (2021).
- Zhou, Y. et al. Graph pangenome captures missing heritability and empowers tomato breeding. *Nature* **606**, 527–534 (2022).
- Yu, X. et al. Chromosome-scale genome assemblies of wild tomato relatives *Solanum habrochaites* and *Solanum galapagense* reveal structural variants associated with stress tolerance and terpene biosynthesis. *Hortic. Res.* **9**, uhac139 (2022).
- Bolger, A. et al. The genome of the stress-tolerant wild tomato species *Solanum pennellii*. *Nat. Genet.* **46**, 1034–1038 (2014).
- Schmidt, M. H.-W. et al. De novo assembly of a new *Solanum pennellii* accession using nanopore sequencing. *Plant Cell* **29**, 2336–2348 (2017).

20. Wang, X. et al. Genome of *Solanum pimpinellifolium* provides insights into structural variants during tomato breeding. *Nat. Commun.* **11**, 1–11 (2020).
21. Takei, H. et al. *De novo* genome assembly of two tomato ancestors, *Solanum pimpinellifolium* and *Solanum lycopersicum* var. *cerasiforme*, by long-read sequencing. *DNA Res* **28**, dsaa029 (2021).
22. Powell, A. F. et al. A *Solanum lycopersicoides* reference genome facilitates insights into tomato specialized metabolism and immunity. *Plant J.* **110**, 1791–1810 (2022).
23. Khan, A. W. et al. Super-pangenome by integrating the wild side of a species for accelerated crop improvement. *Trends Plant Sci.* **25**, 148–158 (2020).
24. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**, 3210–3212 (2015).
25. Chen, J. et al. Tracking the origin of two genetic components associated with transposable element bursts in domesticated rice. *Nat. Commun.* **10**, 1–10 (2019).
26. Stein, J. C. et al. Genomes of 13 domesticated and wild rice relatives highlight genetic conservation, turnover and innovation across the genus *Oryza*. *Nat. Genet.* **50**, 285–296 (2018).
27. Liu, Y. et al. Pan-genome of wild and cultivated soybeans. *Cell* **182**, 162–176.e13 (2020).
28. Mu, Q. I. et al. Fruit weight is controlled by Cell Size Regulator encoding a novel protein that is expressed in maturing tomato fruits. *PLoS Genet.* **13**, e1006930 (2017).
29. Mora-García, S. & Yanovsky, M. J. A large deletion within the clock gene LNK2 contributed to the spread of tomato cultivation from Central America to Europe. *Proc. Natl Acad. Sci. USA* **115**, 6888–6890 (2018).
30. Yuste-Lisbona, F. J. et al. ENO regulates tomato fruit size through the floral meristem development network. *Proc. Natl Acad. Sci. USA* **117**, 8187–8195 (2020).
31. Wellenreuther, M. & Bernatchez, L. Eco-evolutionary genomics of chromosomal inversions. *Trends Ecol. Evol.* **33**, 427–440 (2018).
32. Huang, K. & Rieseberg, L. H. Frequency, origins, and evolutionary role of chromosomal inversions in plants. *Front. Plant Sci.* **11**, 296 (2020).
33. Xia, X. et al. Brassinosteroid signaling integrates multiple pathways to release apical dominance in tomato. *Proc. Natl Acad. Sci. USA* **118**, e2004384118 (2021).
34. Vasav, A. P. & Barvkar, V. T. Phylogenomic analysis of cytochrome P450 multigene family and their differential expression analysis in *Solanum lycopersicum* L. suggested tissue specific promoters. *BMC Genomics* **20**, 1–13 (2019).
35. Eshed, Y. & Zamir, D. An introgression line population of *Lycopersicon pennellii* in the cultivated tomato enables the identification and fine mapping of yield-associated QTL. *Genetics* **141**, 1147–1162 (1995).
36. Gamuyao, R. et al. The protein kinase *Pstol1* from traditional rice confers tolerance of phosphorus deficiency. *Nature* **488**, 535–539 (2012).
37. Zhang, Z. et al. Genome-wide mapping of structural variations reveals a copy number variant that determines reproductive morphology in cucumber. *Plant Cell* **27**, 1595–1604 (2015).
38. Garrison, E. et al. Variation graph toolkit improves read mapping by representing genetic variation in the reference. *Nat. Biotechnol.* **36**, 875–879 (2018).
39. Ameer, A. Goodbye reference, hello genome graphs. *Nat. Biotechnol.* **37**, 866–868 (2019).
40. Zhu, G. et al. Rewiring of the fruit metabolome in tomato breeding. *Cell* **172**, 249–261.e12 (2018).
41. Darwin, S. C., Knapp, S. & Peralta, I. E. Taxonomy of tomatoes in the Galapagos islands: native and introduced species of *Solanum* section *Lycopersicon* (Solanaceae). *Syst. Biodivers.* **1**, 29–53 (2003).
42. Peralta, I. E., Knapp, S. & Spooner, D. M. New species of wild tomatoes (*Solanum* section *Lycopersicon*: Solanaceae) from Northern Peru. *Syst. Bot.* **30**, 424–434 (2005).
43. Bayer, P. E., Golicz, A. A., Scheben, A., Batley, J. & Edwards, D. Plant pan-genomes are the new reference. *Nat. Plants* **6**, 914–920 (2020).

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023

Methods

Plant materials

Briefly, eight wild species from section *Lycopersicon* (*S. galapagense*, *S. pimpinellifolium*, *S. chmielewskii*, *S. neorickii*, *S. corneliomulleri*, *S. peruvianum*, *S. chilense* and *S. habrochaites*), one wild species from section *Lycopersicon* (*S. lycopersicon*) and two domesticated tomatoes (*S. lycopersicon* var. *lycopersicum* cv. M82 and *S. lycopersicum* var. *cerasiforme*) were collected. All seedlings were planted in Anningqu field test station of Xinjiang Academy of Agricultural Sciences.

De novo genome assembly

Methods for library construction and sequencing are provided in the Supplementary Note. Contig-level assemblies for the 11 representative accessions were conducted using a pipeline based on Canu (v.1.5)^{16,44} with the following procedures: longer read reads were selected with the settings corOutCoverage = 35; raw read overlapping was detected using a highly sensitive overlapper MHAP⁴⁵ (v.2.1.2, parameter corMhapSensitivity = normal), and error correction was performed using the Falcon⁴⁶ sense method (option correctedErrorRate = 0.025); error-corrected reads were trimmed of unsupported bases and hairpin adapters to reach their longest supporting range with default parameters, and the draft assemblies were then generated using the top 80% longest trimmed reads. Finally, to ensure base accuracy of assembly results from SMRT molecules, we further polished the consensus genome sequences based on Illumina paired-end reads using Pilon⁴⁷ (v.1.22) with parameter: -mindepth 10 -fix bases.

Scaffolding using Bionano optical maps

For *S. galapagense*, we constructed Bionano optical maps. Young leaves were collected after two days of dark treatment. High-molecular-weight DNA was isolated and labeled with the restriction endonuclease Nb.BssSI, and labeled DNA was imaged with a Bionano Irys system. Molecules with lengths >150 kb, label SNR > 3.0 and average molecule intensity <0.6 were retained for scaffolding. These molecules were de novo assembled into genome maps using IrysSolve v.3.5_12162019 (<https://bionanogenomics.com/support/software-downloads/>). Pairwise comparison was first performed with RefAligner (<https://bionanogenomics.com/support/software-downloads/>) to identify overlaps among these molecules, and consensus maps were constructed. All molecules were then mapped back to the consensus maps and recursively refined and extended.

The Bionano IrysSolve module 'HybridScaffold' was used to perform hybrid assembly between the assembled contigs and genome maps. Assembled contigs were first converted into cmap format and then aligned to the contig cmaps with RefAligner, followed by label rescaling. The rescaled Bionano cmaps were aligned again to the contig cmaps, and sequences were split at the conflict points. Finally, scaffolds were built according to the alignment information. To improve the contiguity of assembly results, PBJelly⁴⁸ (v.15.8.24) was used to fill gaps using the error-corrected PacBio reads.

Pseudomolecule construction

The Hi-C data were mapped to the assemblies using BWA⁴⁹ (v.0.7.10-r789) with default parameters. Only uniquely aligned read pairs with mapping quality >20 were retained for further analysis. Duplicate removal, sorting and quality assessment were performed using HiC-Pro⁵⁰ (v.2.8.1) with default parameters. Only valid interaction pairs of Hi-C reads were fed into LACHESIS (v.1.0)⁵¹ for chromosome-scale scaffold construction. Briefly, contigs or scaffolds for each tomato assembly were broken into fragments with a length of 200 kb and then clustered using valid interaction read pairs by LACHESIS with the following parameters: 'CLUSTER_MIN_RE_SITES = 22, CLUSTER_MAX_LINK_DENSITY = 2, CLUSTER_NONINFORMATIVE_RATIO = 2, ORDER_MIN_N_RES_IN_TRUN = 10, ORDER_MIN_N_RES_IN_SHREDS = 10'. We manually checked the Hi-C interaction heat maps to identify potential

genomic regions containing assembled haplotigs due to heterozygosity, which were then excluded from the assembly. The manual curation step was reperformed several times, until the chromatin interaction signals reflecting putative haplotigs were undetectable.

Evaluation of genome assemblies

Completeness of the assembled tomato genomes was first assessed using BUSCO²⁴ (v.5.2.0) based on the embryophyta_odb9 database. We also assessed the mapping proportions of transcripts assembled with Trinity (v.2.8.5)⁵² to corresponding genome assemblies using BLASTN (v.2.12.0+)⁵³ with minimum alignment length of 300 bp and sequence identity >95%. These assemblies were also evaluated by mapping the Illumina short reads using BWA (default parameters).

Repeat sequence annotation

Both homology-based and de novo strategies were applied to identify repetitive sequences for all the tomato genomes. Four de novo prediction programs were applied: RepeatScout⁵⁴ (v.1.0.5), LTR-FINDER⁵⁵ (v.1.05), MITE-hunter (v.1.0)⁵⁶ and PILER-DF⁵⁷ (v.1.0). Results from these four programs were integrated into a repetitive sequence database, which was then merged with Repbase⁵⁸ (v.19.06) and classified into different categories by the PASTEClassifier.py script included in REPET⁵⁹ (v.2.5). Using this repeat database, repetitive sequences were identified by homolog searching using RepeatMasker⁶⁰ (v.4.0.5) with default parameters. We computed the genetic distance (K) between both ends of an intact LTR-RT using the distmat (default parameters) program in the EMBOSS package (v.6.6.0)⁶¹, and the insertion time (T)

then concatenated. We constructed a phylogenetic tree using phyML (v.3.3.20190909)⁷⁹ with parameters '-model JTT -fe -v 0.576 -a 0.886 -nclasses 4 -search SPR -t e'. The divergence time was estimated using the MCMCtree program in the PAML package⁸⁰ (v.4.7b). Three calibration points (*S. tuberosum* versus *S. lycopersicum* var. *cerasiforme*: 7.0–8.0 Ma; *S. lycopersicoides* versus *S. lycopersicum* var. *cerasiforme*: 2.0–2.7 Ma; and *S. pimpinellifolium* versus *S. lycopersicum* var. *cerasiforme*: 1.0–1.5 Ma)⁸¹ were used to constrain the divergence time.

Analyses of the super-pangenome

To identify homologous relationships among the genomes of 11 tomatoes assembled in this study, *S. lycopersicum* var. *lycopersicum* cv. Heinz 1706 and *S. pennellii*, the longest transcript of each predicted gene in each genome was chosen as a representative. To handle unannotated genes, a common issue during gene prediction, we aligned coding sequences of all predicted genes to each of the 13 tomato genomes using GMAP (v.2015-06-12)⁸². If a gene showed more than 80% alignment coverage and identity, and no gene was predicted within the aligned regions, it was considered to be an unannotated gene in the corresponding genome and was not regarded as 'missing' in the further analysis. An all-against-all comparison was then performed using BLASTP⁵³ (*E*

Genome-wide association studies

We selected the 321 tomato accessions that have been resequenced^{2,10} for GWAS. A total of 43,901,591 SNPs were identified using the GATK (v.4.1.4.1) pipeline⁹⁰ with the *S. galapagense* genome as the reference. Population structure was calculated by principal component analysis in PLINK (v.1.9.0b4.6)⁹¹ using 437,028 SNPs showing less linkage disequilibrium, which was extracted using PLINK with parameters ‘–indep-pairwise 50 5 0.1 (windows, step, r^2)’. The first five principal components were used as cofactors for population structure correction.

A total of 32 tomato flavor-related metabolite traits reported previously² and contents of 362 annotated metabolites from tomato fruits reported previously⁴⁰ were analyzed using EMMAX (v.20120210)⁹² with the default KN kinship, in which the selected principal components were used as cofactors. SNP-based and SV-based GWAS were performed using SNPs or SVs with minor allele frequency >0.01 and missing call rate <0.1. The genome-wide significance thresholds (7.58×10^{-7}) were determined using a uniform threshold of $1/n$, where n is the effective number of independent SNPs and SVs calculated using the Genetic type 1 Error Calculator (v.0.2)⁹³. Phenotypic variation explained (PVE) was calculated by the formula $PVE = [2 \times (\beta^2) \times MAF \times (1 - MAF)] / [2 \times (\beta^2) \times MAF(1-MAF) + ((s.e. \times (\beta))^2) \times 2 \times N \times MAF \times (1 - MAF)]$, where N represents sample size, s.e. is the standard error of the effect number of genetic variants, β is the effect number of genetic variants and MAF is the minor allele frequency of the target marker.

Reporting summary

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

Data availability

All assembled genome sequences and annotations are accessible through our database (<http://caastomato.biocloud.net>). Assemblies for the tomato genomes have also been deposited in the National Center for Biotechnology Information (NCBI) under BioProject accession number [PRJNA809001](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA809001). Raw PacBio, transcriptome and Hi-C sequencing reads have been deposited in the NCBI sequence read archive (<https://www.ncbi.nlm.nih.gov/sra/>) under BioProject accession number [PRJNA756391](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA756391). Tomato whole-genome sequencing data were downloaded from NCBI (BioProjects: [PRJNA259308](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA259308), [PRJNA353161](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA353161), [PRJNA454805](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA454805) and [PRJEB5235](https://www.ncbi.nlm.nih.gov/bioproject/PRJEB5235)). The RepBase database was downloaded from <https://www.girinst.org/server/RepBase/index.php>. Source data are provided with this paper.

Code availability

Custom scripts and codes used in this study are available at GitHub (<https://github.com/HongboDoll/TomatoSuperPanGenome>) and Zenodo (<https://doi.org/10.5281/zenodo.7396707>)⁹⁴.

References

44. Koren, S. et al. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res.* **27**, 722–736 (2017).
45. Berlin, K. et al. Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nat. Biotechnol.* **33**, 623–630 (2015).
46. Chin, C.-S. et al. Phased diploid genome assembly with single-molecule real-time sequencing. *Nat. Methods* **13**, 1050–1054 (2016).
47. Walker, B. J. et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS ONE* **9**, e112963 (2014).
48. English, A. C. et al. Mind the gap: upgrading genomes with Pacific Biosciences RS long-read sequencing technology. *PLoS ONE* **7**, e47768 (2012).

49. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
50. Servant, N. et al. HiC-Pro: an optimized and flexible Hi-C data analysis pipeline. *BMC Bioinformatics* **14**, 115 (2013).

72. Haas, B. J. et al. Automated eukaryotic gene structure annotation using EVIDENCEModeler and the Program to Assemble Spliced Alignments. *Genome Biol.* **9**, 1–22 (2008).
73. Kanehisa, M. & Goto, S. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Res.* **28**, 27–30 (2000).
74. Bairoch, A. & Apweiler, R. The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic Acids Res.* **28**, 45–48 (2000).
75. Pruitt, K. D., Tatusova, T. & Maglott, D. R. NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.* **35**, D61–D65 (2007).
76. Conesa, A. et al. Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics* **21**, 3674–3676 (2005).
77. Tang, H. et al. Screening synteny blocks in pairwise genome comparisons through integer programming. *BMC Bioinformatics* **12**, 102. (2011).
78. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
79. Guindon, S. et al. New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst. Biol.* **59**, 307–321 (2010).
80. Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**, 1586–1591 (2007).
81. Särkinen, T., Bohs, L., Olmstead, R. G. & Knapp, S. A phylogenetic framework for evolutionary study of the nightshades (Solanaceae): a dated 1000-tip tree. *BMC Evol. Biol.* **13**, 214 (2013).
82. Wu, T. D. & Watanabe, C. K. GMAP: a genomic mapping and alignment program for mRNA and EST sequences. *Bioinformatics* **21**, 1859–1875 (2005).
83. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome Biol.* **20**, 238 (2019).
84. Marçais, G. et al. MUMmer4: a fast and versatile genome alignment tool. *Bioinformatics* **36**, 1519–1521 (2020).

