

1 **Multiple reference genome sequences of hot pepper reveal the massive**
2 **evolution of plant disease resistance genes by retroduplication**
3

4 Seungill Kim¹, Jieun Park^{1,2}, Seon-In Yeom³, Yong-Min Kim⁴, Eunyoung Seo¹, Ki-Tae Kim⁵, Myung-
5 Shin Kim¹, Je Min Lee⁶, Kyeongchae Cheong^{2,5}, Ho-Sub Shin¹, Saet-Byul Kim¹, Koeun Han^{1,7},
6 Jundae Lee⁸, Minkyu Park⁹, Hyun-Ah Lee¹, Hye-Young Lee¹, Young-sill Lee¹, Soohyun Oh¹, Joo
7 Hyun Lee¹, Eunhye Choi¹, Eunbi Choi¹, So Eui Lee¹, Jongbum Jeon², Hyunbin Kim², Gobong Choi²,
8 Hyunjung Song², JunKi Lee¹, Sang-Choon Lee¹, Jin-Kyung Kwon^{1,7}, Hea-Young Lee^{1,7}, Namjin Koo⁴,
9 Yunji Hong⁴, Ryan W. Kim⁴, Won-Hee Kang³, Jin Hoe Huh¹, Byoung-Cheorl Kang^{1,7}, Tae-Jin Yang¹,
10 Yong-Hwan Lee^{2,5}, Jeffrey L. Bennetzen⁹ & Doil Choi¹

11
12 ¹Department of Plant Science, Plant Genomics and Breeding Institute, Research Institute for
13 Agriculture and Life Sciences, Seoul National University, Seoul 08826, Korea.

14 ²Interdisciplinary Program in Agricultural Genomics, Seoul National University, Seoul 08826, Korea.

15 ³Department of Agricultural Plant Science, Institute of Agriculture & Life Science, Gyeongsang
16 National University, Jinju 52828, Korea.

17 ⁴Korean Bioinformation Center, Korea Research Institute of Bioscience and Biotechnology, Daejeon
18 34141, Korea.

19 ⁵Department of Agricultural Biotechnology, Seoul National University, Seoul 08826, Korea.

20 ⁶Department of Horticultural Science, Kyungpook National University, Daegu 41566, Korea.

21 ⁷Vegetable Breeding Research Center, Seoul National University, Seoul 08826, Korea.

22 ⁸Department of Horticulture, Chonbuk National University, Jeonju 54896, Korea.

23 ⁹Department of Genetics, University of Georgia, Athens, GA 30602-7223, USA.

24

25

1 **Summary**

2 Transposable elements (TEs) provide major evolutionary forces leading to new genome structure and
3 species diversification. However, the role of TEs in the expansion of disease resistance gene families
4 has been unexplored in plants. Here, we report high-quality *de novo* genomes for two peppers
5 (*Capsicum baccatum* and *C. chinense*) and an improved reference genome (*C. annuum*). Dynamic
6 genome rearrangements involving translocations among chromosome 3, 5 and 9 were detected in
7 comparison between *C. baccatum* and the two other peppers. The amplification of *athila* LTR-
8 retrotransposons, members of the *gypsy* superfamily, led to genome expansion in *C. baccatum*. In-
9 depth genome-wide comparison of genes and repeats unveiled that the copy numbers of NLRs were
10 greatly increased by LTR-retrotransposon-mediated retroduplication. Moreover, retroduplicated NLRs
11 exhibited great abundance across the angiosperms, with most cases lineage-specific and thus recent
12 events. Our study revealed that retroduplication has played key roles in the emergence of new disease-
13 resistance genes in plants.

14

15

1 Introduction

2 Long terminal repeat retrotransposons (LTR-Rs) are a major evolutionary force in animals, fungi,
3 and, especially, plants. They comprise >75% of many plant genomes and cause genomic instability,
4 including genome expansion by amplification using an RNA intermediate¹. Besides genome
5 expansion, LTR-Rs facilitate the creation of new candidate genes called retrogenes by means of
6 retroduplication, in which spliced mRNA is captured, reverse transcribed, and subsequently integrated
7 into the genome by a retrotransposon²⁻⁴. In contrast to transduplication caused by class II transposable
8 elements (TEs)^{5,6}, the distinctive features of retrogenes are (i) intron loss compared to their parental
9 source genes, (ii) the presence of a 3' poly(A) tail, and (iii) flanking direct repeats⁷.

10 The evolutionary forces acting on most plant retrogenes are still largely unclear^{3,8-11}. Although LTR-
11 Rs are the most abundant TEs in all but the tiniest plant genomes, few studies have been reported on
12 the detection of retrogenes generated by LTR-Rs in plants^{6,12,13}. Wang *et al*³ identified 27 retrogenes
13 within LTR-Rs in rice and concluded that retrogenes that originated within LTR-Rs were often not
14 classified as retrogenes, partly because of the rapid destruction of the LTR-R structure by illegitimate
15 recombination¹⁴. Moreover, they suggested that the retrogenes might be very frequent in grass species
16 due to the abundance of LTR-Rs. In agreement with this prediction, recent studies have reported the
17 genome-wide identification of hundreds retrogenes within LTR-Rs in maize¹⁵, rice and sorghum¹⁶ as
18 well as the existence of retrogenes captured by LTR-Rs in *Arabidopsis*⁴. However, most of those
19 retrogenes were classified as pseudogenes or uncharacterized genes.

20 Previous studies reported the massive capture of specific gene families by certain TEs and
21 suggested a correlation between TE-mediated gene duplication and specific gene family expansion^{17,18}.
22 The nucleotide-binding and leucine-rich-repeat proteins (NLRs) represent a highly amplified gene
23 family in plants and provide the majority of functional plant disease resistance loci¹⁹⁻²¹. Comparative
24 genomic analyses have suggested the possibility of LTR-Rs and NLRs co-evolution, partly because
25 they are often co-localized^{20,22,23}. Because the NLRs usually reside in clusters within genomes, NLR
26 expansions have been mainly interpreted as the products of ectopic recombinational duplications^{19,24}.

27 Here, we report high-quality *de novo* genomes of two novel domesticated *Capsicum* species and

1 also improved the quality of the reference pepper genome²⁵. Comparative analyses of the three pepper
2 genomes, with other plant genomes as outgroups, provided insights into genome evolution and species
3 diversification in the genus *Capsicum*. Our analyses unveil an important mechanism for the massive
4 emergence of new plant NLRs by LTR-R-mediated retroduplication and show the dynamic
5 evolutionary processes for functional disease resistance genes in plants.

6

7 **Results**

8 ***De novo* sequencing, assembly and annotation of *Capsicum* genomes**

9 We sequenced and assembled the genomes of *Capsicum baccatum* PBC81 (hereafter, Baccatum)
10 and *C. chinense* PI159236 (hereafter, Chinense) using Illumina HiSeq 2500 with library insert sizes
11 ranging from 200 bp to 10 kb (Supplementary Tables 1-3). The estimated genome sizes of Baccatum
12 and Chinense, based on 19-mer analysis, were 3.9 and 3.2 Gb, respectively (Extended Data Fig. 1).
13 The assembled genomes of Baccatum and Chinense constituted 3.2 and 3.0 Gb (83 and 94% of the
14 estimated genome sizes, respectively) and had scaffold (contig) N50 sizes of 2.0 Mb (39 kb) and 3.2
15 Mb (49 kb), respectively (Supplementary Table 3). We annotated the protein-coding genes in the
16 Baccatum and Chinense assemblies as well as those in the pre-existing *C. annuum* CM334 genome²⁵
17 (hereafter, Annum) for detailed comparative analysis (Extended Data Fig. 2). On average, ~35,000
18 genes were annotated in each species (Supplementary Table 4). A comparison of the updated and
19 previous gene models of Annum revealed ~10,000 genes that did not overlap between the two
20 models, suggesting that most of the non-overlapping genes in the previous version were associated
21 with TEs (Extended Data Fig. 3).

22 A high-density genetic map of each species was generated by genotyping-by-sequencing on F2
23 populations^{26,27} (Extended Data Fig. 4). After breaking up chimeric scaffolds on the basis of genetic
24 map information, we organized the assembled genomes into 12 pseudo-chromosomes. Overall, 87%
25 of Baccatum (2.8 Gb in 2,083 scaffolds) and 89% of Chinense (2.8 Gb in 1,557 scaffolds) in
26 assembled genomes were ordered by the genetic map and inspected for syntenic inferences with the
27 improved pseudomolecules of Annum (Supplementary Table 5). We validated the assembled

1 genomes by reference guided mapping using the refined single-end and paired-end data, and
2 alignment to the assembled transcriptome of each species. In total, we detected more than 98.1% of
3 the filtered raw sequences (>98% identity) and more than 93.4% of the assembled transcriptomes
4 (>98% identity and 80% of query coverage) in the genome assemblies (Supplementary Table 6).
5 Taken together, our analyses provide the *de novo* reference genomes of two new pepper species as
6 well as an improved Annum genome.

7 Repeat annotation was performed with the assembled genomes and the initial contigs covering the
8 estimated genome sizes of the three species (Extended Data Fig. 5 and Supplementary Tables 7-8).
9 Overall, ~85% of the initial contigs were annotated as repeat sequences. LTR-Rs of the Ty3-*gypsy*
10 superfamily accounted for about half of the entire genome in each of the three species (Supplementary
11 Tables 7-8). Among the subgroups of the *gypsy* superfamily, *del* elements comprised the largest
12 fraction, representing 41.5, 34.9 and 41.7% (1,482, 1,337 and 1,343 Mb) in Annum, Baccatum and
13 Chinese, respectively. Furthermore, *athila* elements were more abundant (>2 fold) in Baccatum,
14 indicating that the *athila* subgroup contributed to species-specific genome expansion in the Baccatum
15 lineage (Supplementary Table 8).

16

17 **Speciation and evolution of the *Capsicum* species**

18 A phylogenetic analysis of the peppers with other plant species revealed that the divergence among
19 the three peppers occurred first between Baccatum and a progenitor of the other two peppers ~1.7
20 million years ago (MYA), followed by divergence between Annum and Chinese lineages ~1.1 MYA
21 (Fig. 1a and Supplementary Table 9). To identify genomic changes in the three pepper species, we
22 compared the genome structures, LTR-R insertion patterns, and gene duplication histories across these
23 pepper genomes (Figs. 1b-c and 2).

24 Chromosomal rearrangement is an important force in speciation, often producing unbalanced
25 gametes that reduce hybrid fertility²⁸. We performed an inter-genomic structural comparison and
26 detected translocations among the pepper genomes (Fig. 1b). The results show that chromosomes 3, 5
27 and 9 exhibit translocations that differentiate Baccatum from the other two species (Fig. 1b-c).

1 Collinearity comparisons among *Capsicum* species and two *Solanum* species revealed that the distal
2 region on the long arm of chromosome 9 was conserved in *Baccatum* but translocated to the short arm
3 of chromosome 3 in a shared ancestor of *Annuum* and *Chinense* (Fig. 1c and Extended Data Fig. 6).
4 Furthermore, chromosomes 6 and 4 of *Solanum* were detected in the terminal regions of the long and
5 short arms of chromosomes 3 and 5 in *Annuum* and *Chinense*, respectively. In contrast, the
6 orthologous regions of *Solanum* were mixed in the corresponding blocks of *Baccatum* (Fig. 1c and
7 Extended Data Fig. 6). This indicates that the distal regions of the long and short arms of
8 chromosomes 3 and 5 were translocated in the *Baccatum* lineage. We detected translocations between
9 the terminal regions of the short arm of chromosome 3 in *Baccatum* and the long arm of chromosome
10 9 in *Annuum* and *Chinense*. Consequently, our analyses revealed that translocations have generated
11 hetero karyotypes in both the *Baccatum* and the *Annuum/Chinense* progenitor lineages.

12 To compare LTR-R insertion patterns across the pepper genomes, we identified full length LTR-Rs
13 in each assembled genome and calculated their insertion times²⁹ (Extended Data Fig. 5 and
14 Supplementary Table 10). A peak of LTR-R activity in *Baccatum* appeared around its speciation time
15 1 to 2 MYA (Fig. 2a). Especially, the *athila* family was highly amplified in *Baccatum* around the
16 estimated speciation time, indicating that this subgroup may have explosively increased in *Baccatum*
17 after speciation. In *Chinense*, we observed the recent proliferation of LTR-retrotransposons (Fig. 2a).

18 Gene duplication is a major mechanism generating functional diversity between species by the
19 creation of new genes^{30,31}. We detected recent gene duplication events and characterised the
20 repertoires of duplicated genes in the three pepper genomes during and after speciation (Fig. 2b).
21 Overall, the duplication events were particularly frequent in the *Baccatum* lineage, both during and
22 after the speciation. In particular, NLRs were extensively amplified in *Baccatum* in the last 0-2 MYA,
23 and more recently in the other two peppers (Fig. 2b). Taken together, those results suggest that the
24 chromosomal rearrangements, accumulation of specific LTR-Rs, and differential gene duplications
25 have contributed enormously to genome diversification in the *Capsicum*.

26

27 **Massive creation of new NLRs via LTR-R-driven retroduplication**

1 Previous study suggested that NLRs were amplified in pepper compared to tomato and potato
2 genomes²¹. In particular, coiled-coil NLR subgroup 2 (CNL-G2) was highly expanded in the pepper
3 genome (Supplementary Table 11). To explore the possible mechanism of the NLR proliferation in
4 *Capsicum* spp., we analyzed the NLRs and their flanking sequences. We identified 157, 163 and 116
5 NLRs located inside LTR-Rs in Annum, Baccatum and Chinense, respectively (Supplementary
6 Tables 11-13). Hence, a large proportion (~18%) of the NLRs were amplified by LTR-Rs, with the
7 structures indicating that their retroduplicated origin is still intact. Most of these NLRs (~71%) were
8 in the CNL-G2 category, indicating that the copy number of specific NLR subfamilies was
9 particularly expanded (Fig. 3a). Furthermore, most of the retroduplicated NLRs (~71% of the total
10 and ~66% of the CNL-G2 type) were inside non-autonomous LTR-Rs that contained no *gag* or *pol*
11 protein coding potential (Supplementary Table 14). This suggests that all steps for the retroduplication,
12 presumably including the initial sequence capture process, had to be provided *in trans*. When we
13 compared the retroduplicated and other NLRs in the CNL-G2 category, the number and length of
14 exons were significantly fewer and longer in the retroduplicated NLRs, but not all of these were
15 single-exon genes (Fig. 3b-c). In total, ~44% of the retroduplicated NLRs in each species had multiple
16 exons but all of those had a reduced number of introns compared with their predicted parental
17 sequences (Supplementary Tables 13 and 15).

18 Earlier analyses of the tomato, potato and rice genomes indicated that retroduplication is a general
19 feature of genome evolution in the plant kingdom (Supplementary Table 11). We found that 22, 103
20 and 30 (8, 23, and 6%) of NLRs were inside LTR-Rs in tomato, potato and rice, respectively. Of these,
21 we identified parental sequences with multiple exons for 18, 88 and 22 of the NLRs inside LTR-Rs in
22 tomato, potato and rice, respectively, thus confirming their origin by retroduplication (Supplementary
23 Table 16). Similar to the peppers, NLRs in particular subgroups were primarily retroduplicated in
24 potato but the duplicated subgroups were distinct in each species (Supplementary Table 11). These
25 results indicate that LTR-Rs played a key role in the expansion of NLRs by retroduplication
26 throughout the plant kingdom, and that the detected events are both recent and lineage-specific.

27 In addition to the NLRs, we looked for other genes inside LTR-Rs in the six plant species

1 (Supplementary Table 17). In total, a range from 1,398 genes in rice to 3,898 genes in potato genomes
2 were found to be inside LTR-Rs, suggesting that 4 to 10% of all the genes in these plant species
3 emerged by LTR-R-driven retroduplication. On average, ~45% of them had functional domains
4 including highly amplified families such as MADS-box TFs, cytochrome P450s, and protein kinases,
5 and ~42% of those genes were expressed in one or more investigated tissues by RNA-Seq analysis
6 (Supplementary Table 17).

7

8 **Evolutionary mechanisms for the emergence of disease resistance genes in Solanaceae**

9 The *L* genes encoded by the NLRs are known to render resistance in peppers against
10 *Tobamoviruses* and they belong to the CNL-G4 category, along with *I2* in tomato that provides
11 resistance to race 2 of *Fusarium oxysporum* f. sp. *lycopersici* and *R3a* in potato that provides
12 resistance to the late blight pathogen, *Phytophthora infestans*³²⁻³⁴. Each of those is a single exon gene
13 encoding a peptide of ~1,300 amino acids. Synteny analysis and sequence comparison among pepper,
14 potato and tomato genomes suggested *L*, *I2*, and *R3a* are orthologous genes and the genomic regions
15 containing *L*, *I2* and *R3a* were tightly linked (Extended Data Fig. 7a and Supplementary Table 18).
16 These results suggest the possibility that the genes originated by an early retroduplication, and then
17 underwent divergent evolution in each lineage.

18 We examined the evolutionary history of *L* genes with their parents in the pepper genomes (Fig 4,
19 Extended Data Fig. 7b-c, and Supplementary Tables 19-20). The candidates of a parental gene (P1 to
20 P6) were identified considering similarity, *Ks* values (synonymous substitutions/synonymous site),
21 and alignment coverage to *L* genes. All candidate parental sequences contained multiple exons. When
22 candidate parental sequence P1 was compared with *L* in Annum, the results suggested that *L* was
23 derived from retroduplication in the ancestral lineage of *Capsicum* spp. ~8.9 MYA (Fig. 4). Because *L*
24 has 6.7 kb coding exons, with only an intron in the 3' UTR, and the presence of both flanking direct
25 repeat sequences and a poly(A) "tail" encoded in the DNA, our analysis suggests that *L* emerged
26 through capture and reverse transcription by a long interspersed nuclear element (LINE)-driven
27 retroduplication (Fig. 4 and Extended Data Fig. 8). Sequence comparison of *L* genes in the three

1 genomes and *L4* in *C. chacoense* revealed that the *L* genes were diversified by accumulation of
2 lineage-specific sequence mutations after speciation within *Capsicum* (Fig. 4 and Supplementary
3 Table 21). Consequently, our results suggest that the ancestor of the *L* genes was derived from
4 retroduplication and that subsequent divergent evolution has led to race-specific resistance against
5 diverse strains of *Tobamovirus* in each species of *Capsicum* after speciation (Fig. 4).

6 To analyse the evolutionary processes acting on *R3a* of potato, we first performed a genome-wide
7 search for the *R3a* as well as for candidate parental sequences. Because *R3a* is absent in the current
8 potato reference genome³⁵, we could not carry out accurate comparisons of *R3a* and their homologues.
9 However, *R3a* and its clustered genes originated from wild species, *Solanum demissum*³⁶ were
10 available in a public database. So, we compared these sequences with their closest homologs in the
11 reference potato genome³³. Our analyses revealed that intronless sequences of the ancestral potato *R3a*
12 might have emerged by RNA-based gene duplication in a shared ancestor of potato and tomato (Fig.
13 4). Subsequently, *R3a* and its paralogues were amplified by two rounds of tandem gene duplication
14 after the divergence of potato and tomato (Fig. 4 and Supplementary Table 22). Taken together, our
15 results suggest that retroduplication events are a main evolutionary process in the emergence of new
16 plant disease resistance genes, which can gain function via subsequent sequence variation and tandem
17 duplication.

18

19 **Evolution of potential anthracnose resistance genes in *Baccatum***

20 Pepper anthracnose caused by *Colletotrichum* spp. is one of the most devastating diseases in
21 worldwide pepper production³⁷. Due to the complexity of the interactions between the host and
22 *Colletotrichum* spp. and the lack of resistance in the *Annum* gene pool, a few *Baccatum* varieties
23 were identified as the only breeding resources for anthracnose resistance³⁸. Using pre-existing genetic
24 information³⁹, we identified the pertinent genomic regions and obtained 64 NLRs from a 3.8 Mb
25 region of *Baccatum* chromosome 3 as candidate resistance genes for *C. capsici* (Fig. 5 and
26 Supplementary Table 23). Previous studies reported that the main quantitative trait locus (QTL) for
27 pepper resistance against *C. capsici* was located on chromosome 9³⁹, however, we found that QTL is

1 located on chromosome 3 due to translocation in *Baccatum* and *Annuum* (Fig. 1c). We obtained 35
2 *Baccatum*-specific NLRs (27 in CNL-G2, 5 in CNL-G10 and 3 in CNL-G10) from the 64 NLRs by
3 sequence comparison among the three pepper genomes (Fig. 5). Considering the gene duplication
4 history, 15 of the 35 genes appear to have emerged after generation of the *Baccatum* lineage and all of
5 them belong to the CNL-G2 category. Transcriptome evidence indicated that 10 of those 15 genes are
6 expressed in one or more tissues (Fig. 5 and Supplementary Table 23). Furthermore, half of the 15
7 genes appear to have emerged by retroduplication (Fig. 5). Consequently, our results suggest that
8 retroduplication along with tandem and segmental duplications, has played a major role in the
9 emergence of anthracnose resistance genes in the *Baccatum* lineage.

10

11 **Conclusions**

12 In this study, we generated new and improved genome resources for three *Capsicum* species. Our
13 data provide an accurate and updated gene model of the pre-existing reference pepper genome based
14 on annotation with accumulated knowledge, highlighting the importance of genome improvement
15 after the completion of sequencing project. High-quality pseudo-chromosomes constructed from three
16 pepper genomes enabled precise comparisons of genome structures and evolutionary analyses,
17 providing new insights into interspecies diversification via genome rearrangements and lineage-
18 specific evolution of LTR-Rs and genes in the genus *Capsicum*. Furthermore, we found evidences of
19 the massive evolution of NLRs by LTR-mediated retroduplication in dicot Solanaceae and monocot
20 rice, suggesting that such phenomena are ubiquitous in angiosperms. Our results suggested that at
21 least 6 to 23% of plant NLRs were emerged by LTR-R-driven duplication (Supplementary Table 11).
22 A notable feature of this retroduplication is that distinct subfamilies of NLRs were highly
23 retroduplicated in different plant lineages. We unveiled the emergence of functional disease resistance
24 genes in the Solanaceae family by retroduplication and the subsequent neofunctionalisation of those
25 genes by dynamic evolutionary processes including lineage-specific sequence mutation and tandem
26 duplication. Our data suggest that a large proportion of all genes (~4 to 10%) in plant species might
27 have emerged by LTR-R-driven retroduplication. We observed the lineage-specific amplification of

1 specific gene families by LTR-Rs in various plant species, including such genes as those encoding
2 cytochrome P450s in potato and MADS-box TFs in *Baccharis* (Supplementary Table 17). Taken
3 together, our results provide new insights into the evolution of functional plant disease-resistance
4 genes that belong to the NLR family as well as other high copy number gene families that are present
5 in the plants.

6
7 **Online Content** The Methods and Supplementary Information with any associated references are
8 available in the online version of the paper.

9

10 **URLs**

11 GapCloser v1.12, http://soap.genomics.org.cn/down/GapCloser_release_2011.tar.gz

12 RepeatModeler, <http://www.repeatmasker.org/RepeatModeler.html>

13 Rice gene expression data, <http://rice.plantbiology.msu.edu/>

14 Potato gene expression data, http://solanaceae.plantbiology.msu.edu/pgsc_download.shtml

15

16 **Acknowledgements**

17 This work was supported by a grant from the Agricultural Genome Center of the Next Generation
18 Biogreen 21 Program of RDA (Project No. PJ01127501) and by a grant from the Ministry of Science,
19 ICT, and Future Planning (MSIP) of the Korean government through the National Research
20 Foundation (NRF-2015R1A2A1A01002327) to D.C.

21

22 **Author contributions**

23 D.C. conceived the project, designed the content, and organised the manuscript. S.K. performed data
24 generation/analysis and managed subprojects. S.-B.K., H.-A.L., and H.-Y.L. prepared the DNA and RNA
25 samples. S.K., K.C., K.H., and J.L. performed the *de novo* genome construction of two domesticated pepper
26 genomes and the improvement of pre-existing reference genome. S.K., M.-S.K., J.P., Y.-M.K., N.K., and R.W.K.
27 carried out the gene annotation. S.K., J.-K.L., S.-C.L., T.-J.Y., J.K., H.-Y.L., and B.-C.K. fulfilled the repeat
28 annotation and analyses. S.K., H.-S.S., J.J., J.H.H., and Y.-H.L. implemented the genome structure comparison
29 and evolutionary analyses of TEs and genes. S.K., E.S., J.P., Y.-S.L., S.O., J.H.L., E.C., E.C., S.E.L., G.C., H. S.,

1 and W.-H.K. performed the gene family analyses. S.K., E.S., and J.P. carried out the retroduplication analyses
2 for the NLRs and other gene families. S.K., H.K., H.-S.S., and Y.H. designed and visualized the figures. S.K.,
3 S.-I.Y., Y.-M.K., K.-T.K., J.-M.L., M.P., J.L.B., and D.C. wrote the manuscript.

4

5 **Author information**

6 The genome sequences of *C. chinense* and *C. baccatum* are deposited in the GenBank under the accessions
7 MCIT00000000 and MLFT00000000 (the versions described in this paper are version MCIT01000000 and
8 MLFT01000000). Further information, containing pseudomolecules and annotations is uploaded to our website
9 (<http://peppergenome.snu.ac.kr>). The authors declare no competing financial interests. Correspondence and
10 requests for materials should be addressed to D.C. (doil@snu.ac.kr)

11

12

1 **References**

- 2 1 Bennetzen, J. L. Patterns in grass genome evolution. *Curr. Opin. Plant Biol.* **10**, 176-181 (2007).
- 3 2 Kaessmann, H., Vinckenbosch, N. & Long, M. RNA-based gene duplication: mechanistic and
4 evolutionary insights. *Nat. Rev. Genet.* **10**, 19-31 (2009).
- 5 3 Wang, W. *et al.* High rate of chimeric gene origination by retroposition in plant genomes. *Plant Cell*
6 **18**, 1791-1802 (2006).
- 7 4 Zhu, Z., Tan, S., Zhang, Y. & Zhang, Y. E. LINE-1-like retrotransposons contribute to RNA-based
8 gene duplication in dicots. *Sci. Rep.* **6**, 24755 (2016).
- 9 5 Jiang, N., Bao, Z. R., Zhang, X. Y., Eddy, S. R. & Wessler, S. R. Pack-MULE transposable elements
10 mediate gene evolution in plants. *Nature* **431**, 569-573 (2004).
- 11 6 Morgante, M. *et al.* Gene duplication and exon shuffling by helitron-like transposons generate
12 intraspecies diversity in maize. *Nat. Genet.* **37**, 997-1002 (2005).
- 13 7 Ohshima, K. RNA-mediated gene duplication and retroposons: retrogenes, LINEs, SINEs, and
14 sequence specificity. *Int. J. Evol. Biol.* **2013**, 424726 (2013).
- 15 8 Lisch, D. How important are transposons for plant evolution? *Nat. Rev. Genet.* **14**, 49-61 (2013).
- 16 9 Sakai, H. *et al.* Retrogenes in rice (*Oryza sativa* L. ssp. *japonica*) exhibit correlated expression with
17 their source genes. *Genome Biol. Evol.* **3**, 1357-1368 (2011).
- 18 10 Zhang, C., Gschwend, A. R., Ouyang, Y. & Long, M. Evolution of gene structural complexity: an
19 alternative-splicing-based model accounts for intron-containing retrogenes. *Plant Physiol.* **165**, 412-423
20 (2014).
- 21 11 Zhu, Z., Zhang, Y. & Long, M. Extensive structural renovation of retrogenes in the evolution of the
22 *Populus* genome. *Plant Physiol.* **151**, 1943-1951 (2009).
- 23 12 Hawkins, J. S., Kim, H., Nason, J. D., Wing, R. A. & Wendel, J. F. Differential lineage-specific
24 amplification of transposable elements is responsible for genome size variation in *Gossypium*. *Genome Res.*
25 **16**, 1252-1261 (2006).
- 26 13 Jin, Y. K. & Bennetzen, J. L. Integration and nonrandom mutation of a plasma-membrane proton
27 atpase gene fragment within the *Bs1* retroelement of maize. *Plant Cell* **6**, 1177-1186 (1994).
- 28 14 Ma, J., Devos, K. M. & Bennetzen, J. L. Analyses of LTR-retrotransposon structures reveal recent and
29 rapid genomic DNA loss in rice. *Genome Res.* **14**, 860-869 (2004).
- 30 15 Baucom, R. S. *et al.* Exceptional diversity, non-random distribution, and rapid evolution of
31 retroelements in the B73 maize genome. *PLoS Genet.* **5**, e1000732 (2009).
- 32 16 Jiang, S. Y. & Ramachandran, S. Genome-wide survey and comparative analysis of LTR
33 retrotransposons and their captured genes in rice and sorghum. *PLoS ONE* **8**, e71118 (2013).
- 34 17 Hoen, D. R. *et al.* Transposon-mediated expansion and diversification of a family of *ULP*-like genes.
35 *Mol. Biol. Evol.* **23**, 1254-1268 (2006).
- 36 18 Kong, H. *et al.* Patterns of gene duplication in the plant *SKPI* gene family in angiosperms: evidence
37 for multiple mechanisms of rapid gene birth. *Plant J.* **50**, 873-885 (2007).
- 38 19 Guo, Y. L. *et al.* Genome-wide comparison of nucleotide-binding site-leucine-rich repeat-encoding
39 genes in *Arabidopsis*. *Plant Physiol.* **157**, 757-769 (2011).

- 1 20 Ratnaparkhe, M. B. *et al.* Comparative analysis of peanut NBS-LRR gene clusters suggests
2 evolutionary innovation among duplicated domains and erosion of gene microsynteny. *New Phytol.* **192**,
3 164-178 (2011).
- 4 21 Seo, E., Kim, S., Yeom, S. I. & Choi, D. Genome-wide comparative analyses reveal the dynamic
5 evolution of nucleotide-binding leucine-rich repeat gene family among Solanaceae plants. *Front. Plant Sci.* **7**,
6 1205 (2016).
- 7 22 Hayashi, K. & Yoshida, H. Refunctionalization of the ancient rice blast disease resistance gene *Pit* by
8 the recruitment of a retrotransposon as a promoter. *Plant J.* **57**, 413-425 (2009).
- 9 23 Kuykendall, D., Shao, J. & Trimmer, K. A nest of LTR retrotransposons adjacent the disease
10 resistance-priming gene *NPR1* in *Beta vulgaris* L. U.S. Hybrid H20. *Int. J. Plant Genomics* **2009**, 576742
11 (2009).
- 12 24 Leister, D. Tandem and segmental gene duplication and recombination in the evolution of plant
13 disease resistance gene. *Trends Genet.* **20**, 116-122 (2004).
- 14 25 Kim, S. *et al.* Genome sequence of the hot pepper provides insights into the evolution of pungency in
15 *Capsicum* species. *Nat. Genet.* **46**, 270-278 (2014).
- 16 26 Jeong, H. S., Jang, S., Han, K., Kwon, J. K. & Kang, B. C. Marker-assisted backcross breeding for
17 development of pepper varieties (*Capsicum annuum*) containing capsinoids. *Mol. Breed.* **35**, 226 (2015).
- 18 27 Lee, Y. R., Yoon, J. B. & Lee, J. A SNP-based genetic linkage map of *Capsicum baccatum* and its
19 comparison to the *Capsicum annuum* reference physical map. *Mol. Breed.* **36**, 61 (2016).
- 20 28 Rieseberg, L. H. Chromosomal rearrangements and speciation. *Trends Ecol. Evol. (Amst.)* **16**, 351-358
21 (2001).
- 22 29 SanMiguel, P., Gaut, B. S., Tikhonov, A., Nakajima, Y. & Bennetzen, J. L. The paleontology of
23 intergene retrotransposons of maize. *Nat. Genet.* **20**, 43-45 (1998).
- 24 30 Flagel, L. E. & Wendel, J. F. Gene duplication and evolutionary novelty in plants. *New Phytol.* **183**,
25 557-564 (2009).
- 26 31 Panchy, N., Lehti-Shiu, M. & Shiu, S. H. Evolution of gene duplication in plants. *Plant Physiol.* **171**,
27 2294-2316 (2016).
- 28 32 Simons, G. *et al.* Dissection of the *Fusarium I2* gene cluster in tomato reveals six homologs and one
29 active gene copy. *Plant cell* **10**, 1055-1068 (1998).
- 30 33 Huang, S. W. *et al.* Comparative genomics enabled the isolation of the *R3a* late blight resistance gene
31 in potato. *Plant J.* **42**, 251-261 (2005).
- 32 34 Tomita, R. *et al.* Genetic basis for the hierarchical interaction between *Tobamovirus* spp. and *L*
33 resistance gene alleles from different pepper species. *Mol. Plant Microbe Interact.* **24**, 108-117 (2011).
- 34 35 Potato Genome Sequencing Consortium. Genome sequence and analysis of the tuber crop potato.
35 *Nature* **475**, 189-195 (2011).
- 36 36 Huang, S. *et al.* The *R3* resistance to *Phytophthora infestans* in potato is conferred by two closely
37 linked R genes with distinct specificities. *Mol. Plant Microbe Interact.* **17**, 428-435 (2004).
- 38 37 Than, P. P., Prihastuti, H., Phoulivong, S., Taylor, P. W. & Hyde, K. D. Chilli anthracnose disease
39 caused by *Colletotrichum* species. *J. Zhejiang Univ. Sci. B* **9**, 764-778 (2008).
- 40 38 Mahasuk, P., Struss, D. & Mongkolporn, O. QTLs for resistance to anthracnose identified in two

- 1 *Capsicum* sources. *Mol. Breed.* **36**, 10 (2016).
- 2 39 Lee, J., Do, J. W. & Yoon, J. B. Development of STS markers linked to the major QTLs for resistance
3 to the pepper anthracnose caused by *Colletotrichum acutatum* and *C. capsici*. *Hortic. Environ. Biotechnol.*
4 **52**, 596-601 (2011).
- 5 40 Luo, R. *et al.* SOAPdenovo2: an empirically improved memory-efficient short-read *de novo* assembler.
6 *Gigascience* **1**, 18 (2012).
- 7 41 Boetzer, M., Henkel, C. V., Jansen, H. J., Butler, D. & Pirovano, W. Scaffolding pre-assembled
8 contigs using SSPACE. *Bioinformatics* **27**, 578-579 (2011).
- 9 42 Kajitani, R. *et al.* Efficient *de novo* assembly of highly heterozygous genomes from whole-genome
10 shotgun short reads. *Genome Res.* **24**, 1384-1395 (2014).
- 11 43 Ghosh, S. & Chan, C. K. Analysis of RNA-Seq data using TopHat and Cufflinks. *Methods Mol. Biol.*
12 **1374**, 339-361 (2016).
- 13 44 Kim, S. *et al.* Integrative structural annotation of *de novo* RNA-Seq provides an accurate reference
14 gene set of the enormous genome of the onion (*Allium cepa* L.). *DNA Res.* **22**, 19-27 (2015).
- 15 45 Pruitt, K. D., Tatusova, T., Brown, G. R. & Maglott, D. R. NCBI Reference Sequences (RefSeq):
16 current status, new features and genome annotation policy. *Nucleic Acids Res.* **40**, D130-D135 (2012).
- 17 46 Slater, G. S. & Birney, E. Automated generation of heuristics for biological sequence comparison.
18 *BMC Bioinformatics* **6**, 31 (2005).
- 19 47 Stanke, M., Tzvetkova, A. & Morgenstern, B. AUGUSTUS at EGASP: using EST, protein and
20 genomic alignments for improved gene prediction in the human genome. *Genome Biol.* **7**, S11-S18 (2006).
- 21 48 Haas, B. J. *et al.* Automated eukaryotic gene structure annotation using EVIDENCEModeler and the
22 program to assemble spliced alignments. *Genome Biol.* **9**, R7 (2008).
- 23 49 McDowall, J. & Hunter, S. InterPro protein classification. *Methods Mol. Biol.* **694**, 37-47 (2011).
- 24 50 Ellinghaus, D., Kurtz, S. & Willhoeft, U. LTRharvest, an efficient and flexible software for *de novo*
25 detection of LTR retrotransposons. *BMC Bioinformatics* **9**, 18 (2008).
- 26 51 Steinbiss, S., Willhoeft, U., Gremme, G. & Kurtz, S. Fine-grained annotation and classification of *de*
27 *nov*o predicted LTR retrotransposons. *Nucleic Acids Res.* **37**, 7002-7013 (2009).
- 28 52 Wang, Y. *et al.* MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and
29 collinearity. *Nucleic Acids Res.* **40**, e49 (2012).
- 30 53 Kim, S. B. *et al.* Divergent evolution of multiple virus-resistance genes from a progenitor in *Capsicum*
31 *spp.* *New Phytol.* doi:10.1111/nph.14177 (2016).
- 32 54 Li, L., Stoeckert, C. J., Jr. & Roos, D. S. OrthoMCL: identification of ortholog groups for eukaryotic
33 genomes. *Genome Res.* **13**, 2178-2189 (2003).
- 34 55 Loytynoja, A. Phylogeny-aware alignment with PRANK. *Methods Mol. Biol.* **1079**, 155-170 (2014).
- 35 56 Zhang, Z. *et al.* KaKs_Calculator: calculating Ka and Ks through model selection and model
36 averaging. *Genomics Proteomics Bioinformatics* **4**, 259-263 (2006).
- 37 57 Moniz de Sa, M. & Drouin, G. Phylogeny and substitution rates of angiosperm actin genes. *Mol. Biol.*
38 *Evol.* **13**, 1198-1212 (1996).
- 39 58 Drummond, A. J., Suchard, M. A., Xie, D. & Rambaut, A. Bayesian phylogenetics with BEAUti and
40 the BEAST 1.7. *Mol. Biol. Evol.* **29**, 1969-1973 (2012).

- 1 59 Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**, 1586-1591
2 (2007).
- 3 60 Baucom, R. S., Estill, J. C., Leebens-Mack, J. & Bennetzen, J. L. Natural selection on gene function
4 drives the evolution of LTR retrotransposon families in the rice genome. *Genome Res.* **19**, 243-254 (2009).
- 5 61 The Tomato Genome Consortium. The tomato genome sequence provides insights into fleshy fruit
6 evolution. *Nature* **485**, 635-641 (2012).
- 7 62 Goff, S. A. *et al.* A draft sequence of the rice genome (*Oryza sativa* L. ssp. *japonica*). *Science* **309**,
8 879-879 (2005).
- 9

1 **Methods**

2 **Genome assembly**

3 In total, 425.7 Gb and 526.7 Gb of the Chinese and Baccatum genome sequences were generated in the
4 Illumina HiSeq 2500 system (Supplementary Table 1). To remove unnecessary sequences for genome assembly,
5 preprocessing analysis was performed as previously described²⁵ (Supplementary Table 2). The *de novo* genome
6 assembly of each species was performed with SOAPdenovo2⁴⁰ using the filtered raw sequences with parameters
7 K=77 and K=81 for Chinese and Baccatum, respectively (Supplementary Table 3). The SSPACE software⁴¹
8 was employed for additional scaffolding (-x 0 -m 46 -k 10 -a 0.4 -p 1), and Gapcloser v.1.12 (See URLs) and
9 Platanus⁴² were implemented using default parameters to close gaps.

10

11 **Gene and repeat annotations**

12 Gene annotation was performed for the three pepper genomes as described in Extended Data Figure 2. To
13 annotate protein coding genes, we assembled transcripts using Tophat and Cufflinks⁴³ with the RNA-Seq reads
14 described in Supplementary Table 12 and in a previous study²⁵. The ISGAP pipeline⁴⁴ was used to extract
15 accurate coding sequences from the assembled transcripts. Plant refSeq⁴⁵ and the public protein databases for
16 *Arabidopsis* (TAIR 10), tomato (iTAG 2.3), potato (PGSC v3.4), and pepper (PGA v1.55) were used with
17 Exonerate v2.2.0⁴⁶ to align protein to the pepper genomes. *Ab initio* prediction was carried out with
18 AUGUSTUS⁴⁷ version 3.0 using an in-house training set consisting of full-length cDNA generated from
19 transcriptome analysis and by protein alignment. Consensus gene models were determined with EVM⁴⁸ and the
20 biological description of each gene model was assigned based on the Uniprot database and INTERPRO scan
21 v5.15-54.0⁴⁹.

22 Repeat sequences were annotated in the initial contigs representing the whole genome sizes and the
23 assembled genomes of the three peppers, as shown in Extended Data Figure 5. An integrated repeat library of
24 the three peppers was constructed using RepeatModeler (See URLs). Annotation of intact LTR-Rs was
25 performed using LTRHarvest⁵⁰ (-maxlenltr 2000 and -similar 80) and LTRDigest⁵¹. The subgroup of LTR-Rs in
26 the integrated library was classified by comparing their sequences to those of the intact LTR-Rs using BLASTN
27 (similarity >90%) (Extended Data Fig. 5).

28

29 **Comparison of genome structures**

30 To identify regions that were either conserved or translocated between *Capsicum* and *Solanum* species, we
31 performed collinear analysis with MCScanX⁵² using the gene models of the three peppers, and the tomato and
32 potato genomes described in Supplementary Table 9. We identified regions that were not translocated between
33 the tomato and potato genomes as conserved blocks in the *Solanum* species. The conserved blocks in the
34 *Solanum* species were then compared to the three pepper genomes. Blocks in the pepper genomes that
35 were conserved or translocated between the *Capsicum* and *Solanum* species were determined as shown in
36 Extended Data Figure 6. To investigate the translocated blocks in the three pepper genomes, we examined the
37 gene collinearity for syntenic blocks as shown in Extended Data Figure 6 and Figure 1c.

38

39 **Gene duplication history**

1 To estimate the gene duplication times of the annotated genes in the pepper genomes, we constructed a
2 computational pipeline by modifying a previously described method⁵³. We first performed gene clustering
3 analysis using OrthoMCL⁵⁴ to classify the gene family. We assumed that the genes in the same clusters were in
4 the same family and performed all-by-all alignments of the coding sequences within the clusters in each species
5 using PRANK⁵⁵. For each alignment result, the K_s values were calculated using KaKs Calculator,⁵⁶ and single-
6 linkage clustering for the K_s values was performed using the hclust function in the R package. The molecular
7 clock rate (r) was calculated to be 6.96×10^{-9} substitutions per synonymous site per year⁵⁷. The duplication time
8 was estimated using the formula, K_s value / $2r$.

9 10 **Estimation of divergence time**

11 To estimate the divergence times of the plant genomes, we identified 2,540 single copy genes in the rice,
12 *Arabidopsis* (TAIR10), grape (VvGDB v2.0), tomato (v2.3), and potato (PGSC v3.4) genomes and the three
13 pepper genomes using OrthoMCL clustering⁵⁴ (Supplementary Table 9). Multiple alignments of the single copy
14 genes from the eight genomes were implemented using PRANK⁵⁵ (-f=nexus -codon). The speciation times of
15 the eight plant species were calculated by phylogenetic analysis using the BEAST package⁵⁸.

16 17 **Evolutionary analyses of LTR-Rs**

18 For the intact LTR-Rs, we performed alignment of the sequences between the 5' and 3' LTRs using PRANK.
19 The DNA substitution rates (K) between the 5' and 3' LTRs were calculated using baseml in the PAML
20 package⁵⁹. The insertion times of the LTR-Rs were estimated using the formula, $K/2r$ ($r = 1.3 \times 10^{-8}$)⁶⁰.

21 22 **Identification of retroduplicated NLRs in the plant genomes**

23 To identify NLR genes inside LTR-Rs, we used the rice (MSU RGAP 7), potato (PGSC v3.4), and tomato (v2.3)
24 genomes with the three pepper genomes. We first identified NLRs using a previously constructed pipeline²¹ and
25 extracted the NLRs within putative LTR-Rs predicted by LTRHarvest (Supplementary Table 11). We then
26 compared those results with the repeat annotated by RepeatMasker, and if the NLRs inside LTR-Rs overlapped
27 with other TEs such as DNA transposons or Helitrons, we considered that the LTR-R predicted using
28 LTRHarvest was probably incorrect and was then removed. Because of rapid deletion of LTR-Rs and other
29 unselected DNA in all flowering plants¹⁻¹⁴, we performed an additional identification of NLRs inside LTR-Rs
30 using the annotated repeats including the partial LTR-Rs generated by RepeatMasker. We reasoned that if the
31 NLRs were fully contained within LTR-Rs annotated by RepeatMasker, the NLRs were retroduplicated. To
32 verify intron removal from the retroduplicated NLRs, we determined whether the candidate parental sequences
33 of the NLRs contained multiple exons by aligning the candidate parental sequences with the NLRs using
34 Exonerate⁴⁶, requiring >95% query coverage of the NLRs (Supplementary Table 13). To increase analysis
35 accuracy, we ignored unclear cases where no parental sequences having multiple exons were detected for NLRs
36 inside LTR-Rs.

37 To predict whole genes inside LTR-Rs in the six plant genomes, we performed genome-wide identification of
38 possible structure of LTR-Rs using LTRHarvest, taking into account rapid sequence change between the LTRs (-
39 -similar 75%, minltrlen 100). Like annotation of the NLRs inside LTR-Rs, we extracted genes within directly

1 repeated LTR regions as putatively retroduplicated genes. For the genes inside LTR-Rs, the number of expressed
2 genes in one or more tissues was counted using RNA-Seq data, as described in Supplementary Table 12 and in
3 previous analyses^{25,35,61,62} (See URLs).

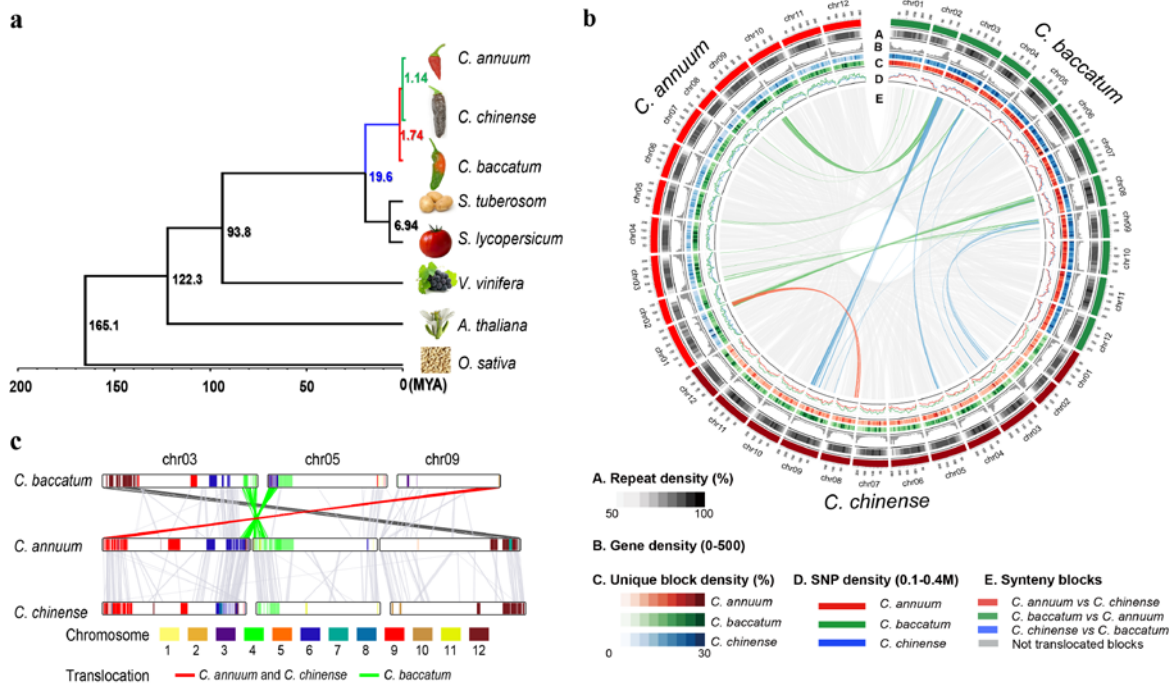
4 **Evolutionary investigation of functional disease-resistance genes in Solanaceae genomes**

6 The *L*, *I2* and *R3a* genes of pepper, tomato and potato were used to investigate evolutionary processes acting on
7 functional disease-resistance genes in the Solanaceae plants. The *L* genes in the *Capsicum* spp. were aligned to
8 paralogues in the pepper genomes using Exonerate⁴⁶ and closest homologs were identified (Supplementary
9 Tables 19-20). All of the closest homologs in each species were found to contain multiple exons and a gene that
10 we named P1 was identified as the most likely parental sequence. By comparison of the sequence divergence
11 between P1 and its closest homologs in the other genomes, we confirmed that P1 was Annum-specific
12 (Extended Data Fig. 7c). The 5' and 3' UTRs of *L1a* annotated based on RNA-Seq evidence were also compared
13 to the UTRs of P1 (Fig. 4).

14 For *R3a* in the potato, we aligned the coding sequences of *R3a* and genes within its cluster downloaded from
15 GenBank (AY849382, AY849383, AY849384 and AY849385) to the potato genome. Because of the absence of
16 those genes in the potato reference genome³⁵, we identified the closest homologs of *R3a* except *R3a* and its
17 clustered genes in the potato genome. The duplication time of the *R3a* family was estimated by comparison of
18 *R3a* and its homologs identified in the potato genome with the clustered genes. The *I2* sequence of tomato
19 downloaded from GenBank was also used to search in tomato reference genome, but *I2* was not found.

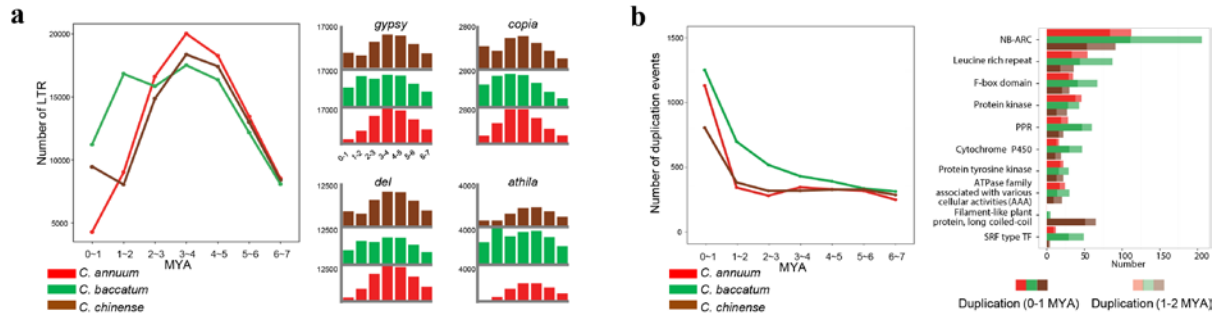
21 **Identification of potential anthracnose resistance genes**

22 To obtain candidate anthracnose resistance genes for *C. capsici*, we extracted NLRs located in the terminal
23 region of the short arm of chromosome 3 of *Baccatum* based on pre-existing genetic information³⁹
24 (Supplementary Table 23). Candidate genes that may provide resistance in *Baccatum* against *C. capsici* were
25 determined based on further inspected for the sequence conservation and the gene duplication time (Fig. 5).



1
2
3
4
5
6
7
8
9
10
11
12
13

Figure 1 | Lineage-divergence and genome structure comparisons of three *Capsicum* species. **a**, The reconstructed phylogenetic tree of eight plant genomes indicates their evolutionary relationships and estimated divergence times. **b**, The circular diagram shows the distribution of repeats, genes, genomic variations, and genome rearrangements in the pepper genomes. The subcategories indicate the density of repeats (A), genes (B), species-specific blocks (C), and SNPs (D) in the pepper genomes. The subcategory E depicts collinear and translocated blocks among the pepper genomes. **c**, A linear comparison of the rearranged blocks in the pepper genomes. Colours in the bars indicate translocated regions when comparing to tomato and potato genomes. The line colours indicate translocations in the ancestral lineage leading to Annuum and Chinense (red), in Baccatum (green) and in the ancestor of Annuum and Chinense or Baccatum (dark grey).

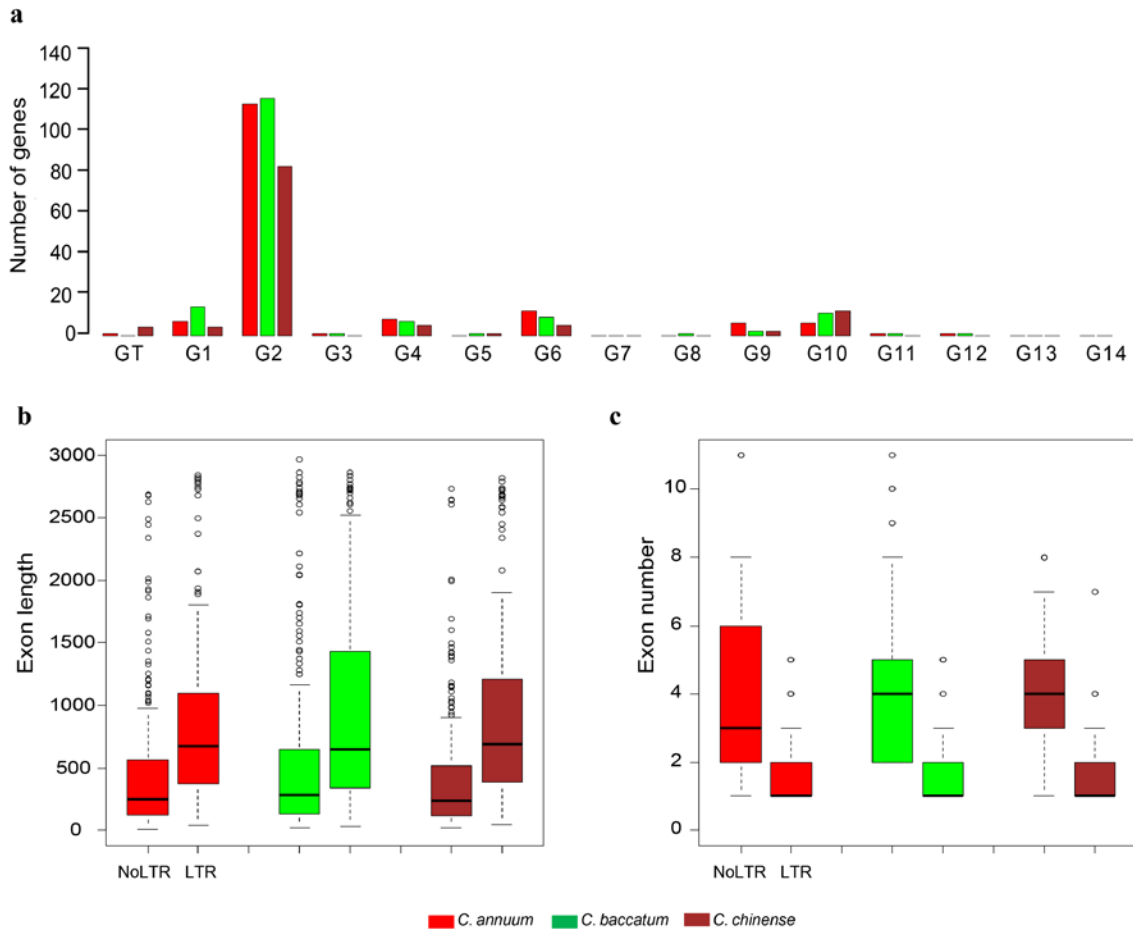


1
2

3 **Figure 2** | Evolutionary history of LTR-Rs and duplications of protein-coding genes in the pepper
4 genomes. **a**, Distribution of LTR-R insertions. The graphs in the left and right panels depict the
5 predicted insertion dates of LTR-R superfamilies (*gypsy*, *copia*) and two specific families (*del*, *athila*).
6 The *x*- and *y*-axes indicate the insertion times and the number of insertions at each time, respectively.
7 **b**, Time-scaled gene duplication history (left panel) and top 10 repertoires of massive gene duplication
8 (right panel). The *x*- and *y*-axes of the graph in the left panel indicate the approximate duplication time
9 (MYA) and the number of gene duplications, respectively. The *x*- and *y*-axes of the histogram in the
10 right panel represent the number of genes and domain description, respectively.

11

1



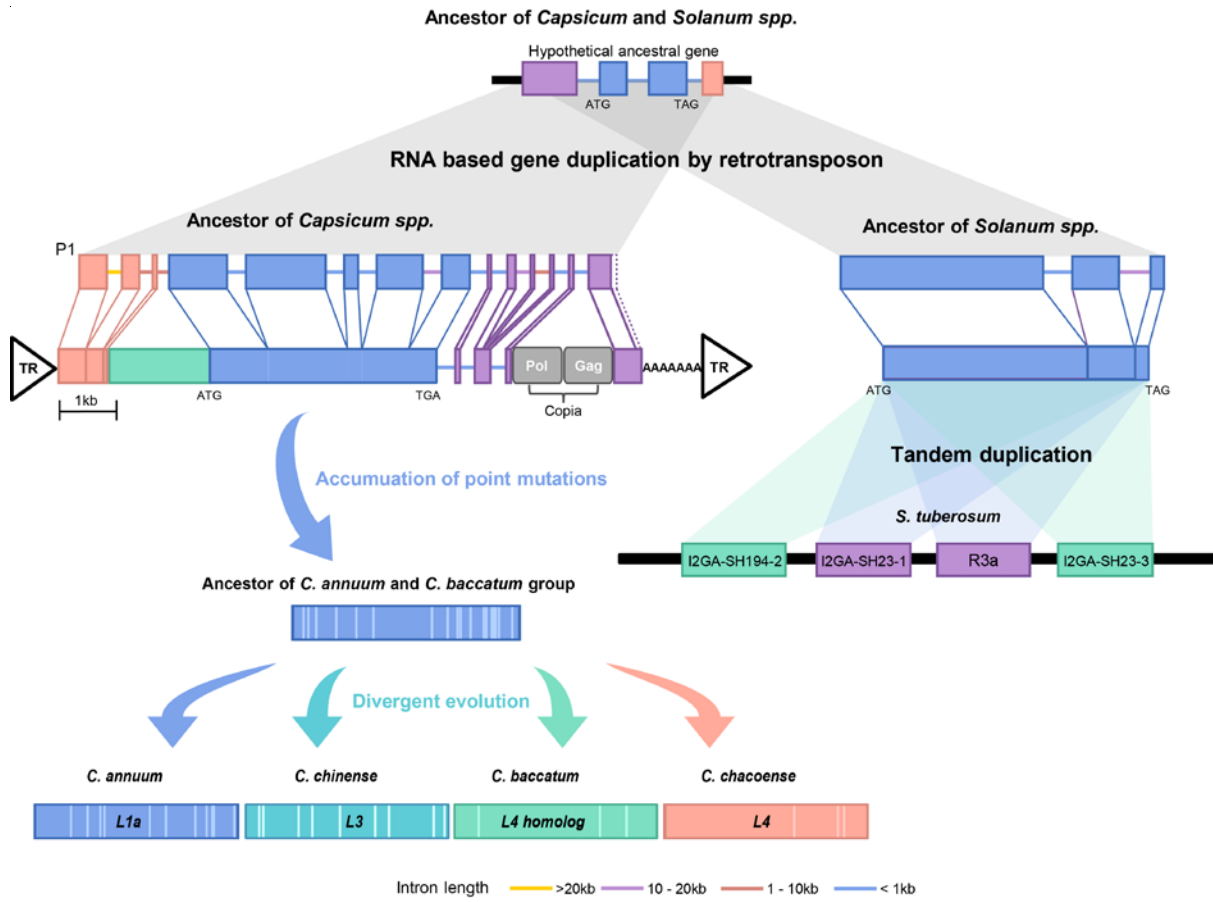
2

3

4 **Figure 3** | Emergence of large NLR gene families by retroduplication. The bar graph indicates the
5 number of retroduplicated NLRs in each subgroup. **a**, The bar graph indicates the number of
6 retroduplicated NLRs in each subgroup. The *x*- and *y*-axes indicate subgroups and the numbers of
7 genes, respectively. **b-c**, The exon lengths and the numbers of normal and retroduplicated NLRs are
8 depicted. **b**, The *x*- and *y*-axes indicate the normal and retroduplicated NLR groups and their exon
9 lengths, respectively. **c**, The *x*- and *y*-axes mean the groups of NLRs and the exon numbers,
10 respectively.

11

1

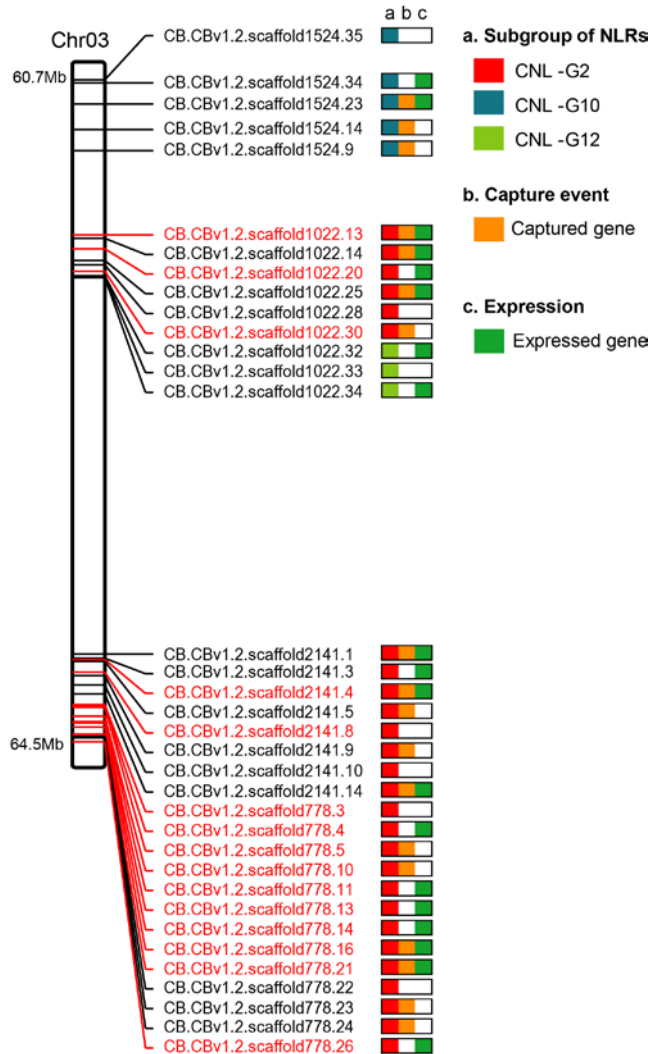


2

3

4 **Figure 4** | Emergence and evolution of *L* and *R3a* genes in the pepper and potato genomes. Models
 5 for the evolution of *L* and *R3a* in the pepper and potato are depicted. The gene names in the *R3a*
 6 cluster are from the previous analysis of Huang *et al*³³. The model proposes that *L* and *R3a* gene
 7 ancestors were first created by retroduplication, followed by the accumulation of point mutations and
 8 tandem duplication, respectively. DNA sequence indicative of a poly(A) tail and flanking terminal
 9 repeat (TR) sequences are depicted in the diagram as genomic evidence for a retroduplicated origin of
 10 *L*.

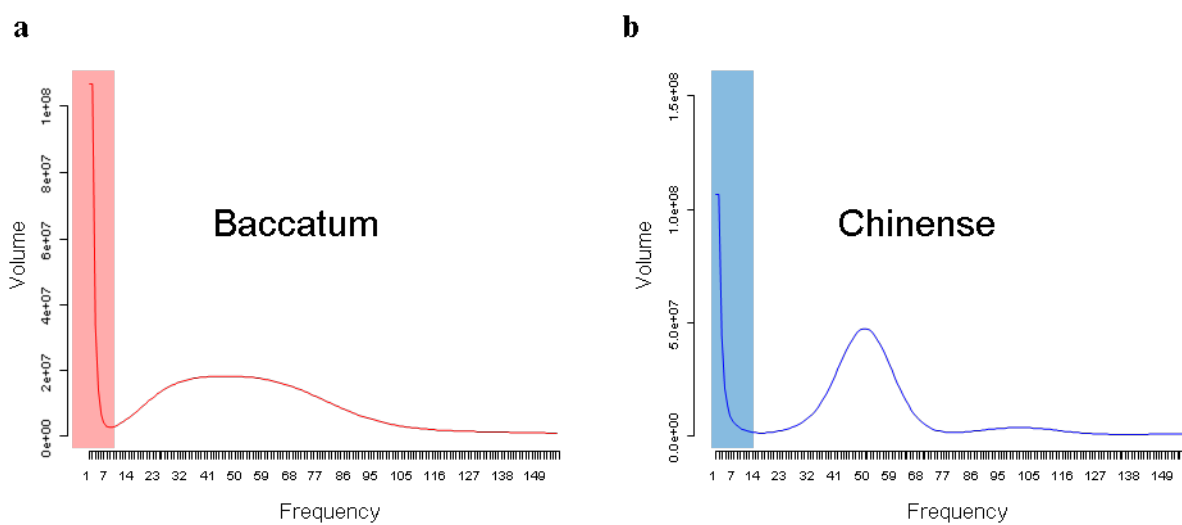
11



1
2

3 **Figure 5** | Potential *C. capsici* anthracnose resistance genes for in chromosome 3 of *Baccatum*.
 4 *Baccatum*-specific NLRs in the major QTL region are visualised on 3.8 Mb of chromosome 3. The
 5 chromosome plot shows the subgroups, proposed retroduplication events, and expression results for
 6 the NLRs. The black and red texts indicate NLR IDs that emerged before and after the speciation of
 7 *Baccatum*, respectively.

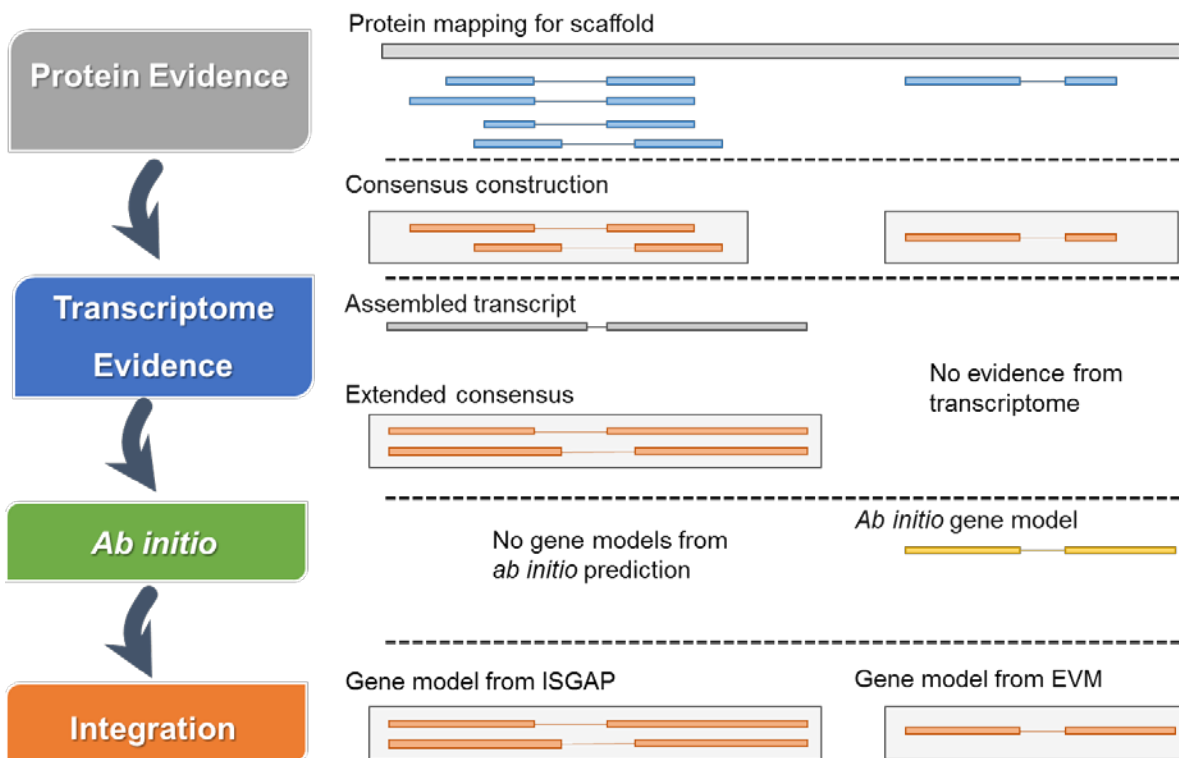
8



1
2
3
4
5

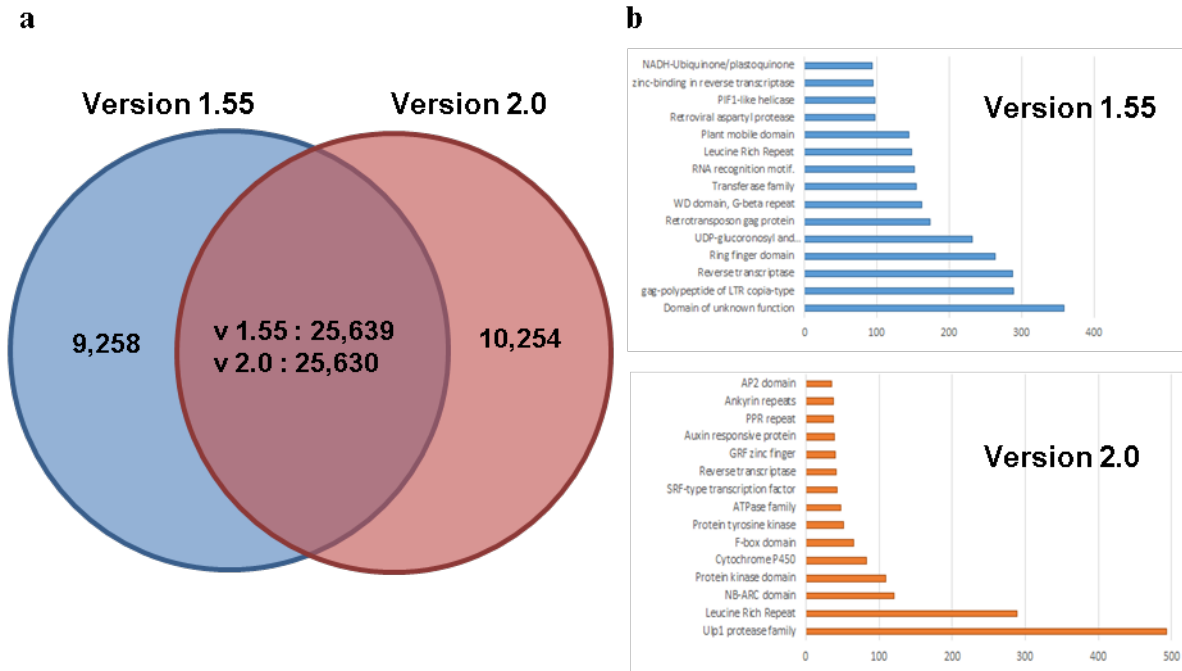
Extended Data Figure 1 | The 19-mer distribution patterns for the Baccatum and Chinense genome sequences. **a-b**, The *x*- and *y*-axes indicate the frequency and volume of 19-mers, respectively. Shaded regions in the graphs indicate low-frequency data as erroneous candidates.

1



2

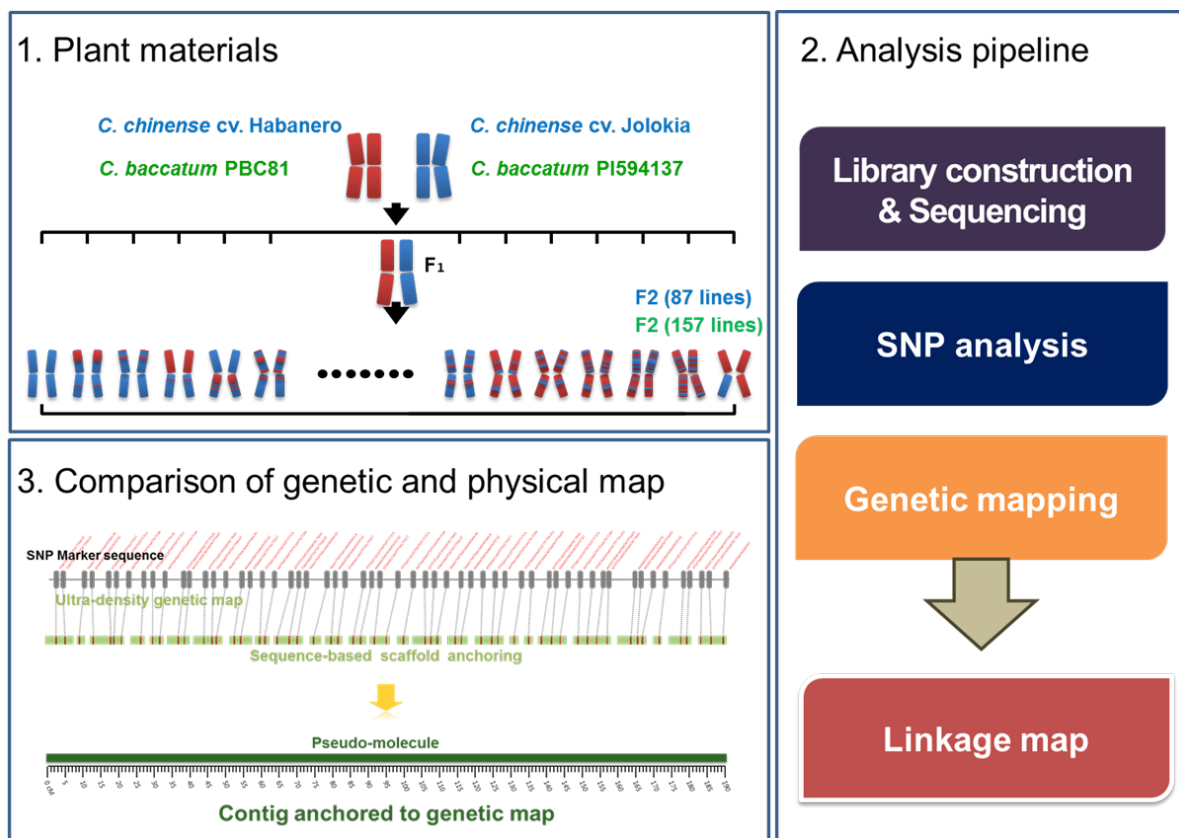
3 **Extended Data Figure 2** | Gene annotation scheme for the pepper genomes. The diagram depicts our
 4 gene annotation process for the pepper genomes using proteins, transcriptome and *ab initio* prediction.
 5 ISGAP⁴⁴ identifies gene structure based on transcriptome and protein evidences. EVM⁴⁸ integrates
 6 results of protein mapping and / or *ab initio* prediction using AUGUSTUS⁴⁷.



1
2
3
4
5
6
7
8
9

Extended Data Figure 3 | Comparison of annotated gene sets between the pre-existing and newly generated versions of the Annum genome. **a**, The Venn diagram indicates the numbers of overlapping and non-overlapping gene models between versions 1.55 and 2.0, using the genomic positions as one consideration. **b**, The bar graphs show the numbers of predicted genes in the top 15 domain descriptions in non-overlapping gene models in v1.55 and v2.0. The *x*- and *y*-axes indicate the number of annotated genes with a particular domain and the category of domain, respectively.

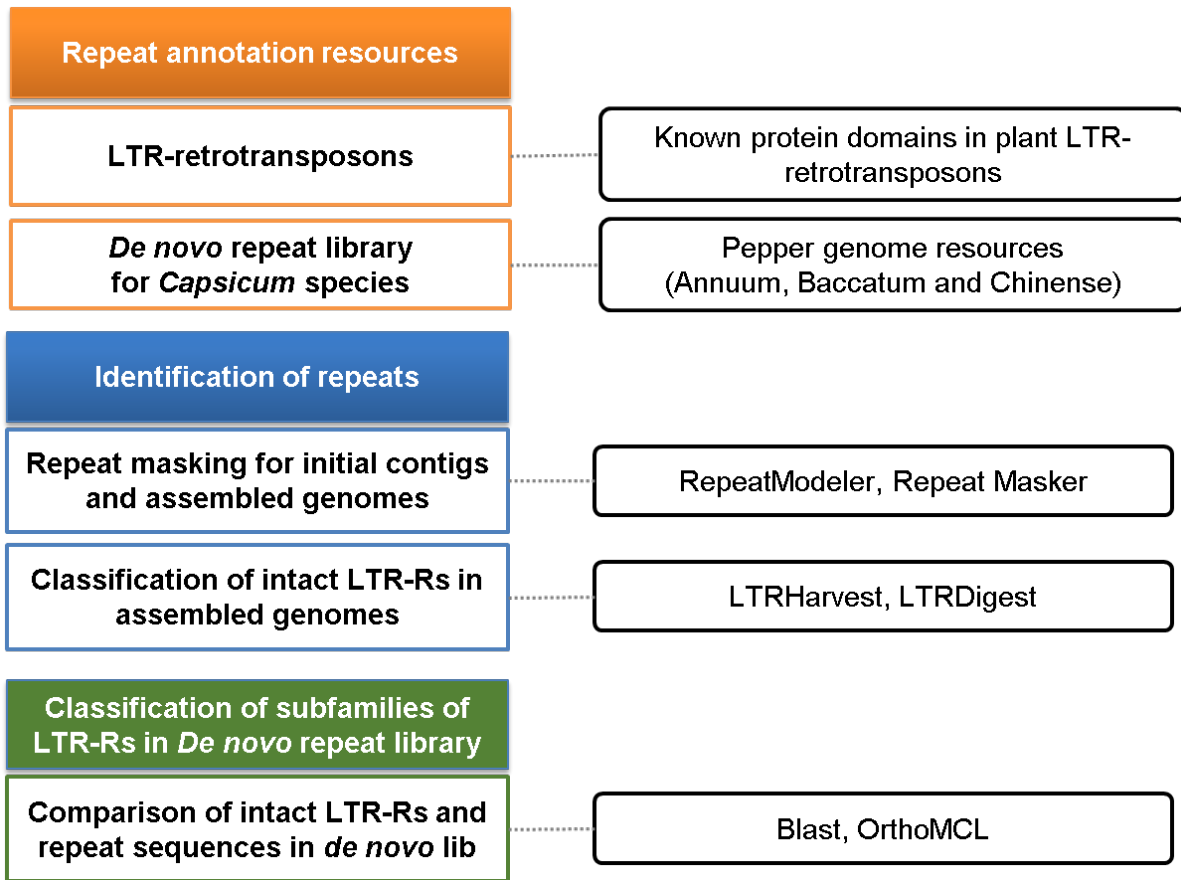
1



2 **Extended Data Figure 4** | Construction of genetic maps and pseudomolecule generation for the
3 pepper genomes. The plant materials used for construction of genetic maps of Baccatum and Chinese
4 are shown in the first diagram. The second diagram depicts the process of genetic map construction
5 for the two pepper genomes. The third diagram depicts the scaffold anchoring using the genetic maps
6 for the construction of pseudo-chromosomes for the two pepper genomes.

7

1

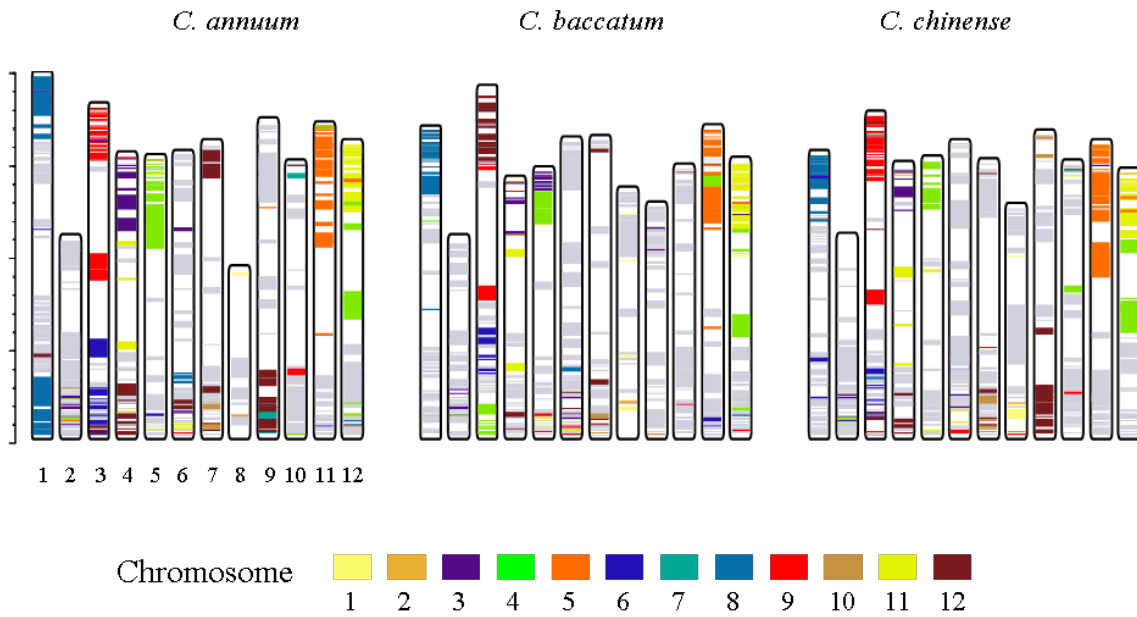


2

3 **Extended Data Figure 5** | The pipeline used for annotation of repeats in the pepper genomes. The
4 diagram depicts three steps in repeat annotation, i) construction of repeat libraries using repeat
5 annotation resources, ii) identification of repeats and intact LTR-Rs in the pepper genomes, and iii)
6 assignment of family information for LTR-Rs in the *de novo* repeat library. The diagrams in the right
7 panel indicate the resources and tools used for the annotation of repeats (See the Methods).

8

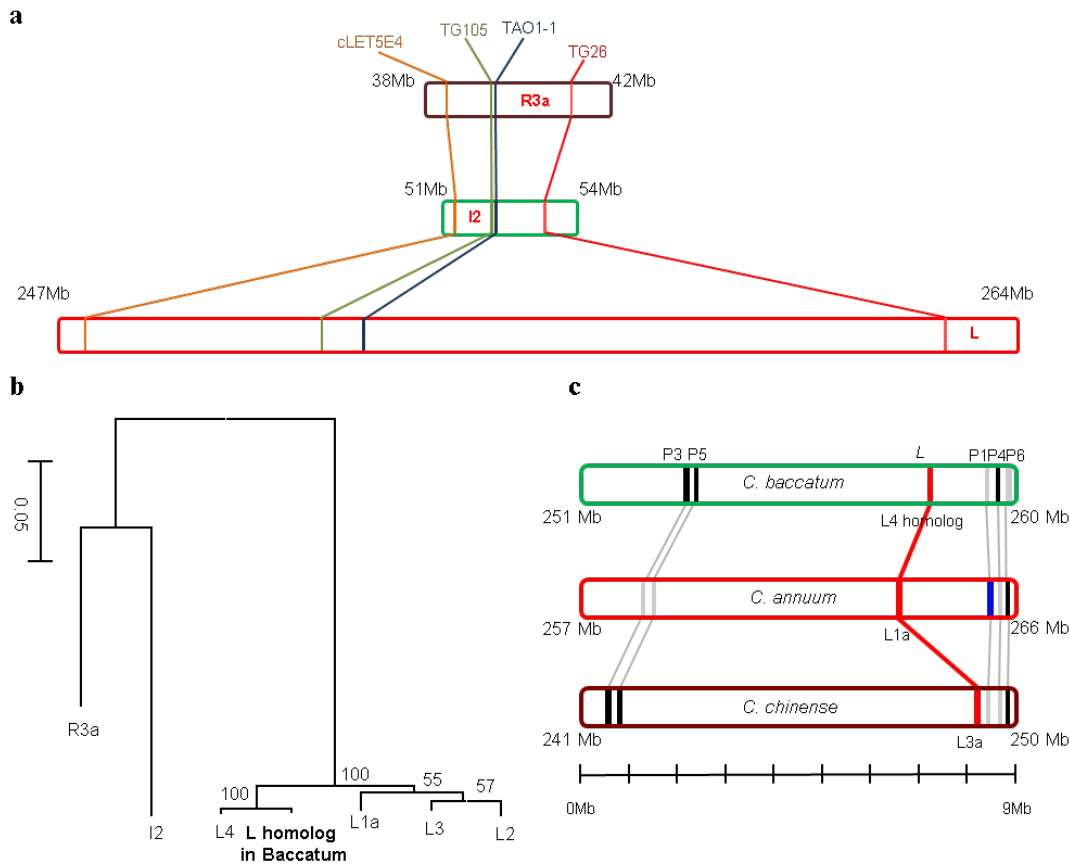
1



2

3 **Extended Data Figure 6** | The collinear blocks in pepper genomes exhibiting conserved syntenic
4 regions between tomato and potato genomes. The colours in the bars indicate collinear blocks in the
5 pepper genomes that are also conserved as blocks within the tomato and potato genomes. Chromatic
6 and grey colours indicate translocated and non-translocated regions, respectively, between the
7 *Capsicum* and *Solanum* genomes.

1



2

3 **Extended Data Figure 7 | Comparative analyses of *R3a*, *I2* and *L* genes and the locations of those**

4 genes in the Solanaceae genomes. **a**, Syntenic regions including *L*, *I2*, and *R3a* in the pepper, tomato

5 and potato genomes are depicted as bar graphs. The marker names were described in a previous

6 study³³. **b**, The phylogenetic tree of *L*, *I2* and *R3a*. **c**, Bar graphs depict the locations of the closest

7 homologs (P1 to P6) of the *L* genes in the pepper genomes as candidate parental sequences. Black and

8 grey lines in the bars mean the presence and absence of the parental sequences in pepper genomes,

9 respectively. The blue band within the red bar indicates the location of P1 in *Annuum*, the most likely

10 parental gene of *L1a*.

11

1



2

3 **Extended Data Figure 8** | Alignments between the 5' and 3' end direct repeat sequences flanking the
4 *L* genes in the pepper genomes.

5