

The *Arabidopsis lyrata* genome sequence and the basis of rapid genome size change

Tina T Hu^{1,15,16}, Pedro Pattyn^{2,3,16}, Erica G Bakker^{4-6,15}, Jun Cao⁷, Jan-Fang Cheng⁸, Richard M Clark^{7,15}, Noah Fahlgren^{5,9}, Jeffrey A Fawcett^{2,3,15}, Jane Grimwood^{8,10}, Heidrun Gundlach¹¹, Georg Haberer¹¹, Jesse D Hollister^{12,15}, Stephan Ossowski^{7,15}, Robert P Ottillar⁸, Asaf A Salamov⁸, Korbinian Schneeberger^{7,15}, Manuel Spannagl¹¹, Xi Wang^{11,15}, Liang Yang¹², Mikhail E Nasrallah¹³, Joy Bergelson⁴, James C Carrington^{5,9}, Brandon S Gaut¹², Jeremy Schmutz^{8,10}, Klaus F X Mayer¹¹, Yves Van de Peer^{2,3}, Igor V Grigoriev⁸, Magnus Nordborg^{1,14}, Detlef Weigel⁷ & Ya-Long Guo⁷

We report the 207-Mb genome sequence of the North American *Arabidopsis lyrata* strain MN47 based on 8.3× dideoxy sequence coverage. We predict 32,670 genes in this outcrossing species compared to the 27,025 genes in the selfing species *Arabidopsis thaliana*. The much smaller 125-Mb genome of *A. thaliana*, which diverged from *A. lyrata* 10 million years ago, likely constitutes the derived state for the family. We found evidence for DNA loss from large-scale rearrangements, but most of the difference in genome size can be attributed to hundreds of thousands of small deletions, mostly in noncoding DNA and transposons. Analysis of deletions and insertions still segregating in *A. thaliana* indicates that the process of DNA loss is ongoing, suggesting pervasive selection for a smaller genome. The high-quality reference genome sequence for *A. lyrata* will be an important resource for functional, evolutionary and ecological studies in the genus *Arabidopsis*.

Genome sizes in angiosperms range from 64 Mb in *Genlisea*¹ to an enormous 149 Gb in *Paris*²⁻⁴. Two major processes increase genome size: polyploidization and transposable elements proliferation. Processes that counteract genome expansion include the loss of entire chromosomes as well as deletion-biased mutations caused by unequal homologous recombination and illegitimate recombination⁵⁻⁹. Recent work comparing two cereals, rice and sorghum, has begun to shed light on some of these processes¹⁰. However, these species are

separated by 60 to 70 million years, making it difficult to disentangle the different evolutionary forces at work.

An exciting opportunity to understand what drives differences in genome size over shorter time scales is offered by the genus *Arabidopsis* in the Brassicaceae family. The genome of the self-incompatible perennial *A. lyrata* is larger than 200 Mb, which is near the family average^{11,12}, whereas the self-compatible annual *A. thaliana* has one of the smallest angiosperm genomes at about 125 Mb, even though the two species diverged only about 10 million years ago¹³⁻¹⁵. Compared to the difference between the two species, there is much less variation within *A. thaliana*¹¹.

A high-quality genome sequence for the partially inbred *A. lyrata* strain MN47 was assembled from approximately 8.3× coverage of dideoxy sequencing reads, making use of information from genetic maps and chromosome painting¹⁶⁻¹⁹ (Online Methods). The final assembly included 206.7 Mb of sequence, 90% of which is included in eight large scaffolds covering the majority of each of the eight chromosomes and another large scaffold of 1.9 Mb representing one of the centromeres. Based on cytological observations²⁰, the centromeric gaps were estimated to span 17.2 Mb. A combination of *de novo* predictions, homology to *A. thaliana* features and RNA sequencing was used to annotate the genome. In *A. lyrata*, we predicted 32,670 protein-coding genes compared to the 27,025 genes in *A. thaliana*²¹.

Because overall sequence identity between *A. lyrata* and *A. thaliana* is greater than 80% (Supplementary Fig. 1), the two genomes could

¹Molecular and Computational Biology, University of Southern California, Los Angeles, California, USA. ²Department of Plant Systems Biology, VIB, Gent, Belgium. ³Department of Plant Biotechnology and Genetics, Ghent University, Ghent, Belgium. ⁴Department of Ecology and Evolution, University of Chicago, Chicago, Illinois, USA. ⁵Center for Genome Research and Biocomputing, Oregon State University, Corvallis, Oregon, USA. ⁶Department of Horticulture, Oregon State University, Corvallis, Oregon, USA. ⁷Department of Molecular Biology, Max Planck Institute for Developmental Biology, Tübingen, Germany. ⁸US Department of Energy Joint Genome Institute, Walnut Creek, California, USA. ⁹Department of Botany and Plant Pathology, Oregon State University, Corvallis, Oregon, USA. ¹⁰HudsonAlpha Genome Sequencing Center, HudsonAlpha Institute for Biotechnology, Huntsville, Alabama, USA. ¹¹Munich Information Center for Protein Sequences/Institute for Bioinformatics and Systems Biology, Helmholtz Center Munich, Neuherberg, Germany. ¹²Department of Ecology and Evolutionary Biology, University of California Irvine, Irvine, California, USA. ¹³Department of Plant Biology, Cornell University, Ithaca, New York, USA. ¹⁴Gregor Mendel Institute, Austrian Academy of Science, Vienna, Austria. ¹⁵Present addresses: Lewis-Sigler Institute for Integrative Genomics, Princeton University, Princeton, New Jersey 08544, USA (T.T.H.), Dow AgroSciences, Portland, Oregon 97224, USA (E.G.B.), Department of Biology, University of Utah, Salt Lake City, Utah, USA (R.M.C.), Graduate University for Advanced Studies, Hayama, Kanagawa, Japan (J.A.F.), Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, Massachusetts, USA (J.D.H.), Center for Genomic Regulation, Barcelona, Spain (S.O.), Department of Plant Developmental Biology, Max Planck Institute for Plant Breeding Research, Cologne, Germany (K.S.) and Department of Molecular Biology, Max Planck Institute for Developmental Biology, Tübingen, Germany (X.W.). ¹⁶These authors contributed equally to this work. Correspondence should be addressed to D.W. (weigel@weigelworld.org) or Y.-L.G. (ya-long.guo@hotmail.com).

Received 13 September 2010; accepted 18 March 2011; published online 10 April 2011; doi:10.1038/ng.807

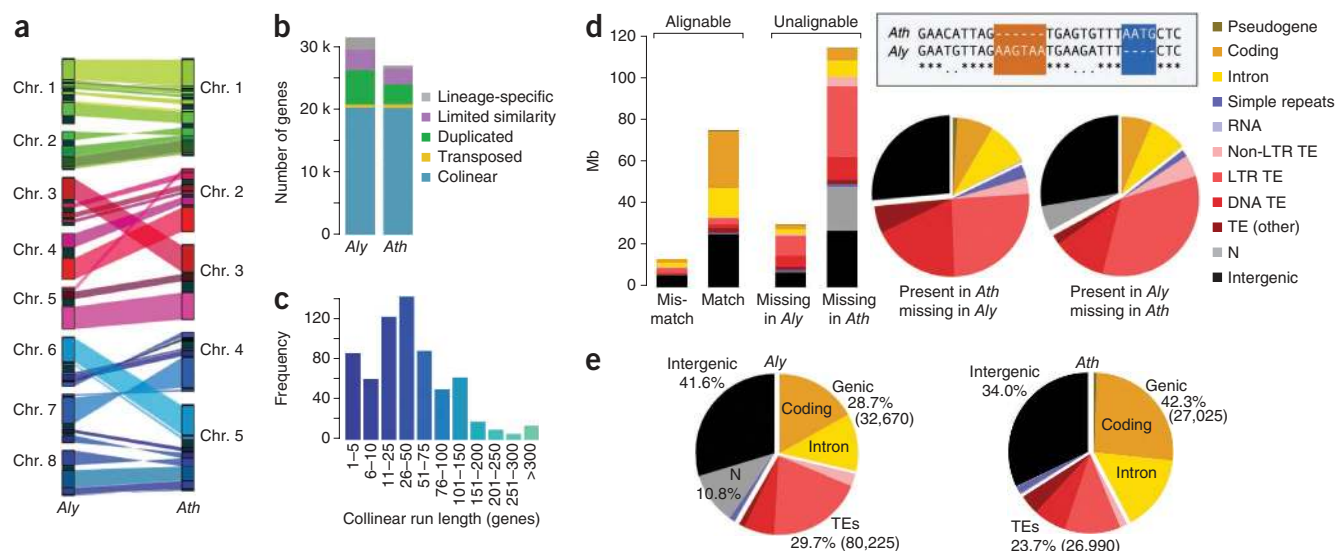


Figure 1 Comparison of *A. lyrata* and *A. thaliana* genomes. **(a)** Alignment of *A. lyrata* (Aly) and *A. thaliana* (Ath) chromosomes. Genomes are scaled to equal size. Only syntentic blocks of at least 500 kb are connected. **(b)** Orthology classification of genes. The mode at 1–5 reflects frequent single-gene transpositions. **(c)** Distribution of run lengths of collinear genes. The mode at 1–5 reflects frequent single-gene transpositions. **(d)** Unalignable sites can be considered as present in one species and absent in the other, as shown in the boxed sequence diagram; matches are indicated by asterisks, and mismatches are indicated by periods. The histogram on the left indicates the absolute number of unalignable sites, and the pie charts in the middle compare their relative distribution over different genomic features. See also **Supplementary Table 3**. **(e)** Genome composition (number of elements in parentheses). TEs, transposable elements.

be easily aligned (**Fig. 1a**). Genetic mapping^{16,18,19} revealed ten major rearrangements, including two reciprocal translocations and three chromosomal fusions, that led to the *A. thaliana* karyotype of five chromosomes, as compared to the ancestral state of eight chromosomes, as found in *A. lyrata* and other Brassicaceae. Although centromeric regions are difficult to assemble, we could identify the syntentic region in *A. thaliana* that corresponds to the chromosome 4 centromere of *A. lyrata*. The entire centromere has been lost, with only two remnants of satellite repeats in the 1.4-kb intergenic region between AT2G26570 and AT2G26580 (**Supplementary Fig. 2**).

Apart from chromosomal-scale changes, approximately 90% of the two genomes have remained syntentic, with the great majority being in highly conserved collinear arrangements (**Fig. 1b**). The run-length distribution of collinear gene pairs is bimodal, with a first peak of fragments of five or fewer collinear gene pairs (**Fig. 1c**), reflecting an abundance of small-scale rearrangements (<10 kb), including single-gene transpositions. Windows containing a breakpoint in collinearity are enriched for transposable elements and other repeats (**Supplementary Table 1**), in agreement with repetitive elements often being associated with chromosomal rearrangements and transposed genes^{22–27}, although they might not necessarily be causal²⁸. Two thirds of the 154 inversions identified between the two species are flanked by inverted repeats (**Supplementary Table 2**).

Despite this overall similarity in gene arrangement, the two genomes are strikingly different in size. A whole-genome alignment revealed that more than 50% (~114 Mb) of the *A. lyrata* genome appears to be missing from the *A. thaliana* reference genome. In contrast, only about 25% (~30 Mb) of the *A. thaliana* genome is absent from *A. lyrata* (**Fig. 1d**, **Supplementary Fig. 1e** and **Supplementary Table 3**). Nevertheless, the distribution across different sequence classes is similar: half of the unalignable sequences are in transposable elements and a quarter are in intergenic regions. The net effect of these changes is that the *A. thaliana* genome is ~80 Mb smaller than the *A. lyrata* genome, with a much higher fraction of genic sequences (42% instead of 29%), even though the total gene count is smaller (**Fig. 1e**).

The apparent shrinkage of the *A. thaliana* genome is not simply because of a few chromosome-scale changes: only 10% of the size difference is attributable to the three missing centromeres, and the rest is due to hundreds of thousands of smaller insertions and deletions

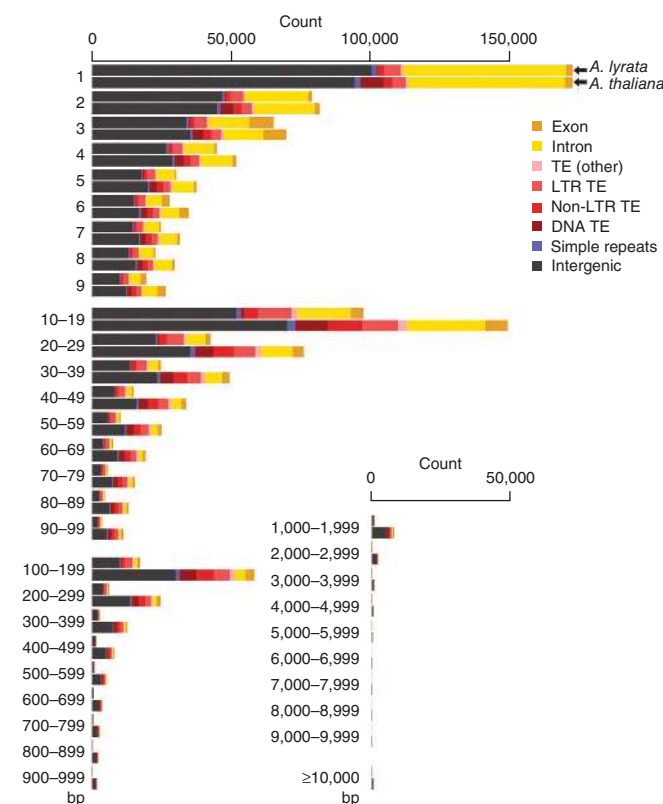


Figure 2 Apparent deletions by size and annotation. *A. lyrata* (Aly) is always shown on top, and *A. thaliana* (Ath) is always shown on the bottom. TE, transposable elements.

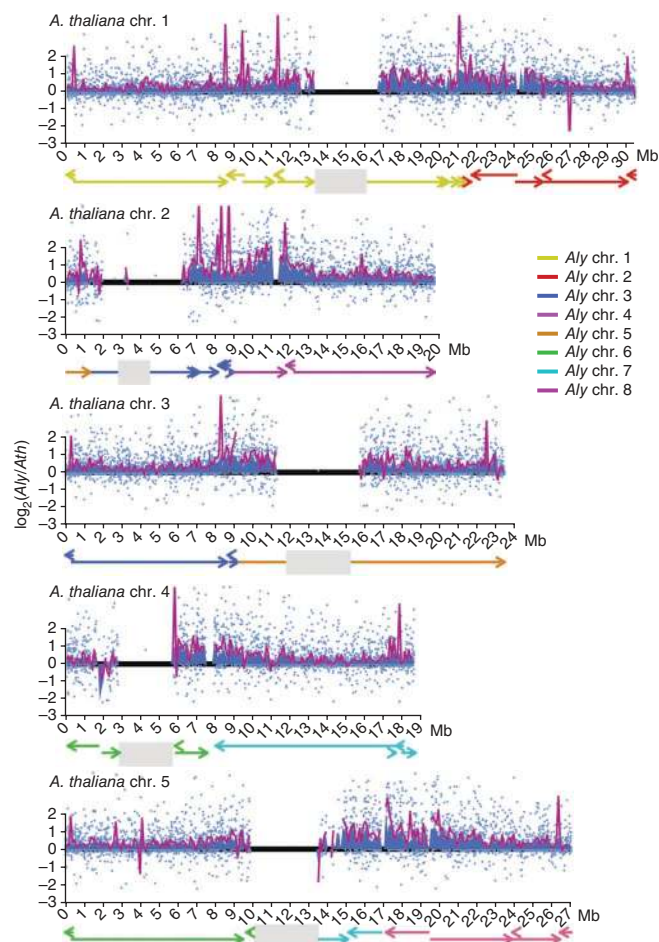
Figure 3 Changes in genomic intervals along the *A. thaliana* genome. Mean ratios for all collinear gene pairs in each 100-kb window are shaded in blue, with individual values shown as light blue dots. The ratio of the absolute length of each non-overlapping 100-kb window is shown as a dark purple line. Centromeres are indicated as gray boxes. *Aly*, *A. lyrata*.

spanning all classes of sites. Notably, whereas large differences much more often correspond to sequences only found in *A. lyrata*, this is not true for very small insertions and deletions (Fig. 2). This is in stark contrast to genomes from other closely related species with similarly sized genomes, such as chimpanzee and human²⁹.

Although rearrangements are correlated with genome shrinkage (rearranged regions are on average shorter in *A. thaliana* than are collinear regions; Fig. 3 and Fig. 4a), unalignable sequences are found throughout the genome. An analysis of collinear gene pairs confirmed that in most cases, intergenic intervals in *A. lyrata* are longer than their counterparts in *A. thaliana* (Fig. 4b). Introns behave similarly, although the difference in length is smaller¹³.

The gene content of *A. thaliana* is ~17% lower than that of *A. lyrata*, but there are no major differences in their Gene Ontology distributions. Similarly, divergence patterns for different gene families between the two species mirror those of within-*A. thaliana* polymorphism levels^{30,31}. The combined gene sets of *A. lyrata* and *A. thaliana* result in 12,951 clusters based on the Markov Cluster Algorithm³² (MCL), with fewer singletons present in *A. thaliana* (Fig. 4c). Among the 8,794 shared multi-gene MCL clusters (Fig. 4d), clusters that are smaller in *A. thaliana* outnumber those that are smaller in *A. lyrata* (1,797 to 612). F-box and NB-LRR genes are examples of gene families with particularly high birth and death rates in plants^{30,31,33–35}. *A. lyrata* has 596 F-box and 187 NB-LRR genes, compared to 502 and 159 genes, respectively, in *A. thaliana*. The trend of fewer genes in *A. thaliana* is supported by a broader comparison of the *Arabidopsis* gene set with those of two other dicots^{36–38}. *A. lyrata* has 114 ortholog clusters³⁹ shared with poplar and grapevine but not with *A. thaliana*, whereas *A. thaliana* has only 45 clusters found in poplar and grapevine but not in *A. lyrata*. Similarly, *A. lyrata* has 875 clusters not detected in any of the other three species, whereas *A. thaliana* has only 156 species-specific clusters (Supplementary Table 4 and Supplementary Fig. 3).

As in other taxa, transposable elements make an important contribution to the change in genome size (Fig. 1d), and transposable



elements comprise a larger fraction of the *A. lyrata* genome (Fig. 1e). Without an outgroup, one cannot infer directly how much such patterns are shaped by different transposable element activity levels or the differential purging of ancestral transposable elements since speciation. To obtain an estimate of relative activity levels, one can exploit the molecular clock to estimate the average age of long terminal repeat (LTR) retrotransposons⁴⁰. Using the experimentally determined mutation rate in *A. thaliana*¹⁴, we calculated the mean and median age in *A. thaliana* to be 3.1 and 2.1 million years, respectively, compared to 1.1 and 0.6 million years, respectively, in *A. lyrata* (Fig. 5a). In agreement with previous estimates⁴¹, this suggests that LTR retrotransposons have been recently more active in *A. lyrata*. A phylogenetic analysis also supports a greater expansion of specific LTR retrotransposon clades in *A. lyrata* (Fig. 5b). Coupled with higher activity levels of transposable elements in *A. lyrata*, we found that transposable elements are differently distributed in the two species, with *A. lyrata* having a higher proportion of genes with a transposable element nearby than *A. thaliana* (Fig. 5c), and this distance is skewed toward larger values in *A. thaliana* (Supplementary Table 5 and Supplementary Fig. 4). Together, these observations are consistent with a model under which selection purges transposable elements

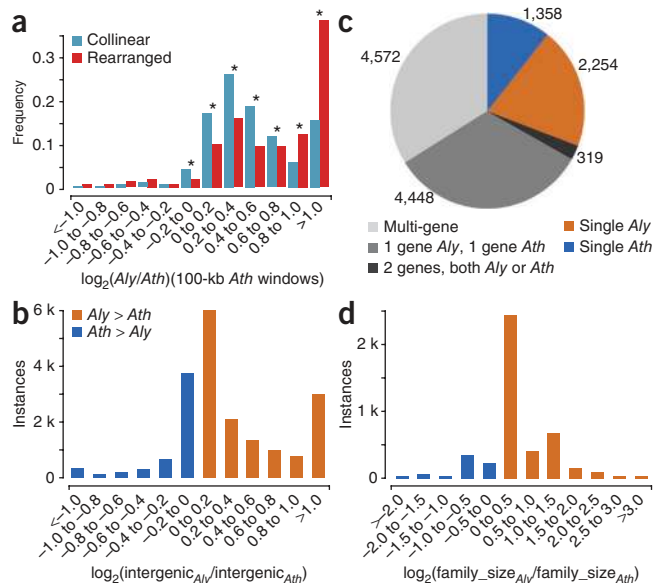
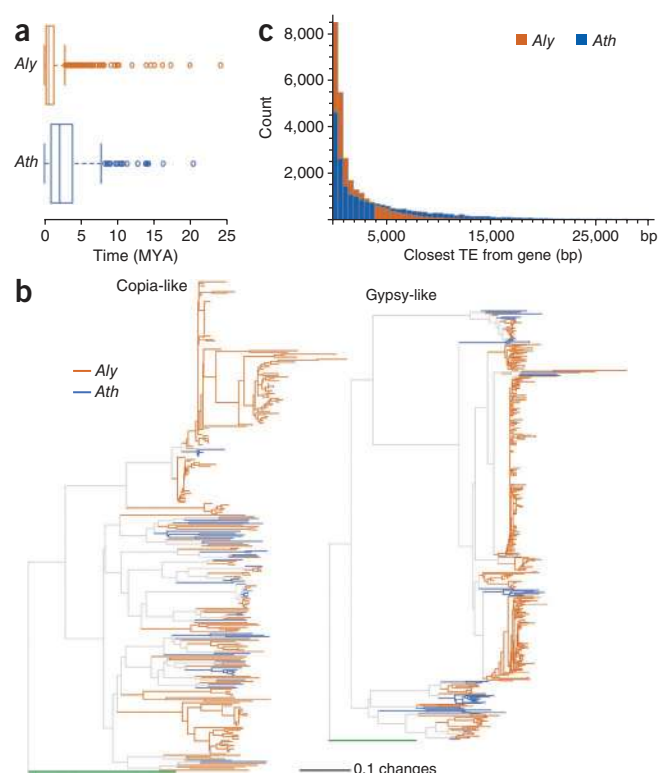


Figure 4 Change in size of collinear and rearranged regions, intergenic regions and gene families. (a) Size comparison of collinear regions, relative to 100-kb windows in *A. thaliana*. Asterisks indicate significant differences (binomial test, $P < 0.001$). (b) Relative size of intergenic intervals. (c) MCL clusters. (d) Relative size of MCL gene families. *Ath*, *A. thaliana*; *Aly*, *A. lyrata*.

Figure 5 Comparison of transposable elements. (a) Estimated insertion times of LTR retrotransposons based on the experimentally determined mutation rate for *A. thaliana*. The whiskers indicate values up to 1.5 times the interquartile range. The difference between the species was highly significant (Wilcoxon rank sum test, $P < 2.2 \times 10^{-16}$). (b) Phylogeny of Ty1/copia-like and Ty3/gypsy-like LTR retrotransposons. *S. cerevisiae* Ty1 and Ty3 that were used as outgroups are indicated in green. (c) Distances of nearest transposable elements from each gene. The difference between the two species was not simply because of fewer transposable elements in the *A. thaliana* genome (Supplementary Table 5 and Supplementary Fig. 4). TE, transposable element; *Ath*, *A. thaliana*; *Aly*, *A. lyrata*; MYA, million years ago.

with deleterious effects on adjacent genes such that transposable elements more distant from genes preferentially survive⁴² and with transposable element elimination having been more efficient in *A. thaliana*. In addition, there is the possibility that transposable elements in *A. lyrata* have experienced less natural selection because they are on average younger.

The evidence presented so far points to *A. thaliana* having suffered a large number of deletions throughout its genome. We can use within-species polymorphisms to shed light on the process by which this has happened. If the *A. thaliana* genome continues to shrink, we would expect fewer segregating insertions than deletions. Using the *A. lyrata* genome as a proxy to determine the derived state among a set of insertion and deletion polymorphisms found throughout the genome of 95 *A. thaliana* plants⁴³, we found a clear excess of deletions over insertions, with 2,685 fixed and 852 segregating deletions compared to 1,941 fixed and 106 segregating insertions. Furthermore, among the fixed differences, deletions are on average longer than insertions (Fig. 6a). If selection were not involved, and if this pattern was only due to mutational bias favoring deletions^{44,45}, deletion and insertion polymorphisms should have similar allele frequencies in the *A. thaliana* population. However, segregating insertions are on average found in fewer individuals than are deletions or SNPs.



Deletions are found in the majority of individuals and many of them are approaching fixation in *A. thaliana* (Fig. 6b). This pattern suggests that deletions are favored over insertions because of selection rather than simple mutational bias, thus leading to a smaller genome.

The pattern of divergence between the two genomes supports this hypothesis. Although more deletions have occurred on the *A. thaliana* than the *A. lyrata* lineage, the bias toward deletions becomes stronger the longer the missing sequence is, and the bias is absent for sequences shorter than 5 bp or so (Fig. 2). This is consistent with a model in which long deletions are selectively favored in *A. thaliana* whereas short deletions are not. We acknowledge that without an outgroup to reconstruct the ancestral state shared by the ancestor of both *A. lyrata* and *A. thaliana*, one cannot accurately determine whether all changes are derived in *A. thaliana*.

In summary, we have presented a high-quality reference genome sequence for *A. lyrata*, which will be a valuable resource for functional, evolutionary and ecological studies in the genus *Arabidopsis*. Several processes contribute to the remarkable difference in genome size between the predominantly selfing *A. thaliana* and the outcrossing *A. lyrata*. In just a few million generations, numerous chromosomal rearrangements have occurred, consistent with theoretical predictions of rearrangements that reduce fitness in heterozygotes being fixed much more easily in strongly selfing species⁴⁶. Though *A. thaliana* has 17% fewer genes than *A. lyrata*, much of the genome size difference seems to be caused by reduced transposable element activity and/or more efficient transposable element elimination in *A. thaliana*,

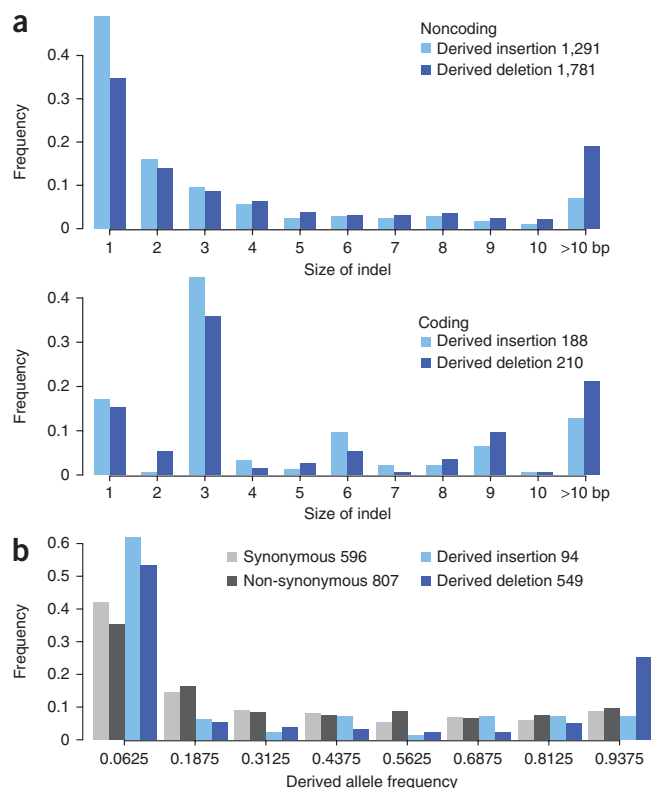


Figure 6 Sizes and allele frequency distribution of insertions and deletions that were either fixed or still segregating in 95 *A. thaliana* individuals⁴³ and that are presumed to be derived based on comparison with the *A. lyrata* allele. (a) Size distribution of fixed insertions and deletions. Insertions and deletions that are multiples of a single codon (3 bp) are overrepresented in coding regions. (b) Allele frequency of segregating noncoding insertion and deletion frequencies compared to that of synonymous and non-synonymous polymorphisms.

especially near genes, as well as shortening of non-transposable element intergenic sequences and introns in *A. thaliana*. Specifically, by making the reasonable assumption that the *A. lyrata* allele presents in the majority of cases the ancestral state, we found that segregating deletions at non-coding sites in *A. thaliana* are skewed toward higher allele frequencies and that both fixed and polymorphic deletions are more common than insertions. Together, this suggests pervasive selection for a smaller genome in *A. thaliana*. Apart from apparent advantages for species with smaller genomes that have been inferred from meta-analyses⁴⁷, the transition to selfing might be an important factor in this process⁴⁶. In addition, a shorter life span may allow a reduction of the genetic repertoire and thus contribute to the smaller genome of *A. thaliana* as well.

What role, if any, genome expansion might play in *A. lyrata* can be addressed once detailed *A. lyrata* polymorphism information as well as closely related outgroup genomes become available, such as the one from *Capsella rubella*, which is currently being assembled. A complete understanding of the processes behind genome contraction and expansion over short time scales will also require better knowledge of mutational events and a deeper understanding of the distribution of, and selection on, noncoding regulatory sequences⁴². For both, high-quality whole-genome sequences of additional *Arabidopsis* relatives will be an important tool.

URLs. MCL, <http://www.micans.org/mcl/>; At_RGenes, http://niblrns.ucdavis.edu/At_RGenes/; RepeatModeler, <http://www.repeatmasker.org/RepeatModeler.html>; EMBOSS, <http://emboss.sourceforge.net/>; RepBase, <http://www.girinst.org/server/RepBase/>; JCVI/TIGR plant repeat database, <http://blast.jcvi.org/euk-blast/index.cgi?project=plant.repeats>; PHYTOZOME portal, <http://www.phytozome.net/alyrata.php>; GenomeThreader, <http://www.genomethreader.org/>.

METHODS

Methods and any associated references are available in the online version of the paper at <http://www.nature.com/naturegenetics/>.

Accession codes. The assembly and annotation (Entrez Genome Project ID 41137) are available from GenBank (accession number ADBK00000000) and from JGI's PHYTOZOME portal (see URLs). Seeds from the MN47 strain have been deposited with the *Arabidopsis* Biological Resource Center under accession number CS22696.

Note: Supplementary information is available on the Nature Genetics website.

ACKNOWLEDGMENTS

The US Department of Energy Joint Genome Institute (JGI) provided sequencing and analyses under the Community Sequencing Program supported by the Office of Science of the US Department of Energy under Contract No. DE-AC02-05CH11231. We are particularly grateful to D. Rokhsar and K. Barry for providing leadership for the project at JGI. We thank J. Borevitz, A. Hall, C. Langley, J. Nasrallah, B. Neuffer, O. Savolainen and S. Wright for contributing to the initial sequencing proposal submitted to the Community Sequencing Program at JGI, C. Lanz and K. Lett for technical assistance, and P. Andolfatto and R. Wing for comments on the manuscript. This work was supported by National Science Foundation (NSF) DEB-0723860 (B.S.G.), NSF DEB-0723935 (M.N.), NSF MCB-0618433 (J.C.C.), NSF IOS-0744579 (M.E.N.), NIH GM057994 (J.B.), grant GABI-DUPLO 0315055 of the German Federal Ministry of Education and Research (K.F.X.M.), ERA-NET on Plant Genomics (ERA-PG) grant ARelatives from the Deutsche Forschungsgemeinschaft (D.W.) and Institute for the Promotion of Innovation through Science and Technology in Flanders (IWT) and the Inter-University Network for Fundamental Research (P6/25, BioMaGNet) (Y.V.d.P.), a Gottfried Wilhelm Leibniz Award of Deutsche Forschungsgemeinschaft (DFG) (D.W.), the Austria Academy of Sciences (M.N.) and the Max Planck Society (D.W. and Y.-L.G.).

AUTHOR CONTRIBUTIONS

J.B., J.C.C., B.S.G., I.V.G., Y.-L.G., K.F.X.M., M.N., Y.V.d.P. and D.W. conceived the study; M.E.N. provided the biological material; J.C., J.-F.C., R.M.C., N.F., J.G. and Y.-L.G. performed the experiments; E.G.B., J.A.F., N.F., H.G., Y.-L.G., G.H., J.D.H., T.T.H., R.P.O., S.O., P.P., A.A.S., J.S., K.S., M.S., X.W. and L.Y. analyzed the data; and Y.-L.G., T.T.H., M.N. and D.W. wrote the paper with contributions from all authors.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Published online at <http://www.nature.com/naturegenetics/>.

Reprints and permissions information is available online at <http://npg.nature.com/reprintsandpermissions/>.

This paper is distributed under the terms of the Creative Commons Attribution Noncommercial-Share Alike license and is freely available to all readers at <http://www.nature.com/naturegenetics/>.

- Greilhuber, J. *et al.* Smallest angiosperm genomes found in Lentibulariaceae, with chromosomes of bacterial size. *Plant Biol.* **8**, 770–777 (2006).
- Gregory, T.R. *et al.* Eukaryotic genome size databases. *Nucleic Acids Res.* **35**, D332–D338 (2007).
- Gaut, B.S. & Ross-Ibarra, J. Selection on major components of angiosperm genomes. *Science* **320**, 484–486 (2008).
- Pellicer, J., Fay, M.F. & Leitch, I.J. The largest eukaryotic genome of them all? *Bot. J. Linn. Soc.* **164**, 10–15 (2010).
- Bennetzen, J.L., Ma, J. & Devos, K.M. Mechanisms of recent genome size variation in flowering plants. *Ann. Bot.* **95**, 127–132 (2005).
- Hawkins, J.S., Proulx, S.R., Rapp, R.A. & Wendel, J.F. Rapid DNA loss as a counterbalance to genome expansion through retrotransposon proliferation in plants. *Proc. Natl. Acad. Sci. USA* **106**, 17811–17816 (2009).
- Piegu, B. *et al.* Doubling genome size without polyploidization: dynamics of retrotransposon-driven genomic expansions in *Oryza australiensis*, a wild relative of rice. *Genome Res.* **16**, 1262–1269 (2006).
- Vitte, C., Panaud, O. & Quesneville, H. LTR retrotransposons in rice (*Oryza sativa*, L.): recent burst amplifications followed by rapid DNA loss. *BMC Genomics* **8**, 218 (2007).
- Woodhouse, M.R. *et al.* Following tetraploidy in maize, a short deletion mechanism removed genes preferentially from one of the two homologs. *PLoS Biol.* **8**, e1000409 (2010).
- Paterson, A.H. *et al.* The *Sorghum bicolor* genome and the diversification of grasses. *Nature* **457**, 551–556 (2009).
- Johnston, J.S. *et al.* Evolution of genome size in Brassicaceae. *Ann. Bot.* **95**, 229–235 (2005).
- Oyama, R.K. *et al.* The shrunken genome of *Arabidopsis thaliana*. *Plant Syst. Evol.* **273**, 257–271 (2008).
- Wright, S.I., Lauga, B. & Charlesworth, D. Rates and patterns of molecular evolution in inbred and outbred *Arabidopsis*. *Mol. Biol. Evol.* **19**, 1407–1420 (2002).
- Ossowski, S. *et al.* The rate and molecular spectrum of spontaneous mutations in *Arabidopsis thaliana*. *Science* **327**, 92–94 (2010).
- Beilstein, M.A., Nagalingum, N.S., Clements, M.D., Manchester, S.R. & Mathews, S. Dated molecular phylogenies indicate a Miocene origin for *Arabidopsis thaliana*. *Proc. Natl. Acad. Sci. USA* **107**, 18724–18728 (2010).
- Kuitinen, H. *et al.* Comparing the linkage maps of the close relatives *Arabidopsis lyrata* and *A. thaliana*. *Genetics* **168**, 1575–1584 (2004).
- Koch, M.A. & Kiefer, M. Genome evolution among cruciferous plants: a lecture from the comparison of the genetic maps of three diploid species—*Capsella rubella*, *Arabidopsis lyrata* subsp. *petraea*, and *A. thaliana*. *Am. J. Bot.* **92**, 761–767 (2005).
- Yogeeswaran, K. *et al.* Comparative genome analyses of *Arabidopsis* spp.: inferring chromosomal rearrangement events in the evolutionary history of *A. thaliana*. *Genome Res.* **15**, 505–515 (2005).
- Lysak, M.A. *et al.* Mechanisms of chromosome number reduction in *Arabidopsis thaliana* and related Brassicaceae species. *Proc. Natl. Acad. Sci. USA* **103**, 5224–5229 (2006).
- Berr, A. *et al.* Chromosome arrangement and nuclear architecture but not centromeric sequences are conserved between *Arabidopsis thaliana* and *Arabidopsis lyrata*. *Plant J.* **48**, 771–783 (2006).
- Swarbreck, D. *et al.* The *Arabidopsis* Information Resource (TAIR): gene structure and function annotation. *Nucleic Acids Res.* **36**, D1009–D1014 (2007).
- Lim, J.K. & Simmons, M.J. Gross chromosome rearrangements mediated by transposable elements in *Drosophila melanogaster*. *Bioessays* **16**, 269–275 (1994).
- Stankiewicz, P. *et al.* Genome architecture catalyzes nonrecurrent chromosomal rearrangements. *Am. J. Hum. Genet.* **72**, 1101–1116 (2003).
- Korbel, J.O. *et al.* Paired-end mapping reveals extensive structural variation in the human genome. *Science* **318**, 420–426 (2007).
- Lee, J., Han, K., Meyer, T.J., Kim, H.S. & Batzer, M.A. Chromosomal inversions between human and chimpanzee lineages caused by retrotransposons. *PLoS ONE* **3**, e4047 (2008).
- Braumann, I., van den Berg, M.A. & Kempken, F. Strain-specific retrotransposon-mediated recombination in commercially used *Aspergillus niger* strain. *Mol. Genet. Genomics* **280**, 319–325 (2008).

27. Woodhouse, M.R., Pedersen, B. & Freeling, M. Transposed genes in *Arabidopsis* are often associated with flanking repeats. *PLoS Genet.* **6**, e1000949 (2010).
28. Ranz, J.M. *et al.* Principles of genome evolution in the *Drosophila melanogaster* species group. *PLoS Biol.* **5**, e152 (2007).
29. The Chimpanzee Sequencing and Analysis Consortium. Initial sequence of the chimpanzee genome and comparison with the human genome. *Nature* **437**, 69–87 (2005).
30. Clark, R.M. *et al.* Common sequence polymorphisms shaping genetic diversity in *Arabidopsis thaliana*. *Science* **317**, 338–342 (2007).
31. Borevitz, J.O. *et al.* Genome-wide patterns of single-feature polymorphism in *Arabidopsis thaliana*. *Proc. Natl. Acad. Sci. USA* **104**, 12057–12062 (2007).
32. Enright, A.J., Van Dongen, S. & Ouzounis, C.A. An efficient algorithm for large-scale detection of protein families. *Nucleic Acids Res.* **30**, 1575–1584 (2002).
33. Michelmore, R.W. & Meyers, B.C. Clusters of resistance genes in plants evolve by divergent selection and a birth-and-death process. *Genome Res.* **8**, 1113–1130 (1998).
34. Thomas, J.H. Adaptive evolution in two large families of ubiquitin-ligase adapters in nematodes and plants. *Genome Res.* **16**, 1017–1030 (2006).
35. Yang, X. *et al.* The F-box gene family is expanded in herbaceous annual plants relative to woody perennial plants. *Plant Physiol.* **148**, 1189–1200 (2008).
36. Tuskan, G.A. *et al.* The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray). *Science* **313**, 1596–1604 (2006).
37. Jaillon, O. *et al.* The grapevine genome sequence suggests ancestral hexaploidization in major angiosperm phyla. *Nature* **449**, 463–467 (2007).
38. Velasco, R. *et al.* A high quality draft consensus sequence of the genome of a heterozygous grapevine variety. *PLoS ONE* **2**, e1326 (2007).
39. Li, L., Stoeckert, C.J. Jr. & Roos, D.S. OrthoMCL: identification of ortholog groups for eukaryotic genomes. *Genome Res.* **13**, 2178–2189 (2003).
40. SanMiguel, P., Gaut, B.S., Tikhonov, A., Nakajima, Y. & Bennetzen, J.L. The paleontology of intergene retrotransposons of maize. *Nat. Genet.* **20**, 43–45 (1998).
41. Devos, K.M., Brown, J.K. & Bennetzen, J.L. Genome size reduction through illegitimate recombination counteracts genome expansion in *Arabidopsis*. *Genome Res.* **12**, 1075–1079 (2002).
42. Hollister, J.D. & Gaut, B.S. Epigenetic silencing of transposable elements: a trade-off between reduced transposition and deleterious effects on neighboring gene expression. *Genome Res.* **19**, 1419–1428 (2009).
43. Nordborg, M. *et al.* The pattern of polymorphism in *Arabidopsis thaliana*. *PLoS Biol.* **3**, e196 (2005).
44. Petrov, D.A., Sangster, T.A., Johnston, J.S., Hartl, D.L. & Shaw, K.L. Evidence for DNA loss as a determinant of genome size. *Science* **287**, 1060–1062 (2000).
45. Petrov, D.A., Lozovskaya, E.R. & Hartl, D.L. High intrinsic rate of DNA loss in *Drosophila*. *Nature* **384**, 346–349 (1996).
46. Charlesworth, B. Evolutionary rates in partially self-fertilizing species. *Am. Nat.* **140**, 126–148 (1992).
47. Knight, C.A., Molinari, N.A. & Petrov, D.A. The large genome constraint hypothesis: evolution, ecology and phenotype. *Ann. Bot.* **95**, 177–190 (2005).

ONLINE METHODS

Sequencing and assembly. *Arabidopsis lyrata* strain MN47 was derived by forced selfing from material collected in Michigan, USA by C. Langley (University of California Davis). It was inbred five times before DNA was extracted for sequencing. Libraries with various insert sizes including fosmids and bacterial artificial chromosomes (BACs) were dideoxy sequenced on ABI 3730XL capillary sequencers. Reads were assembled with Arachne⁴⁸, and collinearity information was integrated with marker information from genetic maps^{16,18,19} to reconstruct the eight linkage groups. Additional details and specifics are presented in the **Supplementary Note**.

Annotation. The genome was annotated using *ab initio* and homology-based gene predictors along with RNA-Seq data (**Supplementary Note**). The complete details are described in the **Supplementary Note**.

MCL cluster analyses. MCL (mcl-06-058 package; see URLs) was used with default parameters (-I 2, -S 6) based on clustering of hits with E value $\leq 10^{-5}$. MCL uses a Markov cluster algorithm that attempts to overcome many of the difficulties with protein sequence clustering, such as the presence of multi-domain proteins, peptide fragments and proteins with very common domains. This method has been used for a variety of animal genomes^{49–51}.

OrthoMCL analysis. Orthologous gene clusters were computed from OrthoMCL comparisons³⁹ of four dicotyledonous species with finished genomes: *A. thaliana* and *A. lyrata*, *Populus trichocarpa*³⁶ and *Vitis vinifera*^{37,38}. A search for potentially missed genes in both *Arabidopsis* genomes resulted in minor adjustments of the OrthoMCL clusters. Instead of 10,573 clusters, 10,878 clusters now contained at least one gene from each the four species, and instead of 5,699 clusters, 5,800 clusters were *Arabidopsis* specific. To determine deleted or newly generated orthologs (by OrthoMCL definition) between the two species, we focused on clusters specific for either *A. lyrata* or *A. thaliana*. For both species, there are two cluster types: those that are supported by members in *P. trichocarpa* and/or *V. vinifera* (supported specific cluster (SSC)) and clusters exclusively found in one of the *Arabidopsis* species (exclusive specific cluster (ESC)). We did not consider 2,939 and 6,103 unclustered genes (singletons) in *A. thaliana* and *A. lyrata*, respectively.

In our initial analysis, we detected 354 SSCs and 161 ESCs for *A. thaliana* and 168 SSCs and 833 ESCs for *A. lyrata*. Whole-genome projects, however, may contain false-positive as well as missed or incomplete partial gene calls that impose difficulties for OrthoMCL to detect orthologous relationships. To ensure that genes from the previously detected SSCs were indeed specific for one of the *Arabidopsis* species, we reevaluated the absence or presence of specific gene calls in the two genome sequences. Previously missed genes detected by GenomeThreader (see URLs) were added to each of the gene sets and the OrthoMCL analysis was repeated.

F-box and NB-LRR gene analysis. Using F-box PF00646.hmm as an HMM profile with hmmsearch (E value $\leq 10^{-5}$), 394 hits were found in *A. thaliana* and 461 hits were found in *A. lyrata*. Alignment of these sequences was optimized with the PF00646 seed using ClustalX 2.0 (ref. 52). The final alignment was produced by aligning with hmmlalign against PF00646.hmm to construct an *Arabidopsis*-specific HMM F-box profile. With this HMM profile, 502 hits were found in *A. thaliana*, and 596 hits were found in *A. lyrata*. hmmlalign was used to align all of these against PF00646.hmm.

A blastp search (E value $\leq 10^{-10}$) was performed with the nucleotide-binding site (NB) domain (based on HMMEMIT, from At_RGenes; see URLs). The NB domains of the retrieved proteins (142 in *A. thaliana* and 162 in *A. lyrata*) were aligned using ClustalX⁵². This alignment was used to develop an *Arabidopsis*-specific HMM profile, which was used to search the complete set of proteins encoded by both the two genomes (cutoff $E \leq 10^{-5}$).

RepeatMasker analyses. To develop *de novo* repeat libraries for both species, we used RepeatModeler (version Beta 1.0.3; see URLs). To reduce false positives, unclassified repeats were compared to annotated genes, and we eliminated all that had at least 80% identity to annotated genes over at least 80 bp (GenBank: green plant GB all (protein); blastx with E value $\leq 10^{-10}$). The remaining RepeatModeler predictions were classified with the 80-80-80 rule⁵³,

grouping repeats if they shared at least 80% identity over at least 80% of the aligned sequence, which had to be at least 80 bp long. The identified repeats were appended to RepBase (*Arabidopsis* library, RM database version 20080611) resulting in a final library with 1,152 repeat units. The final libraries were used to annotate transposable elements using RepeatMasker version 3.2.5.

LTR retrotransposons. Intact LTR retrotransposons were identified *de novo* using LTR_STRUC⁵⁴ with default parameters. Based on the sequence divergence between the two LTRs of the same element, insertion times were estimated. All LTR pairs were aligned using MUSCLE⁵⁵, and the distance K between them was calculated with the Kimura two-parameter model using the distmat program implemented in the EMBOSS package (see URLs). The insertion time T was calculated as $T = K/(2(r))$, with r as the rate of nucleotide substitution. The molecular clock was set based on the observed mutation rate of 7×10^{-9} per site per generation (assumed to equal one year)¹⁴.

Classification and phylogeny of LTR retrotransposons. LTR retrotransposons can be classified into Ty1/copia-like and Ty3/gypsy-like elements⁵⁶. We classified repeats using RepBase (version 13.08; see URLs) and blastn (E value $\leq 10^{-10}$) and by direct comparison against the JCVI/TIGR plant repeat database (see URLs). All intact LTR retrotransposons were compared with blastx (E value $\leq 10^{-10}$) against a conserved 156-amino-acid segment corresponding to the reverse transcriptase domain⁵⁷ of Ty1/copia-like and Ty3/gypsy-like sequences, and this segment was then used for phylogenetic reconstruction using PAUP* version 4.0b10 (ref. 58) and a neighbor-joining method. As the outgroup sequence, we used the reverse transcriptase domain from yeast Ty1 and Ty3 elements, respectively⁵⁷.

Detection and analysis of chromosomal breakpoints. Genome-wide collinearity was detected by running i-ADHoRe⁵⁹ on the core-orthologous genes, allowing for the identification of breakpoints, including inversions and nested inversions. For each inversion, the regions 10 kb upstream and downstream of the delimiting breakpoints were compared to each other using blastn (word size 4), tblastx (word size 1) and SSEARCH^{60–62}. Tblastx outperforms blastn for coding regions. In noncoding regions, SSEARCH is more sensitive than blastn but is computationally less efficient and, hence, is the most useful for comparison of shorter sequences. Only one hit per strand was reported. Therefore, for each pair of inversion flanking regions, all combinations of repeats and protein coding genes were evaluated. Default settings were used for gap penalties. An E value of ≤ 0.01 was considered as indicating similarity between the upstream and downstream regions.

Similarity of syntenic regions. To investigate the nucleotide divergence of intergenic regions around coding genes (**Supplementary Fig. 1b**), we extracted for each syntenic gene pair the 2-kb sequences 5' of the start codon and 3' from the stop codon. If the neighboring gene was closer than 2 kb, the extracted sequence was accordingly trimmed. Coding sequences of syntenic genes were also analyzed. Global alignments of syntenic sequences were generated using the Needleman-Wunsch algorithm as implemented in the EMBOSS package 5.0 (using default parameters). The sequence identity of the coding regions was measured over the full-length alignment. To investigate whether the divergence of an intergenic sequence is affected by the relative orientation to neighboring genes, upstream sequences were split into head-to-tail and head-to-head groups, and downstream sequences were split into tail-to-head and tail-to-tail groups.

Fixed insertions and deletions. To identify fixed insertions and deletions among 1,238 fragments that had been amplified by PCR and sequenced in 95 *A. thaliana* individuals⁴³, two representative sequences for each fragment were first constructed to represent the insertion and deletion states among all segregating indels. The representative sequence consisting of insertions was then queried against the *A. lyrata* genome with both BLAT⁶³ (-maxGap = 100 - extendThroughN - minIdentity = 80) and BLAST⁶⁴ (-e 0.00001 -F F -G -5 -E -1). Based on the longest hit from the union of hits obtained by both methods, the representative sequences for each alignment were profile aligned with the *A. lyrata* allele with MAFFT⁶⁵. Fixed insertions and deletions were identified in the resulting alignment.

Segregating insertions and deletions. A similar procedure to that described above was used to identify the *A. lyrata* allele (presumed ancestral state) for each polymorphic indel in *A. thaliana*. Instead of querying the entire fragment, we queried each insertion allele along with 25 bp flanking each side against the *A. lyrata* genome using BLAT. For each polymorphic indel, we filtered for the best hit that spanned both sides of the indel site (by at least 3 bp) and reported each indel as either a derived insertion (if the *A. lyrata* allele was a deletion in the resulting profile alignment) or a derived deletion (if the *A. lyrata* allele was not a deletion).

Data and seed availability. The assembly and annotation (Entrez Genome Project ID 41137) are available from GenBank (see Accession Codes section) and from JGI's PHYTOZOME portal (see URLs). Seeds of the MN47 strain have been deposited with the *Arabidopsis* Biological Resource Center (see Accession Codes section).

48. Jaffe, D.B. *et al.* Whole-genome sequence assembly for mammalian genomes: Arachne 2. *Genome Res.* **13**, 91–96 (2003).
49. Demuth, J.P., De Bie, T., Stajich, J.E., Cristianini, N. & Hahn, M.W. The evolution of mammalian gene families. *PLoS ONE* **1**, e85 (2006).
50. Prachumwat, A. & Li, W.H. Gene number expansion and contraction in vertebrate genomes with respect to invertebrate genomes. *Genome Res.* **18**, 221–232 (2008).
51. Drosophila 12 Genomes Consortium. *et al.* Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature* **450**, 203–218 (2007).
52. Thompson, J.D., Gibson, T.J., Plewniak, F., Jeanmougin, F. & Higgins, D.G. The CLUSTAL-X Windows interface: flexible strategies for multiple sequence alignment aided by quality analysis tools. *Nucleic Acids Res.* **25**, 4876–4882 (1997).
53. Wicker, T. *et al.* A unified classification system for eukaryotic transposable elements. *Nat. Rev. Genet.* **8**, 973–982 (2007).
54. McCarthy, E.M. & McDonald, J.F. LTR-STRUC: a novel search and identification program for LTR retrotransposons. *Bioinformatics* **19**, 362–367 (2003).
55. Edgar, R.C. MUSCLE: a multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics* **5**, 113 (2004).
56. Xiong, Y. & Eickbush, T.H. Origin and evolution of retroelements based upon their reverse transcriptase sequences. *EMBO J.* **9**, 3353–3362 (1990).
57. Zhang, X. & Wessler, S.R. Genome-wide comparative analysis of the transposable elements in the related species *Arabidopsis thaliana* and *Brassica oleracea*. *Proc. Natl. Acad. Sci. USA* **101**, 5589–5594 (2004).
58. Swofford, D.L. PAUP*: Phylogenetic Analysis Using Parsimony (*and Other Methods): Version 4. (Sinauer Associates, Sunderland, Massachusetts, USA, 2003).
59. Simillion, C., Vandepoele, K., Saeys, Y. & Van de Peer, Y. Building genomic profiles for uncovering segmental homology in the twilight zone. *Genome Res.* **14**, 1095–1106 (2004).
60. Altschul, S.F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).
61. Smith, T.F. & Waterman, M.S. Identification of common molecular subsequences. *J. Mol. Biol.* **147**, 195–197 (1981).
62. Pearson, W.R. Searching protein sequence libraries: comparison of the sensitivity and selectivity of the Smith-Waterman and FASTA algorithms. *Genomics* **11**, 635–650 (1991).
63. Kent, W.J. BLAT—the BLAST-like alignment tool. *Genome Res.* **12**, 656–664 (2002).
64. Altschul, S.F., Gish, W., Miller, W., Myers, E.W. & Lipman, D.J. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
65. Katoh, K., Asimenos, G. & Toh, H. Multiple alignment of DNA sequences with MAFFT. *Methods Mol. Biol.* **537**, 39–64 (2009).