

The complex sequence landscape of maize revealed by single molecule technologies

Yinping Jiao¹, Paul Peluso², Jinghua Shi³, Tiffany Liang³, Michelle C. Stitzer⁴, Bo Wang¹, Michael S. Campbell¹, Joshua C. Stein¹, Xuehong Wei¹, Chen-Shan Chin², Katherine Guill⁵, Michael Regulski¹, Sunita Kumari¹, Andrew Olson¹, Jonathan Gent⁶, Kevin L. Schneider⁷, Thomas K. Wolfgruber⁷, Michael R. May⁸, Nathan M. Springer⁹, Eric Antoniou¹, Richard McCombie¹, Gernot G. Presting⁷, Michael McMullen⁵, Jeffrey Ross-Ibarra¹⁰, Kelly Dawe⁶, Alex Hastie³, David R. Rank², Doreen Ware^{1,11*}

Affiliations:

¹ Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724

² Pacific Biosciences, Menlo Park, CA 94025

³ BioNano Genomics, San Diego, CA 92121

⁴ Department of Plant Sciences and Center for Population Biology, University of California, Davis, Davis, CA 95616

⁵ USDA-ARS, Plant Genetics Research Unit, Columbia, MO 65211

⁶ University of Georgia, Athens, Georgia 30602

⁷ Department of Molecular Biosciences and Bioengineering, University of Hawaii, Honolulu, HI 96822

⁸ Department of Evolution and Ecology, University of California, Davis, CA 95616

⁹ Department of Plant Biology, University of Minnesota, St. Paul, MN 55108

¹⁰ Department of Plant Sciences, Center for Population Biology, and Genome Center, University of California, Davis, CA 95616

¹¹ USDA-ARS, NEA Robert W. Holley Center for Agriculture and Health, Cornell University, Ithaca, New York 14853

* Corresponding author: Doreen Ware (ware@cshl.edu, Doreen.Ware@ARS.USDA.GOV)

ABSTRACT

Complete and accurate reference genomes and annotations provide fundamental tools for characterization of genetic and functional variation. These resources facilitate elucidation of biological processes and support translation of research findings into improved and sustainable agricultural technologies. Many reference genomes for crop plants have been generated over the past decade, but these genomes are often fragmented and missing complex repeat regions. Here, we report the assembly and annotation of maize, a genetic and agricultural model crop, using Single Molecule Real-Time (SMRT) sequencing and high-resolution genome map. Relative to the previous reference genome, our assembly features a 52-fold increase in contig length and significant improvements in the assembly of intergenic spaces and centromeres. Characterization of the repetitive portion of the genome revealed over 130,000 intact transposable elements (TEs), allowing us to identify TE lineage expansions unique to maize. Gene annotations were updated using 111,000 full-length transcripts obtained by SMRT sequencing. In addition, comparative optical mapping of two other inbreds revealed a prevalence of deletions in the region of low gene density region and maize lineage-specific genes.

Maize is the most productive and widely grown crop in the world, as well as a foundational model for genetics and genomics¹. An accurate genome assembly for maize is critical for all forms of basic and applied research, which will enable increases in yield to feed the growing world population. The current assembly of the maize genome, based on Sanger sequencing, was first published in 2009². Although this initial reference enabled rapid progress in maize genomics³, the original assembly is composed of more than a hundred thousand small contigs, many of which are arbitrarily ordered and oriented, significantly complicating detailed analysis of individual loci⁴ and impeding investigation of intergenic regions crucial to our understanding of phenotypic variation^{5,6} and genome evolution^{7,8}.

Here we report a vastly improved *de novo* assembly and annotation of the maize reference genome (Figure 1). Based on 65X Single Molecule Real-Time sequencing we assembled the genome of the maize inbred B73 into 2,958 contigs, where half of the total assembly is made up of contigs larger than 1.2 Mb (Table 1, Extended Data Figure 1, 2a). The assembly of the long reads was then integrated with a high quality optical genome map (Extended Data Figure 1, Extended Data Table 1) to create a hybrid assembly consisting of 625 scaffolds (Table 1). To build chromosome-level super-scaffolds, we combined the hybrid assembly with a minimum tiling path generated from Sanger sequence of bacterial artificial chromosomes (BACs)⁹ and a high-density genetic map¹⁰ (Extended Data Figure 2b). After gap-filling and error correction using short sequence reads, the total size of maize B73 RefGen_v4 pseudomolecules was 2,106 Mb. The new reference assembly has 2,522 gaps, the genome maps covering approximately half (1,115) of these allowed us to estimate their mean length to be (27 kb) (Extended Data Figure 2c).

Comparison of the new assembly to the previous BAC-based Sanger assembly revealed >99.9% sequence identity and a 52-fold increase in the mean contig length, with 84% of the BACs spanned by a single contig from the long reads assembly. Alignment of ChIP-seq data for centromere-specific histone H3 (CENH3)¹¹ revealed that centromeres are accurately placed and largely intact. A number of previously identified¹² megabase-sized mis-oriented pericentromeric regions were also corrected (Extended Data Figure 3a,b). Moreover, the ends of the chromosomes are properly identified on 14 of the 20

chromosome arms based on the presence of tandem telomeric repeats and knob 180 sequences (Extended Data Figure 3a,c).

Our assembly made substantial improvements in the gene space including resolution of gaps and misassemblies, correct order and orientation. We also updated the annotation of our new assembly, resulting in consolidation of gene models (Extended Data Figure 4a). Published full-length cDNA data¹³ improved the annotation of alternative splicing by more than doubling the number of alternative transcripts from 1.6 to 3.3 per gene (Extended Data Figure 5a) with about 70% of genes with support from the full-length transcripts. Our reference assembly also vastly improved the coverage of regulatory sequences, decreasing the number of genes exhibiting gaps in the 3 kb region(s) flanking coding sequence from 20% to <1% (Extended Data Figure 4b, Extended Data 5b). The more complete sequence enabled significant improvements in the annotation of core promoter elements, especially the TATA-box, CCAAT-box, and Y patch motifs (Supplementary information). Quantitative genetic analyses have shown that the genetic variation in regulatory regions explain a large number of phenotypes^{5,6}, suggesting that the new reference will dramatically improve our ability to identify and predict functional genetic variation.

After its divergence from *Sorghum*, the maize lineage underwent genome doubling followed by diploidization and gene loss. Previous work showed that gene loss had bias toward one of the parental genomes¹⁴, but our new assembly and annotation paint a more dramatic picture, revealing that 56% of syntenic sorghum orthologs map uniquely to the dominant maize subgenome (designated A, total size 1.16 Gb), whereas only 24% map uniquely to subgenome B (total size, 0.63Gb). Gene loss in maize has primarily been considered in the context of polyploidy and functional redundancy^{14,15}, but we found that despite its polyploidy, maize has lost a larger proportion (14%) of the 22,048 ancestral gene orthologs than any of the other four grass species evaluated to date (*Sorghum*, rice, *Brachypodium distachyon*, and *Setaria italica*, Extended Data Figure 6). Nearly one-third of these losses are specific to maize, and analysis of a restricted high-confidence set revealed enrichment for genes involved in biotic and abiotic stresses (Extended Data Table 2), e.g., NB-ARC domain disease resistance genes and the serpin protease inhibitor involved in pathogen defense and programmed cell death^{16,17}.

Transposable elements (TEs) were first reported in maize¹⁸ and have since been shown to play important roles in shaping genome evolution and gene regulatory networks of many species¹⁹. The majority of the maize genome is derived from TEs^{2,20}, and careful study of a few regions has revealed a characteristic structure of sequentially nested retrotransposons^{20,21} and the effect of deletions and recombination on retrotransposon evolution²². In the annotation of the original maize assembly, however, fewer than 1% of LTR retrotransposon copies were intact²³. By applying a novel homology-independent annotation pipeline to our assembly (Extended Data Table 3), we identified 1,268Mb (130,604 copies) of structurally intact retrotransposons, of which 661 Mb (70,035 copies) are nested retrotransposon copies disrupted by the insertion of other TEs, 8.7 Mb (14,041 copies) of DNA terminal inverted repeat transposons, and 76 Mb (21,095 copies) of helitrons. To understand the evolutionary history of maize LTR retrotransposons, we also applied our annotation pipeline to the sorghum reference genome, and used reverse transcriptase protein domain sequences accessible due to the improved assembly of the internal protein coding domains of maize LTR retrotransposons to reconstruct the phylogeny of maize and sorghum LTR retrotransposon families. Despite a higher overall rate of diversification of LTR TEs in the maize lineage, consistent with its larger genome size, differences in LTR retrotransposon content between genomes were primarily the result of dramatic expansion of distinct families in both lineages (Figure 2).

Maize exhibits tremendous genetic diversity²⁴, and both nucleotide polymorphisms and structural variations play important roles in its phenotypic variation^{8,25}. However, genome-wide patterns of structural variation in plant genomes are difficult to assess²⁶, and previous efforts have relied on short-read mapping, which misses the vast majority of intergenic spaces where most rearrangements are likely to occur⁸. To investigate structural variation at a genome-wide scale, we generated genome maps (Extended Data Table 1) for two additional maize inbred lines: the tropical line Ki11, one of the founders of the maize nested association mapping (NAM) population²⁷, and W22, which has served as a foundation for studies of maize genetics²⁸. Due to the high degree of genomic diversity among these lines, only 32% of the assembled 2,216 Mb map of Ki11 and 39% of the 2,280 Mb W22 map could be mapped to our new B73 reference via common restriction patterns (Table 2, Figure 3A, Extended Data Figure 7). Within the

aligned regions, approximately 32% of the Ki11 and 24% of the W22 genome maps exhibited clear evidence of structural variation, including 3,410 insertions and 3,300 deletions (Table 2). The average indel size was ~20kb, with a range from 100 bp (the smallest detectable) to over 1 Mb (Figure 3B). More than 90% of the indels were unique to one inbred or the other, indicating a high level of structural diversity in maize. As short-read sequence data are available from both Ki11 and W22⁸, we analyzed 1,451 of the largest (<10kb) deletions and found that 1,083 were supported by a clear reduction in read depth (Figure 3C). The confirmed deletions occurred in regions of low gene density (4.4 genes/Mb compared to genome-wide average of 18.7 genes/Mb). One third (83/257) of the genes missing in Ki11 and W22 lack putative orthologs in four grasses (rice, sorghum, *Brachypodium*, and *Setaria*), consistent with prior data²⁹.

Although maize is often considered to be a large-genome crop, most major food crops have even larger genomes with more complex repeat landscapes³⁰. Our improved assembly of the B73 genome, generated using SMRT Sequencing technology, demonstrates that additional assemblies of other maize inbred lines and similar high-quality assemblies of other repeat-rich and large-genome plants are feasible. Additional high quality assemblies will in turn extend our understanding of the genetic diversity that forms the basis of the phenotypic diversity in maize and other economically important plants.

Supplementary Information

Data Availability

Raw reads, genome assembly sequences, and gene annotations have been deposited at NCBI under BioProject number PRJNA10769 and BioSample number SAMN04296295. PacBio whole-genome sequencing data and Illumina data were deposited in the NCBI SRA database under accessions SRX1472849 and SRX1452310, respectively. The GenBank accession number of the genome assembly and annotation is LPUQ00000000. A genome browser including genome feature tracks and ftp is available from Gramene: http://ensembl.gramene.org/Zea_mays/Info/Index

Contribution

D.W. and Y.J. designed and conceived the research, M.M and K.G. prepared DNA sample for PacBio SMRT sequencing, D.R.R., P.P., E.A., and R.C. performed PacBio SMRT sequencing, B.W., J.S., K.D., T.L. and A.H. generated BioNano genome maps, M.R. generated Illumina sequencing data, Y.J., T.L., J.S., C.C., and A.H. performed the genome assembly, J.C.S., M.S.C., X.W., B.W., Y.J., and S.K. performed gene annotation and evolutionary study, M.C.S., M.R.M., N.M.S., and J.R-I. performed transposable element analysis, J.G., J.S., K.D., K.L.S., T.K.W., G.G.P., and Y.J. performed the analysis of centromeres and telomeres. B.W., Y.J., J.S., T.L., A.H. and K.D. performed the structural variation study. X.W., J.C.S. and Y.J. contributed to the data release. Y.J., J.R-I., K.D., G.G.P. and D.W. wrote the paper. All authors contributed to the revision of the manuscript.

Competing Financial Interests

P.P., C.C. and D.R.R. are full-time employees of Pacific Biosciences. J.S., T.L., and A.H. are employees at BioNano Genomics, Inc., and own company stock options. All other authors declare no competing financial interests.

Acknowledgment

Y.J., B.W., J. S., M.C., X.W., S.K. and D.W. were supported by NSF Gramene grant IOS-1127112, NSF Cereal Gene Discovery grant 1032105, USDA-ARS CRIS 1907-

21000-030-00D and NSF Plant Genome award 1238014. J.R.I. would like to acknowledge support from USDA Hatch project CA-D-PLS-2066-H and NSF Plant Genome award 1238014. K.D. would like to acknowledge support from NSF Plant Genome award 1444514. G.G.P. acknowledges support from NSF grant 1444624 and USDA NIFA project HAW05022-H. M.S.C would also like to acknowledge support from NSF PGRP PRFB 1523793. The authors thank Sergey Koren (National Human Genome Research Institute) for sharing genome assembly related scripts.

Online Methods

Genome assembly

De novo assembly of the genome sequencing data: *De novo* assembly of the long reads from SMRT Sequencing was performed using two assemblers: the Celera Assembler PBcR –MHAP pipeline³¹ and Falcon³² with different parameter settings. Quiver from SMRT Analysis v2.3.0 was used to polish base calling of contigs. The three independent assemblies were evaluated by aligning with the genome maps.

Contamination of contigs by bacterial and plasmid genomes was eliminated using the NCBI GenBank submission system³³. Curation of the assembly, including resolution of conflicts between the contigs and the genome map and removal of redundancy at the edges of contigs, is described in the supplemental material.

Hybrid scaffold construction: To create hybrid scaffolds, conflict-resolved sequence contigs and genome maps were aligned and merged with RefAligner using a threshold P value of 1×10^{-11} ³⁴. To maximize the sequence content in the hybrid scaffolds, all (scaffolded and unscaffolded) sequence contigs were aligned to the hybrid scaffolds using a less stringent threshold P value (1×10^{-8}). The unscaffolded sequence contigs that aligned to the hybrid scaffold and maps at least 50% of their length and that did not overlap with already scaffolded sequence contigs were added to the hybrid scaffolds.

Pseudomolecule construction: Sequences from BACs on the physical map that were used to build the maize V3 pseudomolecules were aligned to contigs using MUMMER package³⁵ with the following parameter settings: “-l(minimum length of a single match) 100 -c(the minimum length of a cluster of matches) 1000”. Next, to make sure only the unique hits were used, the alignment hits were filtered with the following parameters: “-i(the minimum alignment identity) 98 -l(the minimum alignment length) 10000”. The scaffolds were then ordered and oriented into pseudochromosomes using the order of BACs as a guide. For quality control, we mapped the SNP markers from a genetic map built from an intermated maize recombinant inbred line population (Mo17 \times B73)¹⁰. Contigs with markers not located in pseudochromosomes from the physical map were placed into the AGP (A Golden Path) using the genetic map.

Further polishing of pseudomolecules: A round of gap filling was applied to the raw pseudochromosomes using Pbjelly (--maxTrim=0, --minReads=2). The pseudomolecules were then polished again using the Quiver pipeline from SMRT Analysis v2.3.0. To increase the accuracy of the base calls, we performed two lanes of sequencing on each DNA sample (library size = 450bp) using Illumina 2500 Rapid run, which generated about 100-fold 250PE data of Maize B73. Reads were aligned to the assembly using BWA-mem³⁶. To correct base calling in the assembly, SAMtools³⁷ was used to generate the BAM format alignment for the Pilon pipeline³⁸, using reads with sequencing and alignment quality above 20.

Annotation

To generate a comprehensive annotation of transposable elements, we generated a Structural identification pipeline incorporating several tools, including LTRharvest³⁹, LTRdigest⁴⁰, SINE-Finder⁴¹, MGEScan-non-LTR⁴², MITE-hunter⁴³, HelitronScanner⁴⁴, and others (details in Supplementary Information). The scripts, parameters, and intermediate files of each TE superfamily are available at https://github.com/mcstitzer/agpv4_te_annotation/tree/master/ncbi_pseudomolecule. The MAKER-P pipeline was used to annotate protein-coding genes⁴⁵. Evidence included publicly available full-length cDNA⁴⁶, *de novo* assembled transcripts from short read mRNA-seq⁴⁷, Iso-Seq full-length transcripts¹³, and proteins from other species. The gene models were filtered to remove transposons and low-confidence predictions. Additional alternative transcript isoforms were obtained from the Iso-Seq data.

Structural Variation

Leaves were used to prepare high molecular weight DNA and genome maps were constructed as described above for B73. Structural variant calls were generated based on alignment to the reference map B73 V4 chromosomal assembly using the Multiple Local Alignment algorithm (RefSplit)³⁴. A structural variant was identified as an alignment outlier^{34,48}, defined as two well-aligned regions separated by a poorly aligned region with a large size difference between the reference genome and the map or by one or more unaligned sites, or alternately as a gap between two local alignments. A confidence score

was generated by comparing the non-normalized p-values of the two well-aligned regions and the non-normalized log-likelihood ratio⁴⁹ of the unaligned or poorly aligned region. With a confidence score threshold of 3, RefSplit is sensitive to insertions and deletions as small as 700 bp and other changes such as inversions and complex events which could be balanced. Insertion and deletion calls were based on an alignment outlier p-value threshold of 1×10^{-4} . Insertions or deletions that crossed gaps in the B73 pseudomolecules, or that were heterozygous in the genome maps, were excluded. Considering the resolution of the genome map, only insertion and deletions larger than 100 bp were used for subsequent analyses. To obtain high-confidence deletion sequences, sequencing reads from the maize HapMap2 project⁸ for Ki11 and W22 were aligned to our new B73 v4 reference genome using Bowtie2⁵⁰. Read depth was calculated from reads with mapping quality > 20 in 10-kb windows with step size of 1 kb. Windows where read depth dropped below 10 in Ki11 and 20 in W22 (short reads sequencing depths in Ki11 and W22 are respectively 2.32X and 4.04X) in the deleted region were retained for further analysis.

References

- 1 Hake, S. & Ross-Ibarra, J. Genetic, evolutionary and plant breeding insights from the domestication of maize. *eLife* **4**, doi:10.7554/eLife.05861 (2015).
- 2 Schnable, P. S. *et al.* The B73 maize genome: complexity, diversity, and dynamics. *Science* **326**, 1112-1115, doi:10.1126/science.1178534 (2009).
- 3 Edwards, D., Batley, J. & Snowdon, R. J. Accessing complex crop genomes with next-generation sequencing. *Theor Appl Genet* **126**, 1-11, doi:10.1007/s00122-012-1964-x (2013).
- 4 Fouquet, R. *et al.* Maize rough endosperm3 encodes an RNA splicing factor required for endosperm cell differentiation and has a nonautonomous effect on embryo development. *The Plant cell* **23**, 4280-4297, doi:10.1105/tpc.111.092163 (2011).
- 5 Wallace, J. G. *et al.* Association mapping across numerous traits reveals patterns of functional variation in maize. *PLoS Genet* **10**, e1004845, doi:10.1371/journal.pgen.1004845 (2014).
- 6 Rodgers-Melnick, E., Vera, D. L., Bass, H. W. & Buckler, E. S. Open chromatin reveals the functional maize genome. *Proc Natl Acad Sci U S A* **113**, E3177-3184, doi:10.1073/pnas.1525244113 (2016).

- 7 Hufford, M. B. *et al.* Comparative population genomics of maize
domestication and improvement. *Nat Genet* **44**, 808-811,
doi:10.1038/ng.2309 (2012).
- 8 Chia, J. M. *et al.* Maize HapMap2 identifies extant variation from a genome in
flux. *Nat Genet* **44**, 803-807, doi:10.1038/ng.2313 (2012).
- 9 Wei, F. *et al.* The physical and genetic framework of the maize B73 genome.
PLoS Genet **5**, e1000715, doi:10.1371/journal.pgen.1000715 (2009).
- 10 Ganal, M. W. *et al.* A large maize (*Zea mays* L.) SNP genotyping array:
development and germplasm genotyping, and genetic mapping to compare
with the B73 reference genome. *PLoS One* **6**, e28334,
doi:10.1371/journal.pone.0028334 (2011).
- 11 Gent, J. I., Wang, K., Jiang, J. & Dawe, R. K. Stable Patterns of CENH3
Occupancy Through Maize Lineages Containing Genetically Similar
Centromeres. *Genetics* **200**, 1105-1116, doi:10.1534/genetics.115.177360
(2015).
- 12 Schneider, K. L., Xie, Z., Wolfgruber, T. K. & Presting, G. G. Inbreeding drives
maize centromere evolution. *Proc Natl Acad Sci U S A* **113**, E987-996,
doi:10.1073/pnas.1522008113 (2016).
- 13 Wang, B. *et al.* Unveiling the complexity of the maize transcriptome by single-
molecule long-read sequencing. *Nature communications* **7**, 11708,
doi:10.1038/ncomms11708 (2016).
- 14 Schnable, J. C., Springer, N. M. & Freeling, M. Differentiation of the maize
subgenomes by genome dominance and both ancient and ongoing gene loss.
Proc Natl Acad Sci U S A **108**, 4069-4074, doi:10.1073/pnas.1101368108
(2011).
- 15 Lai, J. *et al.* Gene loss and movement in the maize genome. *Genome research*
14, 1924-1931, doi:10.1101/gr.2701104 (2004).
- 16 Fluhr, R., Lampl, N. & Roberts, T. H. Serpin protease inhibitors in plant
biology. *Physiologia plantarum* **145**, 95-102, doi:10.1111/j.1399-
3054.2011.01540.x (2012).
- 17 Francis, S. E., Ersoy, R. A., Ahn, J. W., Atwell, B. J. & Roberts, T. H. Serpins in
rice: protein sequence analysis, phylogeny and gene expression during
development. *BMC genomics* **13**, 449, doi:10.1186/1471-2164-13-449
(2012).
- 18 McClintock, B. The origin and behavior of mutable loci in maize. *Proc Natl
Acad Sci U S A* **36**, 344-355 (1950).
- 19 Slotkin, R. K. & Martienssen, R. Transposable elements and the epigenetic
regulation of the genome. *Nature reviews. Genetics* **8**, 272-285,
doi:10.1038/nrg2072 (2007).
- 20 SanMiguel, P. *et al.* Nested retrotransposons in the intergenic regions of the
maize genome. *Science* **274**, 765-768 (1996).
- 21 Brunner, S., Fengler, K., Morgante, M., Tingey, S. & Rafalski, A. Evolution of
DNA sequence nonhomologies among maize inbreds. *The Plant cell* **17**, 343-
360 (2005).

- 22 Sharma, A., Schneider, K. L. & Presting, G. G. Sustained retrotransposition is mediated by nucleotide deletions and interelement recombinations. *Proceedings of the National Academy of Sciences* **105**, 15470-15474 (2008).
- 23 Baucom, R. S. *et al.* Exceptional diversity, non-random distribution, and rapid evolution of retroelements in the B73 maize genome. *PLoS Genet* **5**, e1000732 (2009).
- 24 Buckler, E. S., Gaut, B. S. & McMullen, M. D. Molecular and functional diversity of maize. *Current opinion in plant biology* **9**, 172-176, doi:10.1016/j.pbi.2006.01.013 (2006).
- 25 Dooner, H. K. & He, L. Maize genome structure variation: interplay between retrotransposon polymorphisms and genic recombination. *The Plant cell* **20**, 249-258, doi:10.1105/tpc.107.057596 (2008).
- 26 Saxena, R. K., Edwards, D. & Varshney, R. K. Structural variations in plant genomes. *Briefings in functional genomics* **13**, 296-307, doi:10.1093/bfgp/elu016 (2014).
- 27 McMullen, M. D. *et al.* Genetic properties of the maize nested association mapping population. *Science* **325**, 737-740, doi:10.1126/science.1174320 (2009).
- 28 Strable, J. & Scanlon, M. J. Maize (*Zea mays*): a model organism for basic and applied research in plant biology. *Cold Spring Harbor protocols* **2009**, pdb emo132, doi:10.1101/pdb.emo132 (2009).
- 29 Swanson-Wagner, R. A. *et al.* Pervasive gene content variation and copy number variation in maize and its undomesticated progenitor. *Genome research* **20**, 1689-1699, doi:10.1101/gr.109165.110 (2010).
- 30 Morrell, P. L., Buckler, E. S. & Ross-Ibarra, J. Crop genomics: advances and applications. *Nature reviews. Genetics* **13**, 85-96, doi:10.1038/nrg3097 (2011).
- 31 Berlin, K. *et al.* Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nature biotechnology* **33**, 623-630, doi:10.1038/nbt.3238 (2015).
- 32 Chin, C. S. *et al.* Phased diploid genome assembly with single-molecule real-time sequencing. *Nat Methods* **13**, 1050-1054, doi:10.1038/nmeth.4035 (2016).
- 33 Clark, K., Karsch-Mizrachi, I., Lipman, D. J., Ostell, J. & Sayers, E. W. GenBank. *Nucleic Acids Res* **44**, D67-72, doi:10.1093/nar/gkv1276 (2016).
- 34 Pendleton, M. *et al.* Assembly and diploid architecture of an individual human genome via single-molecule technologies. *Nat Methods* **12**, 780-786, doi:10.1038/nmeth.3454 (2015).
- 35 Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome biology* **5**, R12, doi:10.1186/gb-2004-5-2-r12 (2004).
- 36 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754-1760, doi:10.1093/bioinformatics/btp324 (2009).
- 37 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-2079, doi:10.1093/bioinformatics/btp352 (2009).

- 38 Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS One* **9**, e112963, doi:10.1371/journal.pone.0112963 (2014).
- 39 Ellinghaus, D., Kurtz, S. & Willhoeft, U. LTRharvest, an efficient and flexible software for de novo detection of LTR retrotransposons. *BMC bioinformatics* **9**, 18, doi:10.1186/1471-2105-9-18 (2008).
- 40 Steinbiss, S., Willhoeft, U., Gremme, G. & Kurtz, S. Fine-grained annotation and classification of de novo predicted LTR retrotransposons. *Nucleic Acids Res* **37**, 7002-7013, doi:10.1093/nar/gkp759 (2009).
- 41 Wenke, T. *et al.* Targeted identification of short interspersed nuclear element families shows their widespread existence and extreme heterogeneity in plant genomes. *The Plant cell* **23**, 3117-3128, doi:10.1105/tpc.111.088682 (2011).
- 42 Rho, M. & Tang, H. MGEScan-non-LTR: computational identification and classification of autonomous non-LTR retrotransposons in eukaryotic genomes. *Nucleic Acids Res* **37**, e143, doi:10.1093/nar/gkp752 (2009).
- 43 Han, Y. & Wessler, S. R. MITE-Hunter: a program for discovering miniature inverted-repeat transposable elements from genomic sequences. *Nucleic Acids Res* **38**, e199, doi:10.1093/nar/gkq862 (2010).
- 44 Xiong, W., He, L., Lai, J., Dooner, H. K. & Du, C. HelitronScanner uncovers a large overlooked cache of Helitron transposons in many plant genomes. *Proc Natl Acad Sci U S A* **111**, 10263-10268, doi:10.1073/pnas.1410068111 (2014).
- 45 Campbell, M. S. *et al.* MAKER-P: a tool kit for the rapid creation, management, and quality control of plant genome annotations. *Plant Physiol* **164**, 513-524, doi:10.1104/pp.113.230144 (2014).
- 46 Soderlund, C. *et al.* Sequencing, mapping, and analysis of 27,455 maize full-length cDNAs. *PLoS Genet* **5**, e1000740, doi:10.1371/journal.pgen.1000740 (2009).
- 47 Law, M. *et al.* Automated update, revision, and quality control of the maize genome annotations using MAKER-P improves the B73 RefGen_v3 gene models and identifies new genes. *Plant Physiol* **167**, 25-39, doi:10.1104/pp.114.245027 (2015).
- 48 Mostovoy, Y. *et al.* A hybrid approach for de novo human genome sequence assembly and phasing. *Nat Methods*, doi:10.1038/nmeth.3865 (2016).
- 49 Cao, H. *et al.* Rapid detection of structural variation in a human genome using nanochannel-based genome mapping technology. *GigaScience* **3**, 34, doi:10.1186/2047-217X-3-34 (2014).
- 50 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357-359, doi:10.1038/nmeth.1923 (2012).

Table 1. Assembly statistics of the maize B73 RefGen_v4 genome.

	No. of contigs (scaffolds)	Mean length (Mb)	N50 size (Mb)	Max length (Mb)	Total assembly length (Mb)
Original genome maps	1,342	1.57	2.47	12.43	2,107
Original contigs from sequence assembly	3,303	0.64	1.04	5.65	2,105
Curated genome maps	1,356	1.56	2.47	12.47	2,114
Curated contigs from sequence assembly	2,958	0.71	1.18	7.26	2,104
Genome maps in hybrid scaffolds	1,287	1.62	2.49	12.47	2,080
Contigs in hybrid scaffolds	2,696	0.77	1.19	7.26	2,075
Hybrid scaffolds	356	5.97	9.73	38.53	2,075
Hybrid scaffolds + non-scaffolded contigs	625	3.45	9.56	38.53	2,105

Table 2. Summary of alignments and structural variations called from the two genome maps.

	Ki11 map vs. B73 RefGen_v4	W22 maps vs. B73 RefGen_v4
Total size of genome map (Mb)	2,216	2,280
Map Aligned to reference genome (bp)	721,910,526	893,468,537
Reference genome covered by map (bp)	694,392,730	904,369,945
Region in B73 with insertion and deletion (bp)	223,253,038	221,057,903
Ratio of region with insertion and deletion	32.15%	27.23%
No. of insertions	1,794	1,616
Average insertion size (bp)	21,510	21,470
No. of deletions	1,701	1,599
Average deletion size (bp)	18,340	20,120
Number of deletion regions potentially effecting genes	636	621

Figure 1. Genome assembly layout. **A)** Workflow for genome construction; **B)** Ideograms of maize B73 version 4 reference pseudomolecules. There are totally 2,522b gaps in the pseudomolecules. The top row of orange rectangles are the 1,115 gaps with size estimation from the genome map, the purple rectangles represent gaps without size estimation. The light gray rectangles in second row correspond to the contigs larger than 1Mb and the dark gray rectangles are the contigs smaller than 1Mb. Over half of the genome is represented by the contigs above 1Mb.

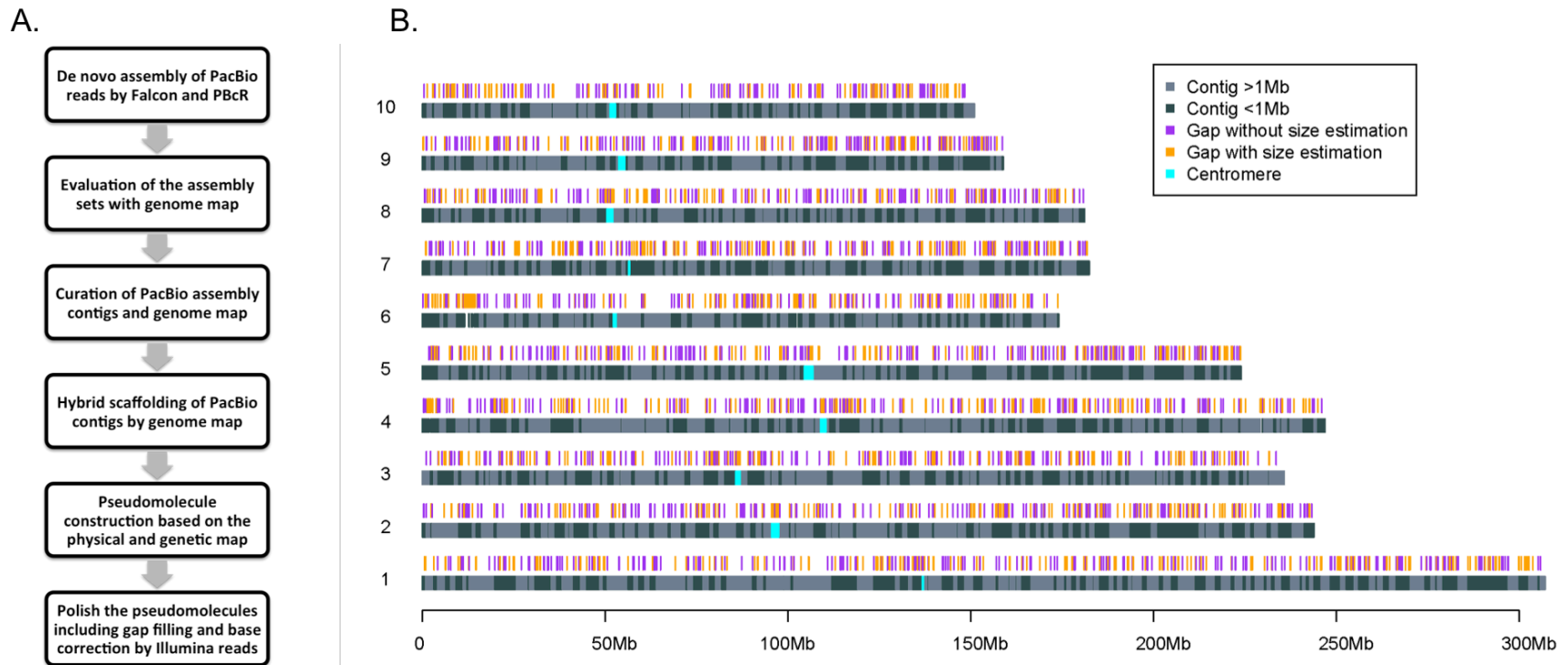


Figure 2. Phylogeny of maize and sorghum LTR retrotransposon families. Both **A) Ty3/Gypsy** and **B) Ty1/Copia** superfamilies are present at higher copy number in maize (red) than in sorghum (blue). Bars (\log_{10} -scaled) depict family copy numbers.

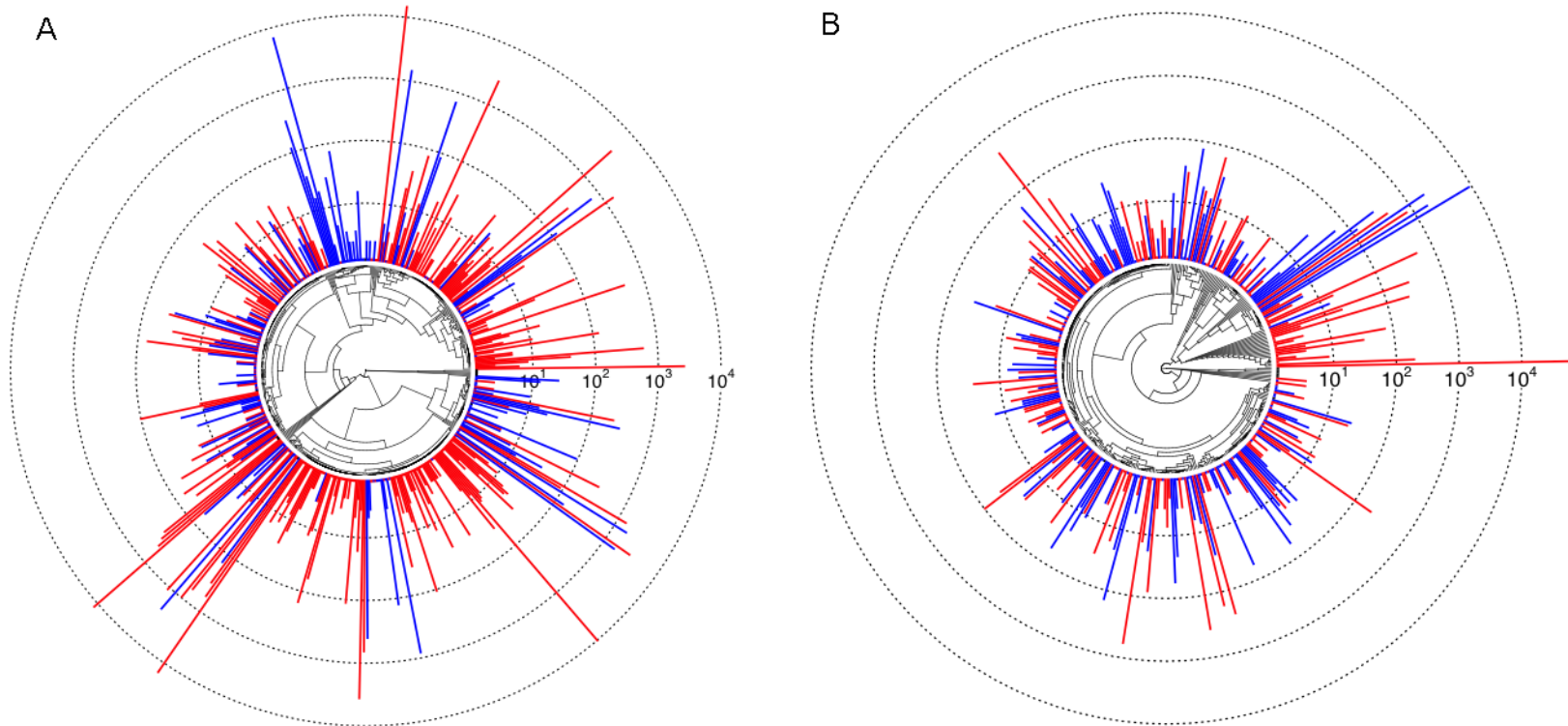


Figure 3. Structure variation from Ki11 and W22. **A).** Alignment and structural variation called from Ki11 and W22 genome map on chromosome 10, **B).** Size distribution of the insertion and deletions in Ki11 and W22, **C).** Example of using short read alignment to determine the missing region in Ki11 and W22.

