


SHORT REPORT

Open Access



Gapless assembly of maize chromosomes using long-read technologies

Jianing Liu¹, Arun S. Seetharam², Kapeel Chougule³, Shujun Ou⁴, Kyle W. Swentowsky⁵, Jonathan I. Gent⁵, Victor Llaca⁶, Margaret R. Woodhouse⁷, Nancy Manchanda⁴, Gernot G. Presting⁸, David A. Kudrna⁹, Magdy Alabady^{5,10}, Candice N. Hirsch¹¹, Kevin A. Fengler⁶, Doreen Ware^{3,12}, Todd P. Michael¹³, Matthew B. Hufford⁴ and R. Kelly Dawe^{1,5*} 

* Correspondence: kdawe@uga.edu

¹Department of Genetics, University of Georgia, Athens, GA 30602, USA
⁵Department of Plant Biology, University of Georgia, Athens, GA 30602, USA

Full list of author information is available at the end of the article

Abstract

Creating gapless telomere-to-telomere assemblies of complex genomes is one of the ultimate challenges in genomics. We use two independent assemblies and an optical map-based merging pipeline to produce a maize genome (B73-Ab10) composed of 63 contigs and a contig N50 of 162 Mb. This genome includes gapless assemblies of chromosome 3 (236 Mb) and chromosome 9 (162 Mb), and 53 Mb of the Ab10 meiotic drive haplotype. The data also reveal the internal structure of seven centromeres and five heterochromatic knobs, showing that the major tandem repeat arrays (CentC, knob180, and TR-1) are discontinuous and frequently interspersed with retroelements.

Keywords: Gapless assembly, Maize genome, Knob structure, Meiotic drive, Long-read technology

Background

Maize is a classic genetic model, known for its excellent chromosome cytology and rich history of transposon research [1]. Transposons make up the majority of the maize genome [2], and their accumulation over millions of years has driven genes far apart from each other and separated genes from their regulatory sequences [3]. There are also large inversions and other structural variations that contribute to fitness [4, 5] and significant variation in genome size caused by tandem repeat arrays [6]. Understanding this remarkable structural diversity is important for the continued improvement of maize, but the high repeat content has impeded progress [2, 5]. Here, we describe an automated assembly merging approach that yields gapless maize chromosomes and dramatically improves contiguity throughout the genome, including centromere and knob regions.

The most challenging genomic regions to assemble are tandem repeat arrays that exceed the read length of the current sequencing technologies. In most eukaryotes, these arrays are enriched in centromeres and ribosomal DNA (rDNA). Maize contains a



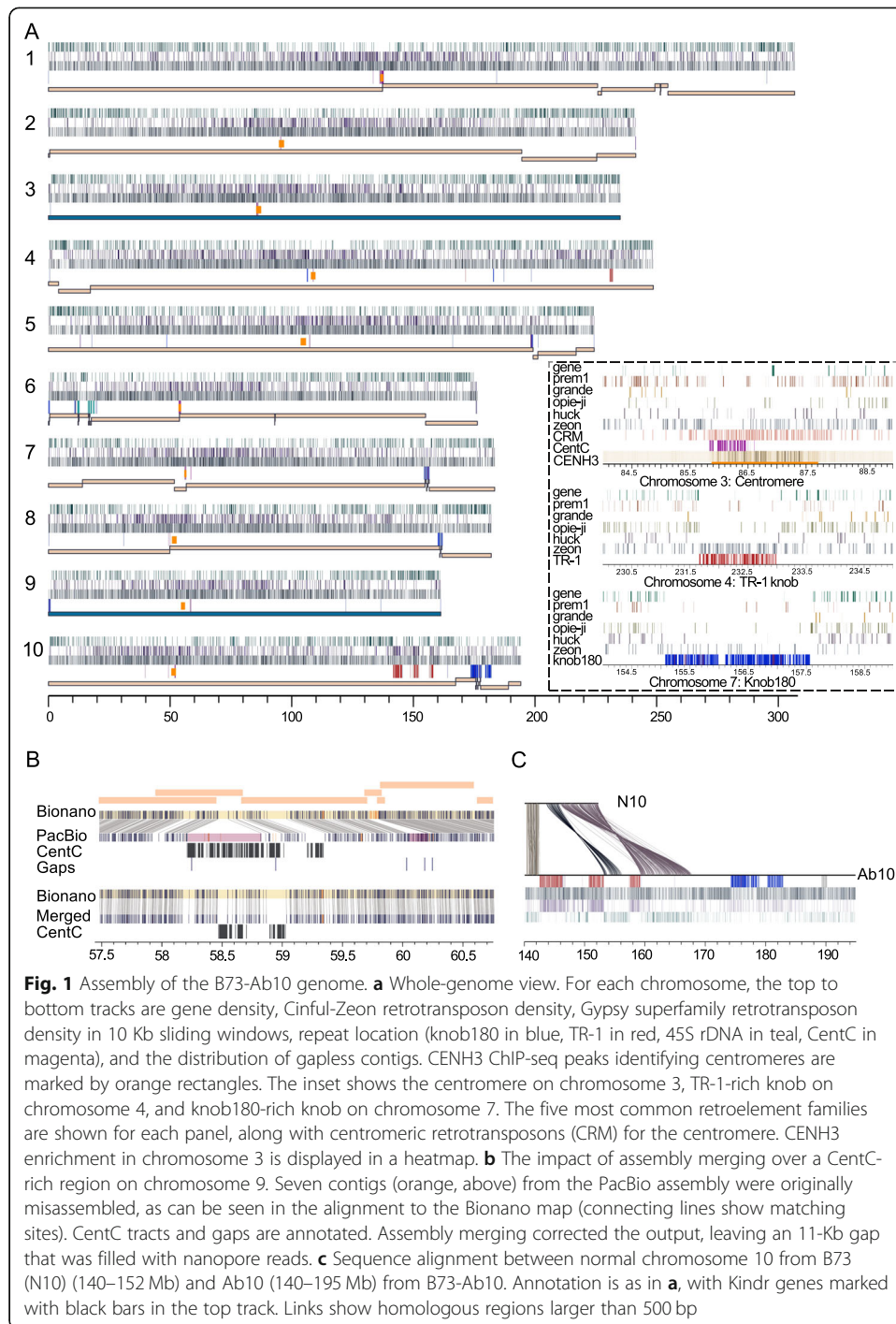
© The Author(s). 2020 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

centromeric repeat of 156 bp [7], a 45S rDNA repeat of 9349 bp, and a 5S rDNA repeat of 341 bp. In addition, maize contains two abundant classes of knob repeats that are found on chromosome arms, the major knob180 repeat (180 bp) [8] and the minor TR-1 repeat (~ 360 bp) [9]. Knob repeats occur in arrays that extend into the tens of megabases and present a significant barrier to full genome assembly. In most maize lines, knobs appear as inert heterochromatic bulges [8], but in lines with a meiotic drive system on Abnormal chromosome 10 (Ab10), they have centromere-like properties and are preferentially segregated to progeny [10]. Ab10 is considerably longer than chromosome 10 and contains two inversions [11], three knobs, and long spans of uncharacterized DNA that include a cluster of *Kinesin driver (Kindr)* genes required for meiotic drive [9]. Meiotic drive systems have been documented in many organisms and often lie within large inversions that contain novel repeat arrays [12], yet no meiotic drive haplotype has been fully sequenced and assembled.

Results and discussion

A new maize inbred, B73-Ab10, was created by backcrossing a line containing Ab10 to the B73 inbred six times and selfing it an additional five times (BC_6F_5). The B73-Ab10 inbred differs from B73 by the end of chromosome 10L which carries the Ab10 haplotype, the end of chromosome 9S which carries a kernel color gene necessary to score meiotic drive, and a 13-Mb internal section of chromosome 6 (coordinates between ~ 155 and 169 Mb). We used DNA from this line to prepare an optical map with the BioNano Saphyr system and sequenced it to high coverage using both PacBio and Nanopore technologies. We then implemented a genome assembly workflow based around the optical map (Additional file 1: Fig. S1). Briefly, the PacBio data were assembled using Canu [13], the Nanopore data assembled using miniasm [14], and the two independent assemblies merged with miniasm and integrated with the optical map as hybrid scaffolds. Hybrid scaffolds were then used to guide further gap closing and create a pseudomolecule assembly (Fig. 1a). Our approach of one-step contig merging and error correction using optical maps as a reference differs from other methods that rely on local assemblies to fill gaps and correct errors [15, 16]. While PacBio provided an overall superior assembly, it tended to fail in large repetitive regions (Additional file 1: Fig. S2A, B) and heterozygous areas (Additional file 1: Fig. S2C) where the Nanopore assembly succeeded due to a longer read length distribution. This was particularly evident in TR-1, knob180, and subtelomeric arrays as well as other tandemly duplicated regions (Additional file 1: Fig. S2B). Alignments of the optical map to the independent assemblies [17] and standard genome completeness measures demonstrate that the approach is highly accurate (Additional file 1: Table S1 and 2).

The final assembly has a contig N50 of 162 Mb (Table 1), which far exceeds the contiguity of any prior maize genome assembly [2, 5]. Of particular note is the complete 236-Mb assembly of chromosome 3, which was assembled gaplessly without manual intervention—a first for any chromosome from a large complex genome. While the human X-chromosome was also assembled gaplessly [18], this outcome required extensive manual inspection and correction. The entire B73-Ab10 genome is represented by 63 contigs where 90% are longer than 20.4 Mb (the N90). In addition to the expected gaps in repeat arrays, there were two gaps associated with residual heterozygosity on chromosome 9. Regions of heterozygosity reduce effective coverage and lead to



assembly chimeras that are broken during hybrid scaffolding. We filled these heterozygosity-associated breaks by choosing the dominant Bionano path and performing local assemblies over the gaps. Nanopore reads were also used to span a gap within a CentC array to complete the chromosome 9 telomere-to-telomere assembly. Aside from these manual interventions, some efforts to manually improve within-knob assemblies, and a correction to the *Kindr* gene complex region of Ab10, the assembly was automated. Our success in assembling chromosomes 3 and 9 can be attributed to the fact

Table 1 Assembly metrics of the B73-Ab10 genome

	Contigs			Pseudomolecules		
	N50 (Mb)	N90 (Mb)	Max Size (Mb)	Contig Number	Total Length (Mb)	Gap ^a Length (Mb)
Nanopore	2.0	0.5	8.3	1673	2161.1	93.2
PacBio	41.2	7.1	156.3	216	2162.7	2.6
Merged	162.0	20.4	235.9	63	2162.8	1.3

^aGaps longer than 10 Ns

that these chromosomes have the fewest cytologically visible repeat arrays [19]. All remaining gaps in the assembly are marked at the edges by tandem repeats (Fig. 1, Additional file 1: Fig. S2D).

Seven of the ten functional centromeres as defined by ChIP-seq of CENP-A/CENH3 [7] were assembled without gaps (Additional file 1: Table S3). Alignment of partial BAC-based assemblies of B73 centromeres showed excellent agreement overall (Additional file 1: Fig. S3). Only a subset of maize centromeres are composed of long CentC arrays, and even within those arrays, the majority of reads (65%) can be uniquely mapped, reflecting a high degree of sequence polymorphism (Additional file 1: Table S3 and S4). Three centromeres have no CentC at all and are composed of transposons of different forms. These include known centromeric retroelements (CRM) [7] as well as other common retrotransposons. We found no tendency for CENH3 to interact with CentC and CRM over any of the other repeats present (Additional file 1: Table S4). The lack of sequence specificity can be seen on centromere 3, where CENH3 localized over a 771-kb CentC array as well as a variety of other transposons in flanking sequence (Fig. 1a, inset).

Prior maize assemblies have succeeded in obtaining only small fragments of knob repeat arrays. In contrast, a knob180-rich knob on chromosome 9 (850 Kb), a TR-1-rich knob on chromosome 4 (1.3 Mb), and three TR-1-rich knobs (4.2 Mb, 2.6 Mb, and 2.1 Mb) on Ab10 were fully assembled in the B73-Ab10 assembly (Fig. 1a and Additional file 1: Table S5). The data show that knobs, like centromeres [2, 7], often contain more transposons than tandem repeats (Fig. 1c). Centromeric retrotransposons target areas with CENP-A/CENH3 [2, 7] and occupy on average 31.9% of functional centromeres, including within CentC arrays (Fig. 1a and Additional file 1: Table S3 and S6). The new knob assemblies reveal that the Cinfu-Zeon family of *Gypsy* elements [20] preferentially target knobs. Cinfu-Zeon elements occupy 27.0% of the assembled TR-1-rich knobs and 8.2% of the knob180-rich knobs, but only 3.8% percent of CentC arrays (Fig. 1a and Additional file 1: Table S6 and S7). Cinfu-Zeon elements are also abundant in other heterochromatic regions throughout the genome (Fig. 1a).

In addition to revealing the internal structure of knobs, the data provide the first complete view of the Ab10 haplotype that provides the selective force for the accumulation and maintenance of knobs [10]. The meiotic drive haplotype on Ab10 contains three fully assembled TR-1 knobs, a much larger knob180 knob that was not assembled, and two large inversions (4.4 and 8.3 Mb) that are homologous to normal chromosome 10 (Fig. 1c). These major structural differences help to explain why recombination between the Ab10 haplotype and normal chromosome 10 is suppressed

[21]. Ab10 also contains 22.4 Mb of novel sequence with no synteny to other regions of the maize genome or related grass genomes. Within this domain is the complete cluster of nine *Kindr* genes that are integral components of the drive system [10], as well as hundreds of other expressed genes, many of which have only one exon or overlap with transposons and are likely non-functional (Additional file 1: Table S8). Additional meiotic drive functions associated with the movement of knobs at meiosis and their delivery to egg cells [22] remain to be identified in this newly discovered sequence.

Conclusions

Gapless genome assemblies remove all uncertainty about the order, spacing, and orientation of genes and their regulators. We have shown that this can be achieved using long reads and well-known assembly algorithms, with significant improvements in contiguity obtained by integrating independent assemblies around an optical map scaffold. Given that most contigs end in telomeres, centromeres, or knobs, we presume that virtually all of the genes and associated regulatory information are represented in this genome assembly. The assembly merging pipeline also revealed the internal structure of repetitive domains that were previously known only by cytological techniques, thereby opening these regions to annotation and future epigenomic profiling. Similar results should be achievable for other complex genomes, although higher sequence coverage, longer reads, and/or additional scaffolding information may be needed for species with polyploidy or higher levels of heterozygosity.

Methods

PacBio assembly

High molecular weight DNA was extracted from young leaves using the protocol of Doyle and Doyle [23] with minor modifications. Young maize leaves flash frozen at -80°C were ground to a fine powder in liquid N₂ followed by very gentle extraction in CTAB buffer (that included proteinase K, PVP-40, and beta-mercaptoethanol) for 1 h at 50°C . After centrifugation, the supernatant was gently extracted twice with 24:1 chloroform:isoamyl alcohol. The upper phase was adjusted to 1/10th volume with 3 M KAc, gently mixed, and DNA precipitated with isopropanol. DNA was collected by centrifugation, washed with 70% EtOH, air dried for 20 min, and dissolved thoroughly in $1\times$ TE at room temperature.

Sequencing libraries were constructed following PacBio's template prep protocols (Procedure & Checklist—Preparing gDNA Libraries Using the SMRTbell Express Template Preparation Kit 2.0, PN 101-693-800 Version 01) for the Express Template Prep Kit 2.0 (Cat# 100-939-900) and sequenced using Sequel SMRTLink V5.1 and Sequel binding and sequencing chemistry v2.1. The longest 50X out of 62X PacBio raw sequences were error-corrected using falcon_kit pipeline v0.7 [24] without repeat masking by TANmask and REPmask (-e 0.75 -l 3000 --min_cov 2 --max_n_read 200). The error-corrected reads (43X, N50 = 22.3 Kb) were then trimmed and assembled with Canu [13] (v1.8) with the following parameters: correctedErrorRate=0.065 corMhap-Sensitivity=normal ovlMerThreshold=500 utgOvlMerThreshold=150. The read error correction process that is necessary for PacBio assembly may have homogenized some repeats and limited the assembly in long repeat regions. The accuracy of the Canu-

generated contigs was increased by aligning the raw PacBio reads to the assembly using pbmm2 (v1.2.0) from pb-assembly [24] and running the PacBio consensus algorithm tool Arrow (v2.3.3) (<https://github.com/PacificBiosciences/GenomicConsensus>) with default parameters to generate sequenced polished contigs. The contig assembly was further polished using 73X PE150 Illumina sequence by first aligning the reads to the Arrow polished assembly using minimap2 [25], followed by running the assembly tool Pilon [26] (v1.22) to correct individual base errors and small indels using the following parameters: --fix bases --minmq 30.

Nanopore assembly

Two different DNA extraction methods were used to generate high molecular weight (HMW) DNA for Oxford Nanopore (ONT) sequencing. CTAB DNA was prepared as described above for the PacBio assembly. Nuclear DNA was prepared using the protocol of Luo and Wing [27] with minor modifications. Young leaves flash frozen at -80°C were ground with liquid nitrogen and incubated with NIB buffer (10 mM Tris-HCL, PH8.0, 10 mM EDTA PH8.0, 100 mM KCL, 0.5 M sucrose, 4 mM spermidine, 1 mM spermine) on ice for 15 min. After filtration through miracloth, Triton X-100 (Sigma) was added to tubes at a 1:20 ratio, placed on ice for 15 min, and centrifuged to collect nuclei. Nuclei were washed with NIB buffer (containing Triton X-100) and re-suspended in 40 ml of the same buffer and centrifuged again. After removal of all liquid, 10 ml of Qiagen G2 buffer was added followed by gentle resuspension of nuclei; then, 30 ml G2 buffer with RNase A (to a final concentration of 50 mg/ml) was added. Tubes were incubated at 37°C for 30 min. Proteinase K (Invitrogen), 30 mg, was added and incubated at 50°C for 2 h followed by centrifugation for 15 min at 8000 rpm, at 4°C , and the liquid gently poured into a new tube. After gentle extraction with chloroform:isoamyl alcohol (24:1), DNA was precipitated with two thirds volume isopropanol. The DNA pellet was washed with 70% EtOH, air dried for 20 min, and dissolved in TE at room temperature.

DNA from both the CTAB and nuclear prep was used to generate either a rapid (SQK-RAD004) or one-dimensional (1D; SQK-LSK109) sequencing library for ONT. The resulting libraries were run on either a MinION or a GridION sequencer running for 48 h. All bases were called on the GridION using Guppy (v2.1.3), and the resulting fastq files were used for genome assembly. A total of 121 Gb ($\sim 50\times$) of ONT sequence was generated over 27 MinION R9.4 flowcells. The data were filtered for reads > 10 Kb using seqtk (<https://github.com/lh3/seqtk>), resulting in an estimated $30\times$ coverage ($N50 = 29,311$ bp) of the maize genome. The resulting uncorrected reads were aligned (overlap) with minimap2 (v2.13;-x ava-ont -t 64) [25], and an assembly graph (layout) was generated with miniasm (v0.3; -f <reads> <overlaps>) [14]. The resulting graph was inspected using Bandage [28]. The fact that the Nanopore assembly was carried out with uncorrected reads may have contributed to its better performance in long repeat regions (Additional file 1: Fig. S2). A consensus genome assembly was generated by mapping reads > 10 Kb to the assembly with minimap2 and then running racon (v1.3.1) [29]; the consensus process was repeated three times. The contig assembly was further polished using 73X PE150 Illumina sequence by first aligning the reads to the consensus assembly using minimap2 [14] followed by running the assembly tool Pilon (v1.18) [26] two times using 73X PE150 Illumina sequence.

Optical map assembly

Ultra-high molecular weight DNA was isolated from maize seedlings using a modified version of the Bionano Genomics Plant Tissue DNA Isolation Base protocol. Approximately 0.5 g of healthy aerial tissue was collected from young B73-Ab10 etiolated seedlings grown in soil-free conditions for 2 weeks. The leaves were treated with a 2% formaldehyde Bionano fixing solution, washed, chopped, and homogenized using a Qia-gen TissueRuptor in homogenization buffer. Free nuclei were pelleted at 2000×g, washed, isolated by gradient centrifugation, and embedded in a low melting point agarose plug. The nuclei were lysed by treating with proteinase K and RNase A treatments as described previously [30], and washed four times in Wash Buffer and five times in TE buffer. The purified high molecular weight nuclear DNA was recovered by melting the plug, digesting it with agarase, and subjecting the resulting sample to drop dialysis against TE.

The Bionano Saphyr platform was used in combination with the Direct Label and Stain (DLS) process to generate chromosome-level sequence scaffolds [31]. Direct labeling was performed using the Direct Labeling and Staining Kit (Bionano Genomics, San Diego CA) according to the manufacturer's protocol, except that 1 µg of DNA was used and DNA Stain was added to a final concentration of 1 µl per 0.1 µg of final DNA. The labeled sample was loaded into a Saphyr chip, and molecules separated, imaged, and digitized using a Saphyr and Compute server. Data visualization, map assembly, and hybrid scaffold construction were performed using Bionano Access (v1.3) and Bionano Solve (v3.4.0). A subset of 1,580,077 molecules with a minimum size of 150 Kb and combined length of 424,488 Mb were assembled without pre-assembly using the non-haplotype, no-CMPR-cut parameters without extend-split.

Assembly merging and gap closing

We developed a pipeline to integrate independent contig assemblies and curate assembly errors using Bionano maps as an anchor. The pipeline consists of five steps: (1) conflict resolution, (2) assembly error curation, (3) contig merging, (4) hybrid assembly and contig overlap removal, and (5) manual curation and gap filling (Additional file 1: Fig. S1). The first four steps were automated. A gapless chromosome 3 was generated upon contig merging in the third step, and the complete assembly of chromosome 9 required manual curation. While contig merging with miniasm can be applied to any two sequence assemblies, the availability of de novo assembled Bionano maps is necessary to perform conflict-cutting in step 1, contig error correction in step 2, and hybrid scaffolding in step 4 of the pipeline.

Step 1: Conflicts between the optical map and DNA sequence assemblies were resolved using Bionano Solve software (<https://bionanogenomics.com/support-page/data-analysis-documentation/>). Sequence assembly can occasionally connect two regions that share a repetitive sequence but do not belong together (making a chimeric contig). These appear as conflicts between Bionano maps and sequence assemblies when they are aligned. Optical maps were aligned to in silico digested representations of the DNA sequence assemblies using RefAligner (v3.4.0), and conflicts identified with the AssignAlignType.pl script. Conflicts with a chimeric

quality score higher than the default threshold were split using `cut_conflicts.pl` (using default parameters from `optArguments_nonhaplotype_noES_DLE1_saphyr.xml`), and a sequence file was produced with custom script `cut_conflict_NGS.py`. Removing chimeric joins increases the chance of complementary contig merging in step 3.

Step 2: Assembly errors in the conflict-resolved PacBio contigs were identified and automatically curated with ONT contigs. In this step, PacBio and ONT contigs were aligned to rescaled optical maps and structural discrepancies detected using the structural variant calling pipeline from BionanoSolve (v3.4.0). Homozygous insertions and deletions with a confidence of at least 0.1 and size larger than 1 Kb were classified as true assembly errors in the PacBio contigs. On the condition that no structural discrepancies were found in the corresponding ONT contigs, the ONT contigs were used to replace the erroneous sequences in PacBio contigs using custom script `SV_fix.py`.

Step 3: ONT contigs were used to close gaps and improve contiguity of the PacBio contig assembly. ONT contigs were mapped to PacBio contigs with `minimap2` [25] (v2.13; `-k28 -w28 -A1 -B9 -O16,41 -E2,1 -z200 -g100000 -r100000 --max-chain-skip 100`), and overlap regions merged using `miniasm` [14] (v0.3; `-1 -2 -r0 -e1 -n1 -h250000 -g100000 -o25000`). This step creates PacBio/ONT hybrid contigs that are called unitigs. The unitigs were then combined with the remaining contigs from the PacBio backbone assembly to create a merged contig assembly. After this step, a gapless chromosome 3 was generated (a region of heterozygosity from 164.5 to 166.2 Mb on chromosome 3 was automatically resolved). The merged contigs were then aligned to Bionano maps, where overlaps between adjacent contigs were detected and merged with `minimap2` (v2.13) and `miniasm` (v0.3) using the custom script `Overlap_merge.py`. This step only identifies large overlaps (roughly > 200 Kb) that can be detected at the level of de novo Bionano label alignment. Identifying all overlaps, including smaller overlaps, requires hybrid scaffolding with the optical map (step 4). If proceeding to step 4, overlap merging in step 3 is optional.

Step 4: Bionano maps were integrated with the sequence contigs by hybrid scaffolding using the `hybridScaffold.pl` script from BionanoSolve (v3.4.0) with default parameters from `optArguments_nonhaplotype_noES_DLE1_saphyr.xml`. This step orders and orients sequence contigs and facilitates the resolution of remaining overlaps between contigs. As the optical maps are aligned and rescaled with the sequence maps repeatedly during hybrid scaffolding, more accurate overlaps between contigs are identified and annotated as 13N gaps. These overlaps were removed through contig merging with `miniasm` (v0.3), as described in step 3. Due to the extreme repetitiveness in the 45S rDNA repeat region on chromosome 6, both the contig assemblies and hybrid scaffolding in this area are erroneous. Therefore, we left the contigs in the NOR un-merged and marked the incorrectness with 13N gaps.

Step 5: Manual curation was performed to correct assembly errors, close gaps in repetitive and heterozygous regions, and assemble telomeres.

Repeat assembly manual curation. In highly repetitive regions, erroneous read joins at the tips of contigs were not detected as conflicts or assembly errors in steps 1 or 2

due to the limited resolution of Bionano alignment. In these regions, we trimmed and removed the unaligned regions to reveal eligible ends for overlap merging using *miniasm* (v0.3). These modifications extended the contiguity of repeat arrays at the edges of longer contigs. Contigs composed exclusively of knob and CentC repeats arrays lack pan-genome anchor markers and are not present in the pseudomolecules.

Chromosome 9 manual curation. Seven gaps, ranging from 2 to 236 Kb, were present in the chromosome 9 assembly after hybrid scaffolding. Two large gaps of 236 Kb and 41 Kb were caused by heterozygosity (76.29–76.80 Mb), one 21 Kb gap was due to repetitiveness in a CentC array (58.43–58.67 Mb), and the remaining four gaps were smaller than 7 Kb (two of these were in the 843-Kb knob on the tip of 9S). The four small gaps were first filled by running three iterations of LR Gapcloser (Sep 24, 2018 commit) [32] at default settings using PacBio error-corrected reads. To resolve the 236-Kb gap caused by heterozygosity, all contigs anchored to chromosome 9 were re-scaffolded using the longest chromosome 9 Bionano map as the sole anchor. This reduced the 236-Kb gap to 58 Kb. Local assemblies were run with *Flye* (v2.6) [33] using ONT reads surrounding gaps to fill the remaining 58-Kb and 41-Kb gaps. *Flye*-assembled contigs were integrated with the flanking contigs by unitigging with *miniasm* (v0.3) and aligned to Bionano maps for inspection. An 8-Kb gap remained, which was filled with a single ONT read that spans it. The gap in the CentC array was filled by manually selecting two long ONT reads (> 50 Kb) that spanned the gap, creating a consensus at the overlap and placing the resulting sequence in the gap.

Kindr complex manual curation. The assembly over the ~1-Mb tandem array of *Kindr* genes (each within an ~100-Kb repeat) was erroneous due to collapsing in the PacBio sequence contig and improper scaffolding. We manually selected the most contiguous ONT contig over this region, carried out hybrid scaffolding for the scaffold containing *Kindr*, placed an excluded contig in the correct area, and removed an overlap region through contig merging.

Telomere manual curation. Fifteen telomeres were assembled by extending the ends of scaffolds with the longest uniquely mapped ONT read that contained telomeric repeats TTTAGGG/CCCTAAA (≥ 1 Kb). The regions with newly assembled telomeres include 1L, 2L, 3S, 3L, 4S, 4L, 5L, 6L, 7S, 7L, 8S, 8L, 9S, 9L, and 10S.

The final scaffolds were polished with PacBio subreads using tools from *pb-assembly* [24]. Read alignment was performed with *pbbmm2* (v1.2.0), and polishing was executed with *GCpp* (v1.0.0) at default parameters. Scaffolds were further polished with 73X PE150 Illumina reads using *Pilon* (v1.23) with default parameters [26]. The error-corrected PacBio reads and Illumina reads often mapped incorrectly in highly repetitive regions (Additional file 1: Fig. S2B,C,D). Regions with excessive incorrect mapping are expected to be overpolished, whereas regions with few correctly mapped reads are expected to retain a higher frequency of sequencing errors.

AGP construction

The pseudomolecules were constructed from the hybrid scaffolds using *ALLMAPS* (v0.8.12) [34]. Both pan-genome anchor markers [35] and the IBM (Intermated B73 \times Mo17)

genetic map [36] were used with equal weights for ordering and orienting the scaffolds. Pan-genome anchor markers were obtained from the CyVerse Data commons [37] and processed to generate a bed file with 50 bp upstream and downstream of B73 V3 coordinates. The extracted markers were mapped to a HiSat2 (v2.1.0)^{29,30} indexed assembly of B73-Ab10 by disabling splicing (`--no-spliced-alignment`) and forcing global alignment (`--end-to-end`). Very high read and reference gap open and extension penalties (`--rdg 10000,10000` and `--rfg 10000,10000`) were also used to ensure full-length mapping of marker sequence. The final alignment was then filtered for mapping quality greater than 30 and tag XM:0 (unique mapping) to retain only high-quality, uniquely mapped marker sequences. The mapped markers were merged with the predicted distance information to generate a CSV input file for ALLMAPS. Only scaffolds with more than 20 uniquely mapped markers, with a maximum of 100 markers per scaffold, were used for pseudomolecule construction. The IBM genetic markers were downloaded from MaizeGDB (https://www.maizegdb.org/complete_map?id=887740) [38] and were processed to generate a bed file similar to pan-genome markers. For the markers with coordinates, 50-bp flanking regions were extracted from the B73 v4 genome. For markers without coordinates, marker sequences were used as-is, and those missing both coordinates and sequences were discarded. Mapping of the markers was done similar to the method described above for the pan-genome anchor markers, with all uniquely mapped markers retained. The genetic distance information for these markers was converted to a CSV file before use in ALLMAPS. ALLMAPS was run with default options, and the pseudomolecules were finalized after inspecting the marker placement plot and the scaffold directions. Of the 50 Bionano scaffolds anchored with sequence contigs, 26 with uniquely mapped genetic markers were included in the pseudomolecules. Among the 24 unplaced scaffolds with a total size of 19.4 Mb, 22 are composed entirely of knob180 and/or TR-1 arrays (17.7 Mb).

Comparing PacBio and Nanopore assemblies in repetitive and heterozygous regions

To determine how tandem repeats and regions of heterozygosity impacted the assemblies, we identified tandemly repeated areas by chromosome self-alignment with minimap2 (v2.17; `-PD -k19 -w19 -m200`) and heterozygous regions by manual inspection using Bionano Access software. PacBio gap coordinates were projected onto the final assembly using minimap2 (v2.17; `-cx asm5 --cs`), followed by coordinate liftover using paftools.js [25]. Gaps that were complemented by Nanopore contigs were identified as gaps present in the PacBio assembly but absent in the final assembly. The PacBio adjusted gap coordinates, complemented gaps, and final assembly gaps were mapped to tandem repeats and heterozygous regions with bedtools [39] (v2.28.0; `window -r 500000 -l 500000`). The co-occurrence of PacBio gaps with tandem repetitiveness and heterozygous regions was assessed by two-tailed Fisher's exact test using bedtools fisher (v2.28.0) at default settings.

To assess read coverage over gap areas, a total of 36.9X error-corrected PacBio reads (≥ 10 Kb), 20.7X error-corrected Nanopore reads (≥ 10 Kb), and 30X PE150 Illumina reads were mapped to the final assembly. Long-read mapping was performed using

minimap2 (v2.17) with default parameters, and short-read mapping was carried out with bwa (v0.7.17) at default settings. Read gap regions were defined as areas mapped with fewer than 3 reads for PacBio and Illumina datasets, and fewer than 2 reads for the Oxford Nanopore dataset. Basepair-level genome coverage was calculated with bedtools genomecov (v2.28.0; -bga), and regions with fewer reads than the cutoff were extracted. The length distributions of PacBio and Oxford Nanopore reads mapped to a tandem repeat (chr8: 31–33.5 Mb) and heterozygous area (chr3: 164–167.6 Mb) were obtained with SAMTools (v1.9).

RNA-seq

Ten tissues were sampled throughout development for evidence-based gene annotation including the following: (1) primary root and (2) coleoptile at 6 days after planting; (3) base of the 10th leaf, (4) middle of the 10th leaf, and (5) tip of the 10th leaf at the Vegetative 11 (V11) growth stage; (6) meiotic tassel and (7) immature ear at the V18 growth stage; (8) anthers at the Reproductive 1 (R1) growth stage; and (9) endosperm and (10) embryo at 16 days after pollination. For each tissue, two biological replicates were harvested, and each biological replicate was made up of tissue from three individual plants. Endosperm and embryo tissues were harvested from 50 kernels per plant (150 total per biological replicate). Tissues 1–5 above were collected from greenhouse-grown plants, and tissues 6–10 were from field-grown plants. Greenhouse-grown plants were planted in Metro-Mix300 (Sun Gro Horticulture) with no additional fertilizer and grown under greenhouse conditions (27 °C/24 °C day/night and 16 h/8 h light/dark) at the University of Minnesota Plant Growth Facilities. Field-grown plants were planted at the Minnesota Agricultural Experiment Station located in Saint Paul, MN, with 30-in. row spacing at ~52,000 plants per hectare. RNA was extracted using the Qiagen RNeasy plant mini kit following the manufacturer's suggested protocol.

Total RNA samples were assayed by Bioanalyzer to determine RNA integrity and normalized in 25 µl of nuclease-free water prior to library preparation. Sequencing libraries were prepared using KAPA's Stranded mRNA-seq kit (#KK4821) according to the manufacturer's instructions. The mRNA was enriched using oligo-dT beads, fragmented, and converted to double-stranded cDNA using random hexamer priming and amplification. Libraries were pooled at equimolar ratios and sequenced on NextSeq 500 instruments using the PE75 protocol.

Gene annotation

For evidence-based predictions, genome-guided transcript assemblies were generated from five different assemblers, viz, Trinity (v2.6.6) [40, 41], StringTie (v1.3.4a) [42], Strawberry (v1.1.1) [43], Cufflinks (v2.2.1) [42, 44], and Class2 [42, 44, 45], and the best set of transcripts were identified and annotated as genes using Mikado (v1.2.4) [46]. Briefly, the RNA-seq reads from each library were mapped to a STAR (v2.5.3a) [47] indexed B73-Ab10 genome using a 2-pass mapping approach (the initial round of alignments provides splice information for the subsequent round of mapping reads). Default options were used for mapping with few post-processing options enabled (print all SAM format attributes --outSAMattributes All; downstream compatibility --outSAMmapqUnique 10; and number of mis-matches

--outFilterMismatchNmax 0). Individually mapped RNA-seq libraries were then pooled, sorted, and indexed using SAMTools (v1.9) [48], for use with the transcript assembly programs. For all genome-guided transcriptome assemblers, default options were used except, if it allowed minimum transcript length setting, it was set to 100 bp (Trinity using --min_contig_length 100, StringTie using -m 100, and Strawberry using -t 100), and if it allowed RNAseq strandedness, it was set to stranded (Trinity using --SS_lib_type FR, Cufflinks using --library-type fr-firststrand). For Trinity, maximum intron size was also set to 10,000 (--genome_guided_max_intron 10000). All assemblers generated a GFF3 as the final output except for Trinity, for which assembled transcripts in fasta format were mapped back to the gmap (v2019-05-12) indexed genome to generate a GFF3 file (by setting the output format option -f to gff3_match_cdna). Portcullis (v1.1.2) [49] was used to generate a high confidence set of splice junctions for the B73-Ab10 genome from the merged mapped reads. Mikado was configured to use all transcript assemblies (with strandedness marked as True for all except for Trinity, and with equal weights), portcullis-generated splice sites, and a plants.yaml scoring matrix. Preliminary transcripts prepared by Mikado, through merging all transcripts and removing the redundant copies, were processed using TransDecoder (v5.5.0) [50] (to identify open reading frames) and blastx (v2.9.0) [40] against SwissProt viridiplantae proteins (for identifying full-length transcripts). Default options were used for TransDecoder, and for blastx, maximum target sequences were set to 5 (-max_target_seqs 5) and output format to xml (-outfmt 5). These were provided as input for Mikado for picking and annotating the best transcripts for each locus. The obtained GFF3 file was used to extract transcripts and proteins using the gffread utility from the Cufflinks package.

Additional structural improvements for the Mikado-generated transcripts were completed using the PASA (v2.3.3) [51] genome annotation tool. The inputs for PASA included 2,019, 896 maize EST derived from genbank, 83,087 Mikado transcripts, 69,163 B73 full-length cDNA from genbank, and 46,311 maize iso-seq transcripts from 11 developmental tissues that were filtered for intron retention [52]. PASA was run with default options, with a first step of aligning transcript evidence to the masked B73-Ab10 genome using GMAP (v.2018-07-04) [53] and Blat (v.36) [54]. The full-length cDNA and Iso-seq transcript IDs were passed in a text file (-f FL.acc.list) during the PASA alignment step. Valid near perfect alignments with 95% identity were clustered based on genome mapping location and assembled into gene structures that included the maximal number of compatible transcript alignments. PASA assemblies were then compared with B73-Ab10 Mikado transcript models using default parameters. PASA updated the models, providing UTR extensions, novel and additional alternative isoforms. PASA-generated models were passed through the MAKER-P (v3.0) [55] annotation pipeline as model_gff along with all the transcript and protein sequences to obtain Annotation Edit Distance (AED) [56] scores to assess the quality of annotations. Transposon element (TE) related genes were filtered using the TESorter tool [40, 57], which uses the REXdb (viridiplantae_v3.0 + metazoa_v3) database of TEs. Finally, the gene annotations were verified for translation errors using the EnsemblCompara pipeline [58].

BUSCO assessment

The gene space completeness of the B73-Ab10 genome assembly was assessed using the GenomeQC [59] tool, which provides a summary of the number of complete, fragmented, and missing Benchmarking Universal Single-Copy Orthologs (BUSCO) in the

assembly. The Embryophyta database (embryophyta_odb9; consisting of 1440 conserved, single-copy plant genes) and the genome assembly in the fasta file format were provided as input to the tool to calculate the BUSCO metrics.

TE annotation

The manually curated transposable element library (maizeTE11222019) derived from the Maize TE Consortium (MTEC; <https://github.com/oushujun/MTEC>) was used as the base TE library. Novel TEs of the maize Ab10 genome not included in the MTEC library were structurally identified using the EDTA pipeline (v1.6.5) [60] with parameters “-species maize -curatedlib maizeTE11222019.” The MTEC library augmented with Ab10-specific TEs was used to annotate TE fragments using RepeatMasker. Coding sequences of the maize B73 v4 assembly were downloaded from MaizeGDB and used to remove gene sequences in the EDTA-generated TE library. Whole-genome TE annotations were generated using the EDTA-augmented MTEC library (-anno 1). The LTR Assembly Index (LAI) [61] scores of genome assemblies were calculated using LAI (beta3.2) within the LTR_retriever (v2.8) [62] package with parameters “-iden 94.8550 -totLTR 76.34.”

Centromere and repeat analyses

The overall accuracy of the centromere assemblies was assessed by aligning previous BAC-based B73 centromere assemblies [37] to the B73-Ab10 genome using Bionano RefAligner (v3.4.0) with default parameters. Although the BAC-based assemblies do not traverse CentC arrays, there is excellent overall agreement in sequence and contiguity (Additional file 1: Fig. S3).

Active centromere locations were determined by identifying the CENH3 ChIP-seq-enriched regions in the final assembly using genomic reads as a control. The SE150 Illumina ChIP-seq reads were obtained from SRA (SRX2737618) [63], and the 73X PE150 Illumina genomic reads were subsampled to 30X with seqtk (<https://github.com/lh3/seqtk>). Both the ChIP-seq reads and the genomic reads were trimmed with Trim Galore (v0.4.5; <https://github.com/FelixKrueger/TrimGalore/>) with default parameters and aligned to the final assembly with BWA-MEM (v0.7.17) [64]. Epic2 [65] was employed to call peaks with the CENH3 ChIP-seq alignment set as treatment, genomic read alignment as control, MAPQ (mapping quality) as 20, effective genome size as 0.8, bin size as 5000, and gap size as 0. The effective genome size of the final genome was calculated as the fraction of unique 150-mers over total 150-mers using Jellyfish (v2.2.6) [66] (-m 150 -s 2193M -out-counter-len 1 -counter-len 1). The coordinates of active centromeres were identified as islands with a score above 250 and a fold change higher than 4.

The coordinates of repeat arrays were identified by blasting the knob180 and CentC consensus sequences [63], a TR-1 consensus (TTCTTTATATTCCTAACTTTTAA GCAACTGTATGGTGGAAAAAGGTGTCTTACAACCTTAACCTATGTTTGGACA GTTCTCTCGTGCAATTTGGCTAAATTTCCCATGGTCTTTATTTATTTTGGAG AAACGATGTGGTATAATGATGTGCGATGTTTTACTTGAGTGGACATAAACAC CATTAGGTATGCCTTGAATAGAGGGGATTATTGGAAACCTGGTATCACAAA AGGTCATTAGCTAGCCCAATAACGTCTTCATCCACTAGTTATACTCTAATAC CCTCTAGTGTGAATACAATGCCCAATATCATAGAAACGTCATTTGAGGT TAAAAGGTGATCTATTGTTTTGAA), subtelomeric repeat (NCBI CL569186.1), and

ribosomal DNA intergenic spacer sequences (NCBI AF013103.1) against the B73-Ab10 genome. Knobs were defined as repeat clusters (≥ 500 Kb) that are composed of at least 10% repeat consensus sequences (knob180 and TR-1) with no more than 100 Kb spacing between repeat units. This definition of knob180 knobs excludes the subtelomeric knob180 arrays. CentC arrays are defined as repeat clusters (≥ 100 Kb) that are composed of at least 10% CentC consensus sequences.

Non-overlapping repeat units were quantified in each repeat array with custom script `repeat_analyses.py`. Five major families of the long terminal repeat (LTR)-retrotransposons in knobs, CentC arrays, and active centromeres were individually quantified with `bedtools` (v2.28.0) [39]. The *Opie-Ji* family includes *Opie*, *Ji*, *Ruda*, and *Giepum*, and the *Prem1* family is composed of *Prem1*, *Xilon*, *Diguus*, and *Tekay* [67]. Centromeric retrotransposons CRM1 and CRM2 were quantified together and annotated as CRM in active centromeric regions.

To assess the enrichment of mappable repeat elements in functional centromeres, each of the elements was first classified into uniquely mappable or non-uniquely mappable groups. A cutoff of MAPQ20 was applied to the alignment file, and `bedtools` (v2.28.0) was used to estimate genome coverage at the base pair level (-bga). Non-uniquely mapped locations (≤ 2 or ≥ 101 aligned reads) were merged into islands with a maximum interval of 1 Kb. CENH3 ChIP-seq enrichment for the unique and non-unique fractions of CentC, CRM, and five major LTR retrotransposon families was then individually assessed. ChIP enrichment was calculated by normalizing ChIP-seq against the input genome-seq alignment bam files using a RPKM normalization method with `deepTools` (v3.2.1) [68]. Default options were used except for the following parameters: `--operation ratio --scaleFactorsMethod None --normalizeUsing RPKM`.

Supplementary information

Supplementary information accompanies this paper at <https://doi.org/10.1186/s13059-020-02029-9>.

Additional file 1: Figure S1. Workflow for the B73-Ab10 assembly pipeline. **Figure S2.** Complementation of PacBio assembly gaps by Nanopore contigs. **Figure S3.** The alignment of BAC-based assemblies of B73 centromeres to the merged assembly in optical map format. **Table S1.** Assembly statistics and gaps in B73-Ab10 assemblies.

Table S2. Accuracy of genome assemblies as assessed by comparison to Bionano maps. **Table S3.** Coordinates and composition of centromeres defined by CENH3 ChIP-seq in the B73-Ab10 assembly. **Table S4.** CENH3 enrichment and mappability of Illumina reads in active centromeres. **Table S5.** Repetitive components in B73-Ab10 assemblies. **Table S6.** Composition of CentC arrays. **Table S7.** Composition of knob180 and TR-1 knobs. **Table S8.** Gene and transposon distributions in the Ab10 haplotype and corresponding N10 regions.

Additional file 2. Review history.

Acknowledgements

We thank Shane Poplawski and Rachel Gominsky for help with Nanopore sequencing.

Peer review information

Yixin Yao was the primary editor on this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

Review history

The review history is available as Additional file 2.

Authors' contributions

RKD created the Ab10-B73 inbred. KAF, TPM, and VL carried out the PacBio assembly, Oxford Nanopore assembly, and Optical map assembly, respectively. JL developed the assembly merging pipeline, measured the assembly metrics, assessed the assembly complementation and accuracy, and analyzed the repeat space structure. AS constructed the pseudomolecules. AS and KC performed the gene annotation. SO carried out the transposable element annotation. KWS and JIG analyzed the gene and transposon distributions in the Ab10 haplotype. NM carried out the BUSCO analyses and metric assessment. MRW coordinated the data transfer and deposition. RKD, GGP, DAK, MA, CNH, DW, TPM,

and MBH supervised the project. JL made the figures. All authors helped write the manuscript. The authors read and approved the final manuscript.

Funding

This study was supported by the National Science Foundation grants MCB-1412063 and IOS-1744001.

Availability of data and materials

The B73-Ab10 inbred can be obtained as PI 690316 at the Germplasm Resources Information Network (GRIN), Ames, IA. All genomic sequence and Bionano data can be obtained at the NCBI SRA under Bioproject PRJEB35367 [69]. The RNA-seq data is deposited in EBI (accession number E-MTAB-8641) [70]. The code used in this study is available at the GitHub repository <https://github.com/dawelab/Ab10-Assembly> [71].

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹Department of Genetics, University of Georgia, Athens, GA 30602, USA. ²Genome Informatics Facility, Iowa State University, Ames, IA 50011, USA. ³Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724, USA. ⁴Department of Ecology, Evolution, and Organismal Biology, Iowa State University, Ames, IA 50011, USA. ⁵Department of Plant Biology, University of Georgia, Athens, GA 30602, USA. ⁶Corteva Agriscience™, 8325 NW 62nd Ave, Johnston, IA 50131, USA. ⁷USDA-ARS Corn Insects and Crop Genetics Research Unit, Ames, IA 50011, USA. ⁸Molecular Biosciences and Bioengineering, University of Hawaii, Honolulu, HI 96822, USA. ⁹Arizona Genomics Institute, School of Plant Sciences, University of Arizona, Tucson, AZ 85721, USA. ¹⁰Georgia Genomics and Bioinformatics Core Laboratory, University of Georgia, Athens, GA 30602, USA. ¹¹Department of Agronomy and Plant Genetics, University of Minnesota, St. Paul, MN 55108, USA. ¹²USDA ARS NAA Robert W. Holley Center for Agriculture and Health, Agricultural Research Service, Ithaca, NY 14853, USA. ¹³Informatics Department, J. Craig Venter Institute, La Jolla, CA, USA.

Received: 2 February 2020 Accepted: 23 April 2020

Published online: 20 May 2020

References

1. Nannas NJ, Dawe RK. Genetic and genomic toolbox of *Zea mays*. *Genetics*. 2015;199:655–69.
2. Jiao Y, Peluso P, Shi J, Liang T, Stitzer MC, Wang B, et al. Improved maize reference genome with single-molecule technologies. *Nature*. 2017;546:524–7.
3. Ricci WA, Lu Z, Ji L, Marand AP, Ethridge CL, Murphy NG, et al. Widespread long-range cis-regulatory elements in the maize genome. *Nat Plants*. 2019; Available from: <https://doi.org/10.1038/s41477-019-0547-0>.
4. Pyhäjärvi T, Hufford MB, Mezouk S, Ross-Ibarra J. Complex patterns of local adaptation in teosinte. *Genome Biol Evol*. 2013;5:1594–609.
5. Yang N, Liu J, Gao Q, Gui S, Chen L, Yang L, et al. Genome assembly of a tropical maize inbred line provides insights into structural variation and crop improvement. *Nat Genet*. 2019;51:1052–9.
6. Bilinski P, Albert PS, Berg JJ, Birchler JA, Grote MN, Lorant A, et al. Parallel altitudinal clines reveal trends in adaptive evolution of genome size in *Zea mays*. *PLoS Genet*. 2018;14:e1007162.
7. Wolfgruber TK, Sharma A, Schneider KL, Albert PS, Koo D-H, Shi J, et al. Maize centromere structure and evolution: sequence analysis of centromeres 2 and 5 reveals dynamic loci shaped primarily by retrotransposons. *PLoS Genet*. 2009; 5:e1000743.
8. Peacock WJ, Dennis ES, Rhoades MM, Pryor AJ. Highly repeated DNA sequence limited to knob heterochromatin in maize. *Proc Natl Acad Sci U S A*. 1981;78:4490–4.
9. Ananiev EV, Phillips RL, Rines HW. A knob-associated tandem repeat in maize capable of forming fold-back DNA segments: are chromosome knobs megatransposons? *Proc Natl Acad Sci U S A*. 1998;95:10785–90.
10. Dawe RK, Lowry EG, Gent JI, Stitzer MC, Swentowsky KW, Higgins DM, et al. A kinesin-14 motor activates neocentromeres to promote meiotic drive in maize. *Cell*. 2018;173:839–50.e18.
11. Mroczek RJ, Melo JR, Luce AC, Hiatt EN, Dawe RK. The maize Ab10 meiotic drive system maps to supernumerary sequences in a large complex haplotype. *Genetics*. 2006;174:145–54.
12. Dyer KA, Charlesworth B, Jaenike J. Chromosome-wide linkage disequilibrium as a consequence of meiotic drive. *Proc Natl Acad Sci U S A*. 2007;104:1587–92.
13. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res*. 2017;27:722–36.
14. Li H. Minimap and miniiasm: fast mapping and de novo assembly for noisy long sequences. *Bioinformatics*. 2016;32: 2103–10.
15. Du H, Liang C. Assembly of chromosome-scale contigs by efficiently resolving repetitive sequences with long reads. *Nat Commun*. 2019;10:5360.
16. Vollger MR, Dishuck PC, Sorensen M, Welch AE, Dang V, Dougherty ML, et al. Long-read sequence and assembly of segmental duplications. *Nat Methods*. 2019;16:88–94.
17. Udall J, Dawe RK. Is it ordered correctly? Validating genome assemblies by optical mapping. *Plant Cell*. 2017; Available from: <https://doi.org/10.1105/tpc.17.00514>.

18. Miga KH, Koren S, Rhie A, Vollger MR, Gershman A, Bzikadze A, et al. Telomere-to-telomere assembly of a complete human X chromosome. *bioRxiv*. 2019:735928 [cited 2019 Nov 3]. Available from: <https://www.biorxiv.org/content/10.1101/735928v3.abstract>.
19. Albert PS, Gao Z, Danilova TV, Birchler JA. Diversity of chromosomal karyotypes in maize and its relatives. *Cytogenet Genome Res*. 2010;129:6–16.
20. Sanz-Alferez S, SanMiguel P, Jin Y-K, Springer PS, Bennetzen JL. Structure and evolution of the Cinfu retrotransposon family of maize. *Genome*. 2003;46:745–52.
21. Rhoades MM. Preferential segregation in maize. *Genetics*. 1942;27:395–407.
22. Hiatt EN, Dawe RK. Four loci on abnormal chromosome 10 contribute to meiotic drive in maize. *Genetics*. 2003;164:699–709.
23. Doyle JJ, Doyle JL. A rapid DNA isolation procedure from small quantities of fresh leaf tissues. *Phytochem Bull*. 1987;19:11–5.
24. Chin C-S, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, et al. Phased diploid genome assembly with single-molecule real-time sequencing. *Nat Methods*. 2016;13:1050–4.
25. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. 2018;34:3094–100.
26. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS One*. 2014;9:e112963.
27. Luo M, Wing RA. An improved method for plant BAC library construction. *Methods Mol Biol*. 2003;236:3–20.
28. Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: interactive visualization of de novo genome assemblies. *Bioinformatics*. 2015;31:3350–2.
29. Vaser R, Sović I, Nagarajan N, Šikić M. Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Res*. 2017;27:737–46.
30. Deschamps S, Zhang Y, Llaca V, Ye L, Sanyal A, King M, et al. A chromosome-scale assembly of the sorghum genome using nanopore sequencing and optical mapping. *Nat Commun*. 2018;9:4844.
31. Belsler C, Istace B, Denis E, Dubarry M, Baurens F-C, Falentin C, et al. Chromosome-scale assemblies of plant genomes using nanopore long reads and optical maps. *Nat Plants*. 2018;4:879–87.
32. Xu G-C, Xu T-J, Zhu R, Zhang Y, Li S-Q, Wang H-W, et al. LR_Gapcloser: a tiling path-based gap closer that uses long reads to complete genome assembly. *Gigascience*. 2019;8 Available from: <https://doi.org/10.1093/gigascience/giy157>.
33. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using repeat graphs. *Nat Biotechnol*. 2019;37:540–6.
34. Tang H, Zhang X, Miao C, Zhang J, Ming R, Schnable JC, et al. ALLMAPS: robust scaffold ordering based on multiple maps. *Genome Biol*. 2015;16:3.
35. Lu F, Romy MC, Glaubitz JC, Bradbury PJ, Elshire RJ, Wang T, et al. High-resolution genetic mapping of maize pan-genome sequence anchors. *Nat Commun*. 2015;6:6914.
36. Lee M, Sharopova N, Beavis WD, Grant D, Katt M, Blair D, et al. Expanding the genetic map of maize with the intermated B73 x Mo17 (IBM) population. *Plant Mol Biol*. 2002;48:453–61.
37. CyVerse Data Commons [Internet]. [cited 2019 Nov 12]. Available from: http://datacommons.cyverse.org/browse/iplant/home/shared/commons_repo/curated/Daniel_Laspisa_B73_RefGen_v4CEN_Feb_2019.
38. Portwood JL 2nd, Woodhouse MR, Cannon EK, Gardiner JM, Harper LC, Schaeffer ML, et al. MaizeGDB 2018: the maize multi-genome genetics and genomics database. *Nucleic Acids Res*. 2019;47:D1146–54.
39. Quinlan AR. BEDTools: the Swiss-Army tool for genome feature analysis. *Curr Protoc Bioinformatics*. 2014;47:11.12.1–34.
40. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol*. 1990;215:403–10.
41. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat Biotechnol*. 2011;644–52 Available from: <https://doi.org/10.1038/nbt.1883>.
42. Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*. 2015;33:290–5.
43. Liu R, Dickerson J. Strawberry: fast and accurate genome-guided transcript reconstruction and quantification from RNA-Seq. *PLoS Comput Biol*. 2017;13:e1005851.
44. Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, et al. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc*. 2012;562–78 Available from: <https://doi.org/10.1038/nprot.2012.016>.
45. Song L, Sabuncian S, Florea L. CLASS2: accurate and efficient splice variant annotation from RNA-seq reads. *Nucleic Acids Res*. 2016;44:e98.
46. Venturini L, Caim S, Kaithakottil GG, Mapleson DL, Swarbreck D. Leveraging multiple transcriptome assembly methods for improved gene structure annotation. *Gigascience*. 2018;7 Available from: <https://doi.org/10.1093/gigascience/giy093>.
47. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013;29:15–21.
48. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*. 2009;2078–9 Available from: <https://doi.org/10.1093/bioinformatics/btp352>.
49. Mapleson D, Venturini L, Kaithakottil G, Swarbreck D. Efficient and accurate detection of splice junctions from RNA-seq with Portcullis. *Gigascience*. 2018;7 Available from: <https://doi.org/10.1093/gigascience/giy131>.
50. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. De novo transcript sequence reconstruction from RNA-Seq: reference generation and analysis with Trinity. *Nat Protoc*. 2013;8 NIH Public Access; [cited 2019 Dec 10]. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3875132/>.
51. Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith RK Jr, Hannick LI, et al. Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. *Nucleic Acids Res*. 2003;31:5654–66.
52. Wang B, Regulski M, Tseng E, Olson A, Goodwin S, McCombie WR, et al. A comparative transcriptional landscape of maize and sorghum obtained by single-molecule sequencing. *Genome Res*. 2018;28:921–32.
53. Wu TD, Watanabe CK. GMAP: a genomic mapping and alignment program for mRNA and EST sequences. *Bioinformatics*. 2005;21:1859–75.
54. Kent WJ. BLAT—the BLAST-like alignment tool. *Genome Res*. 2002;12:656–64.

55. Campbell MS, Law M, Holt C, Stein JC, Moghe GD, Hufnagel DE, et al. MAKER-P: a tool kit for the rapid creation, management, and quality control of plant genome annotations. *Plant Physiol.* 2014;164:513–24.
56. Eilbeck K, Moore B, Holt C, Yandell M. Quantitative measures for the management and comparison of annotated genomes. *BMC Bioinformatics.* 2009;10:67.
57. Zhang R-G, Wang Z-X, Ou S, Li G-Y. TESorter: lineage-level classification of transposable elements using conserved protein domains. Available from: <https://doi.org/10.1101/800177>.
58. Vilella AJ, Severin J, Ureta-Vidal A, Heng L, Durbin R, Birney E. EnsemblCompara GeneTrees: complete, duplication-aware phylogenetic trees in vertebrates. *Genome Res.* 2009;19:327–35.
59. Manchanda N, Portwood JL, Woodhouse MR, Seetharam AS, Lawrence-Dill CJ, Andorf CM, et al. GenomeQC: a quality assessment tool for genome assemblies and gene structure annotations. *bioRxiv.* 2019:795237 [cited 2019 Dec 11]. Available from: <https://www.biorxiv.org/content/10.1101/795237v1.abstract>.
60. Ou S, Su W, Liao Y, Chougule K, Ware D, Peterson T, et al. Benchmarking transposable element annotation methods for creation of a streamlined, comprehensive pipeline. *bioRxiv.* 2019:657890 [cited 2019 Sep 24]. Available from: <https://www.biorxiv.org/content/10.1101/657890v1>.
61. Ou S, Chen J, Jiang N. Assessing genome assembly quality using the LTR Assembly Index (LAI). *Nucleic Acids Res.* 2018; Available from: <https://doi.org/10.1093/nar/gky730>.
62. Ou S, Jiang N. LTR_retriever: a highly accurate and sensitive program for identification of long terminal repeat retrotransposons. *Plant Physiol.* 2018;176:1410–22.
63. Gent JI, Wang N, Dawe RK. Stable centromere positioning in diverse sequence contexts of complex and satellite centromeres of maize and wild relatives. *Genome Biol.* 2017;18:121.
64. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv.* 2013; [q-bio.GN]. Available from: <http://arxiv.org/abs/1303.3997>.
65. Stovner EB, Sætrum P. epic2 efficiently finds diffuse domains in ChIP-seq data. *Bioinformatics.* 2019;35:4392–3.
66. Marçais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics.* 2011;27:764–70.
67. SanMiguel P, Vitte C. The LTR-retrotransposons of maize. In: Bennetzen JL, Hake S, editors. *Handbook of maize: genetics and genomics.* New York: Springer New York; 2009. p. 307–27.
68. Ramírez F, Dündar F, Diehl S, Grüning BA, Manke T. deepTools: a flexible platform for exploring deep-sequencing data. *Nucleic Acids Res.* 2014;42:W187–91.
69. Liu, J., Seetharam, A., Chougule, K., Ou, S., Swentowsky, K., Gent, J., Llaca, V., Woodhouse, M., Manchanda, N., Presting, G., Kudrna, D., Alabady, M., Hirsch, C., Fengler, K., Ware, D., Michael, T., Hufford, M., Dawe, K. Genome sequence and assembly of abnormal chromosome 10 genome Ab10. NCBI Sequence Read Archive. <https://www.ncbi.nlm.nih.gov/bioproject/PRJEB35367> (2020).
70. Liu, J., Seetharam, A., Chougule, K., Ou, S., Swentowsky, K., Gent, J., Llaca, V., Woodhouse, M., Manchanda, N., Presting, G., Kudrna, D., Alabady, M., Hirsch, C., Fengler, K., Ware, D., Michael, T., Hufford, M., Dawe, K. RNA-seq samples of ten tissues for B73 abnormal 10 (B73_Ab10). European Bioinformatics Institute. <https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-8641> (2020).
71. Liu, J., Seetharam, A., Chougule, K., Ou, S., Swentowsky, K., Gent, J., Llaca, V., Woodhouse, M., Manchanda, N., Presting, G., Kudrna, D., Alabady, M., Hirsch, C., Fengler, K., Ware, D., Michael, T., Hufford, M., Dawe, K. Ab10 genome assembly. Github. <https://github.com/dawelab/Ab10-Assembly> (2020).

Publisher's Note

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

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

